

PrivHAR: Recognizing Human Actions From Privacy-preserving Lens

Carlos Hinojosa^{1,2}, Miguel Marquez¹, Henry Arguello¹, Ehsan Adeli²,
Li Fei-Fei², and Juan Carlos Niebles²

¹ Universidad Industrial de Santander, Colombia

² Stanford University, USA

<https://carloshinojosa.me/project/privhar/>

Abstract. The accelerated use of digital cameras prompts an increasing concern about privacy and security, particularly in applications such as action recognition. In this paper, we propose an optimizing framework to provide robust visual privacy protection along the human action recognition pipeline. Our framework parameterizes the camera lens to successfully degrade the quality of the videos to inhibit privacy attributes and protect against adversarial attacks while maintaining relevant features for activity recognition. We validate our approach with extensive simulations and hardware experiments.

Keywords: Privacy-preserving lens design, human action recognition (HAR), adversarial training, deep optics.

1 Introduction

We are at the beginning of a new era of smart systems. From health care to video games, computer vision applications have provided successful solutions to real-world problems [3,24,28]. For decades, cameras have been engineered to imitate the human vision system, and machine learning algorithms are always constrained to be optimized using high-quality images as inputs. However, the abundance and growing uses of smart devices are also causing a social dilemma: we want intelligent systems (e.g., in our home) to recognize relevant events and assist us in our activities, but we also want to ensure they protect our privacy.

There have been some previous studies dealing with such a social dilemma. For instance, some early works rely on hand-crafted strategies, e.g., pixelation [52], blurring [33], face/object replacement [7], and person de-identification [1], to degrade sensitive information. More recently, Ren *et al.* [40] proposed an adversarial training strategy to learn to anonymize faces in videos and then perform activity detection. Similarly, using adversarial training, [56,57] proposed to optimize privacy attributes and recognition performance. However, all these methods rely on software-level processing of original high-resolution videos, which may already contain privacy-sensitive data. Hence, there is a possibility of these original videos being snatched by an attacker. Instead of developing new algorithms

carloshinojosa@saber.uis.edu.co

Fig. 1. Traditional HAR pipeline uses standard cameras that acquire visual details from the scene leading to privacy issues. We introduce PrivHAR, an adversarial optimization framework that learns a lens' phase mask to encode human action features and perform HAR while obscuring privacy-related attributes.

or designing software-level solutions that still rely on high-resolution images and videos as input, we believe that the privacy-preserving problem in computer vision should be addressed directly within the camera hardware, *i.e.*, sensible visual data should be protected before the images are acquired in the sensor.

Currently, few works have been developed in this direction. For instance, [41,42] proposed to use low-resolution cameras to create privacy-preserving anonymized videos and perform human action recognition. Also, Pittaluga *et al.* [37] introduced a defocusing lens to provide a certain level of privacy over a working region. On the other hand, several works used depth cameras to protect privacy and perform human action recognition [20,2]. These approaches rely on a fixed optical system; thus, their main contribution included designing an algorithm for a specific input type. More recently, [17] proposed to jointly design the lens of a camera and optimize a deep neural network to achieve two goals: privacy protection and human pose estimation. However, the formulation of the optimization in this work poses different problems: the privacy-preserving loss is not bounded as it maximizes an ℓ_2 term to enforce degradation, which may cause instability in the optimization; authors only used one human pose estimation model (OpenPose [6]) in all experiments and its not clear if the method works with other pose estimators. More importantly, to test and measure privacy, authors performed adversarial attacks after training the network; hence such attacks were not considered in the lens design.

In this paper, we address the problem of privacy-preserving human action recognition and propose a novel adversarial framework to provide robust privacy protection along the computer vision pipeline, see Fig. 1. We adopt the idea of end-to-end optimization of the camera lens and vision task [17,31,44] and propose an optimization scheme that: (1) Incorporates adversarial defense objectives into the learning process across a diversity of canonical privacy categories, including face, skin color, gender, relationship, and nudity detection. (2) Encourages distortions in the videos without compromising the training stability by including the structural similarity index (SSIM) [18] in our optimization loss.

(3) To further preserve the temporal information in the distorted videos, we use temporal similarity matrices (TSM) and constrain the structure of the temporal embeddings from the private videos to match the TSM of the original video.

We test our approach with two popular human action recognition backbone networks. To experimentally test our privacy-preserving human action recognition network (PrivHAR) and lens design, we built a proof-of-concept optical system in our lab. Our testbed acquires distorted videos and their non-distorted version simultaneously. Our experimental results in hardware match the simulations. While we do observe a trade-off between Human Action Recognition (HAR) accuracy and image distortion level, our proposed PrivHAR system offers robust protection with reasonable accuracy.

2 Related Work

Human action recognition is a challenging task [35] and has many applications, such as video surveillance, human-computer interfaces, virtual reality, video games, and sports action analysis. Therefore, developing privacy-preserving approaches for HAR is even more challenging and has not been widely explored.

Human Action Recognition (HAR). Nowadays, there are multiple approaches in the computer vision literature for addressing the HAR problem. Some prior work relies on 2D CNNs to conduct video recognition [8,22,53,9]. A major drawback of 2D CNN approaches is not properly modeling the temporal dynamics. On the other hand, 3D CNN-based approaches use spatial and temporal convolutions over the 3D space to infer complicated spatio-temporal relationships. For instance, C3D [49] is a 3D CNN based on the VGG model that learns spatio-temporal features from a frame sequence. However, 3D CNNs are typically computationally heavy, making the deployment difficult. Therefore, many efforts on HAR focus on proposing new efficient architectures; for example, by decomposing 3D filters into separate 2D spatial and 1D temporal filters [23,50] or extending efficient 2D architectures to 3D counterparts. Moreover, RubiksNet [13] is a hardware-efficient architecture for HAR based on a shift layer that learns to perform shift operations jointly in spatial and temporal context. We build our proposed PrivHAR using both C3D and RubiksNet.

Privacy protection in Computer Vision. Currently, few works address the privacy-preserving HAR problem. We divide prior work into software-level and hardware-level protection, where we consider the latter more robust to attacks.

Software-level Privacy-preserving HAR. Most prior privacy-preserving works apply different computer vision algorithms to the video data after their acquisition. The literature has relied on domain knowledge and hand-crafted approaches, such as pixelation, blurring, and face/object replacement, to protect sensitive information [1,7,33]. These methods can be useful in settings when we know in advance what to protect in the scene. More recent works propose a more general approach that learns privacy-preserving encodings through adversarial training [56,5,36]. These methods learn to degrade or inhibit privacy attributes while maintaining important features to perform inference tasks and provide more robust videos to adversarial attacks. Ren *et al.* [40] use adversarial training to learn a video anonymizer and remove facial features for activity detec-

Fig. 2. Our proposed end-to-end framework. Our **optical component** consists of a camera with two thin convex lenses and a phase mask between them. We achieve robust privacy protection by training an adversarial framework under three goals: (1) to learn to add aberrations to the lens surface such that the acquired videos are distorted to obscure private attributes while still preserving features to (2) achieve high **video action recognition** accuracy, and (3) being robust to **adversarial attacks**.

tion. Similarly, Wu *et al.* [56,57] proposed an adversarial framework that learns a degradation transform for the video inputs using a 2D convolution layer. These works optimize the trade-off between action recognition performance and the associated privacy budget on the degraded video. Although these software-level approaches preserve privacy, the acquired images are not protected.

Hardware-level Privacy-preserving HAR. These approaches rely on the camera hardware to add an extra layer of security by removing sensitive data during the imaging sensing. Prior hardware-level privacy-preserving approaches use low-resolution cameras to anonymize videos, i.e., the videos are intentionally captured to be in special low-quality conditions that only allow for the recognition of some events or activities while avoiding the unwanted leak of the identity information for the human subjects in the video [39,42]. In [38,37], two optical designs were proposed that filter or block sensitive information directly from the incident light-field before sensor measurements acquisition, enabling k-anonymity and privacy protection by using a camera with defocusing lens. In particular, they show how to select a defocus blur that provides a certain level of privacy over a working region within the sensor size limits; however, only using optical defocus for privacy may be susceptible to reverse engineering. In addition, the authors did not test their method on the action recognition task. More recently, [55] proposed a coded aperture camera system to perform privacy-preserving HAR directly from encoded measurements without the need for image restoration. However, it was only tested for indoor settings in a small dataset.

3 Privacy-preserving Action Recognition

We are interested in human action recognition from privacy-preserving videos. We propose a framework to accomplish three goals: 1) to learn the parameters of a robust privacy-preserving lens by backpropagating the gradients from the action recognition and adversarial branches to the camera; 2) to learn the pa-

parameters of an action recognition network to perform HAR on the private videos with high accuracy; 3) to obtain private videos that are robust to adversarial attacks. Our framework (Fig. 2) consists of three parts: optical, action recognition, and an adversarial component.

The optical component consists of a camera with two thin convex lenses and a phase mask between them. Our simulated camera takes a video $\mathbf{V}_x \in \mathbb{R}_+^{w \times h \times 3 \times T} = f\mathbf{X}_t g_{t=1}^T$ as input, which has $w \times h$ pixels and T frames, and outputs the corresponding distorted video $\mathbf{V}_y \in \mathbb{R}_+^{w \times h \times 3 \times T}$. Formally, $\mathbf{V}_y = O(\mathbf{V}_x)$, where we denote our designed camera as the function $O(\cdot)$, which distorts every single frame $\mathbf{X}_t \in \mathbb{R}_+^{w \times h \times 3}$, and produces the respective *private* frames $\mathbf{Y}_t \in \mathbb{R}_+^{w \times h \times 3}$. Then, the distorted video \mathbf{V}_y passes through the action recognition component where a convolutional neural network C predicts the class labels. Besides, \mathbf{V}_y also passes through the adversarial component where an attribute estimator network A tries to predict the private information (attributes) from the distorted video. All three components consist of neural networks with trainable parameters, and the whole framework is trained adversarially. At the end of the optimization process, we obtain the optimal camera lens parameters θ_o , and the optimal action recognition parameters θ_c . Hence, the loss function of our adversarial framework is formulated as follows:

$$\theta_o, \theta_c = \arg \min_{\theta_o, \theta_c} L(O) + L(C) - L(A); \quad (1)$$

where $L(O)$, $L(C)$, and $L(A)$ are the loss functions for our optical component, action recognition component, and adversarial component, respectively.

During inference, we can construct a camera lens using the optimal parameters θ_o that acquires degraded images, on which our network C can perform HAR. Since we develop our protection directly in the optics (camera lens), it provides an extra layer of protection and, hence, is more difficult for a hacker to attack our system to reveal the person's identity. One could also deploy a less secure software-only approach implementing image degradation post-acquisition. A hybrid solution consists of designing an embedded chipset responsible for distorting the videos immediately after the camera sensor.

3.1 Optical Component

The main goal of the optical component in our adversarial framework PrivHAR (Fig. 2) is to design a phase mask to visually distort videos (hence obscuring privacy-sensitive attributes), encode the physical characteristics and preserve human action features to perform HAR. We adopted a similar strategy as the authors in [44,17] to couple the modeling and design of two essential operators in the imaging system: wave propagation and phase modulation.

Image Formation Model. We model the image acquisition process using the point spread function (PSF) defined in terms of the lens surface profile to emulate the wavefront propagation and train the parameters of the refractive lens. Considering by the Fresnel approximation and the paraxial regime [14], for incoherent illumination, the PSF can be described by

$$H(u^0; v^0) = jF^{-1} f F f t_L(u; v) t(u; v) W(u; v) g^{-T}(f_u; f_v) g f^2; \quad (2)$$

where $W(u; v)$ is the incoming wavefront, $T(\cdot)$ represents the transfer function with $(f_u; f_v)$ as the spatial frequencies, $t(u; v) = \exp(-ik \phi(u; v))$ with $\phi(u; v)$

as the lens phase mask and $k = 2\pi/\lambda$ as the wavenumber, $t_L(\cdot)$ denotes the light wave propagation phase with $t_L(u; v) = \exp\left[-i\frac{k}{2z}(u^2 + v^2)\right]$ with z as the object-lens distance, $\mathcal{F}\mathcal{F}g$ denotes the 2D Fourier transform, and $(u^\rho; v^\rho)$ is the spatial coordinate on the camera plane. The values of $\phi(\cdot)$ are modelled via the Zernike polynomials with $\phi(u; v) = R_n^m(\rho)\sqrt{2} \cos(\arctan(v/u))$, where $R(\cdot)$ represents the radial polynomial function [27], m and n are nonnegative integers with $n \geq m \geq 0$. To train the phase mask values using our PrivHAR, we discretize the phase mask $\phi(\cdot)$ as:

$$\phi = \sum_{j=1}^J \alpha_j \mathbf{Z}_j, \quad (3)$$

where \mathbf{Z}_j denotes the j -th Zernike polynomial in Noll notation, and α_j is the corresponding coefficient [4]. Each Zernike polynomial describes a wavefront aberration [27]; hence the phase mask ϕ is formed by the linear combination of all aberrations. In this regard, the optical element parameterized by ϕ can be seen as an optical encoder, where the coefficients α_j determine the data transformation. Therefore, our adversarial training finds a set of coefficients $\alpha_o = \mathcal{F} \sum_{j=1}^J \alpha_j \mathbf{Z}_j$ that provides the maximum visual distortion of the scene but allows to extract relevant features to perform HAR. Using the defined PSF-based propagation model (assuming that image formation is a shift-invariant convolution of the image and PSF), the acquired private images for each RGB channel can be modelled as:

$$\mathbf{Y}_\cdot = G_\cdot(\mathbf{H}_\cdot \cdot \mathbf{X}_\cdot) + \mathbf{N}_\cdot; \quad (4)$$

where $\mathbf{X}_\cdot \in \mathbb{R}_+^{w \times h}$ represents the discrete image from the \cdot channel, with each pixel value in $[0; 1]$; \mathbf{H}_\cdot denotes the discretized version of the PSF [14] in Eq. (2) for the channel \cdot , $\mathbf{N}_\cdot \in \mathbb{R}^{w \times h}$ represents the Gaussian noise in the sensor, and $G_\cdot(\cdot) : \mathbb{R}^{w \times h} \rightarrow \mathbb{R}^{w \times h}$ is the camera response function, which is modeled as a linear function. Please see our supplementary document for a schematic diagram of the light propagation in our model.

Loss Function. To encourage image degradation, we train our network to minimize the quality of the acquired image by our camera $\mathbf{Y} = \mathcal{F} \mathbf{Y}_\cdot \mathcal{G}_{\cdot=1}^3$ in comparison with the original image $\mathbf{X} = \mathcal{F} \mathbf{X}_\cdot \mathcal{G}_{\cdot=1}^3$. Instead of maximizing the ℓ_2 norm error between the two images as previous works did [17], we use the structural similarity index (SSIM) [54] in our optimization loss to measure quality. The ℓ_2 norm does not have an upper bound; hence maximizing it to enforce degradation causes instability in the optimization. On the other hand, the SSIM function is bounded, which leads to better stability during training. Specifically, the SSIM value ranges between 0 and 1, where values near 1 (better quality) indicate more perceptual similarity between the two compared images. Then, we define the loss function for our camera lens optimization as:

$$L(O) = \text{SSIM}(\mathbf{X}; \mathbf{Y}); \quad (5)$$

Since we encourage distortion in the camera's output images/videos, the $L(O)$ loss is minimized in our adversarial training algorithm, see Algorithm 1.

3.2 Action Recognition Component

We can use any neural network architecture in our adversarial framework to perform human action recognition. In this work, without loss of generality, we adopt two HAR CNN architectures: the well-known C3D [49], and the Rubkismet

[13], a more recent and efficient architecture for HAR. For a set of private videos, we assume that the output of the classifier C is a set of action class labels S_C . Then, we can use the standard cross-entropy function H as the classifier's loss.

On the other hand, since our degradation model distorts each frame of the input video separately (2D convolution), part of the temporal information could be lost, decreasing the performance of the HAR CNN significantly. To preserve temporal information, we use temporal similarity matrices (TSMs). TSMs are useful representations for human action recognition and have been employed in several works [21,12,34,45] due to their robustness against dynamic view changes of the camera when paired with appropriate feature representation. Unlike previous works, we propose using TSMs as a proxy to keep the temporal information (features) similar after distortion: we build a TSM for the original and private videos and compare their structures. Specifically, we take the embeddings $\hat{\mathbf{e}}$ from the last convolutional layer of our HAR CNN architecture and compute the TSM values using the negative of the squared euclidean distance, i.e., $(\mathbf{T}_m^0)_{n_1 n_2} = -k\hat{\mathbf{e}}_{n_1} - \hat{\mathbf{e}}_{n_2}k^2$. Then, we calculate the mean square error (MSE) between the \mathbf{T}_m^0 and the TSM from the input video \mathbf{T}_m , which was computed similarly using the last convolutional layer of the corresponding pretrained HAR CNN (non-privacy) network. We define the action recognition objective as:

$$L(C) = H(S_C; C(V_y)) + MSE(\mathbf{T}_m; \mathbf{T}_m^0); \quad (6)$$

where V_y denotes the set of E private videos: $V_y = f\mathbf{V}_y^e g_{e=1}^E = fO(\mathbf{V}_x^e)g_{e=1}^E$.

3.3 Adversarial Component and Training Algorithm

The attacks that an adversarial agent could perform to our privacy-preserving pipeline depends on the definition of privacy. There are different ways to measure privacy and this is, in general, not a straightforward task. For example, in smart homes with video surveillance, one might often want to avoid disclosure of the face or identity of persons. Therefore, an adversarial agent could try to attack our system by training a face detection network. However, there are other privacy-related attributes, such as race, gender, or age, that an adversarial agent could also want to attack too. In this work, we define the adversarial attack as a classification problem, where a CNN network A takes a private video \mathbf{V}_y as input and tries to predict the corresponding private information. Therefore, the goal of our adversarial training is to try that the predictions from A diverges from the set of class labels S_A that describe the private information within the scene. To train the attribute estimator network, we also use the cross-entropy H function and define the adversarial loss as:

$$L(A) = H(S_A; A(V_y)); \quad (7)$$

Algorithm 1 summarizes the proposed adversarial training scheme. Before performing the adversarial training, we first train each framework component separately without privacy concern to obtain the optimal performance on each task. Specifically, we train the optical component O by minimizing $1 - L(O)$ to acquire videos without distortions, i.e., \mathbf{V}_y videos are very similar to the corresponding input \mathbf{V}_x . We also train the HAR network C by minimizing $H(S_C; C(V_x))$, obtaining the highest action recognition accuracy (the upper bound). Finally, we train the attribute estimator network A by minimizing $H(S_A; A(V_x))$, thus obtaining the highest classification accuracy (the upper

Algorithm 1: Our Adversarial Training Algorithm.

Input : Video Dataset $V_x = \{V_x^e\}_{e=1}^E$. Hyperparameters $\alpha, \beta, \gamma, \lambda, \eta, \epsilon$
Output: $\theta_o, \theta_c, \theta_a$
Function $\text{Train}(V_x; \alpha, \beta, \gamma, \lambda, \eta, \epsilon)$

```

1  for every epoch do
2      for every batch of videos  $V_x^B$  do
3           $V_y^B = O(V_x^B)$  . Acquire private videos
4           $\theta_o \leftarrow \theta_o - \eta \nabla_{\theta_o} (L(O) + \alpha L(C) + \beta L(A))$ 
5           $\theta_c \leftarrow \theta_c - \eta \nabla_{\theta_c} L(C)$ 
6           $\theta_a \leftarrow \theta_a - \eta \nabla_{\theta_a} L(A)$ 
7  return  $\theta_o, \theta_c, \theta_a$ 
```

bound). After initialization, we start the adversarial training shown in Algorithm 1, where, for each epoch and every batch, we first acquire the private videos with our camera O . Then, we update the parameters of the camera θ_o by freezing the attribute estimator network parameters θ_a and minimizing the weighted sum $L(O) + \alpha L(C) + \beta L(A)$, shown on line 4 of the algorithm. Similarly, we update the parameters of the HAR network θ_c by freezing the attribute estimator network parameters and using the private videos acquired on line 3 to minimize $L(C)$. Finally, we perform the adversarial attack by minimizing $L(A)$ and updating the parameters of the attribute estimator network θ_a while the camera and HAR network parameters are fixed. Contrary to the prior work [17], our training scheme jointly models the privacy-preserving optics with HAR and adversarial attacks during training.

4 Experimental Results

Datasets. Given the lack of a public dataset containing both human actions and privacy attribute labels on the same videos, we follow the same approach as authors in [56] to train our proposed adversarial framework. Specifically, we perform cross-dataset training using three datasets: the HMDB51 [25], the VISPR [32], and the PA-HMDB51[56]. The VISPR dataset contains 22,167 images annotated with 68 privacy attributes which include: semi-nudity, face, race, gender, skin color, among others. The attributes of a specific image are labeled as "present" or "not-present". The HMDB51 dataset comprises 6,849 video clips from 51 action categories, with each category containing at least 101 clips. The Privacy-annotated HMDB51 (PA-HMDB51) is a small subset of the HMDB51 dataset, containing 515 videos, with privacy attribute labels. For each video in PA-HMDB51, there are five attributes annotated on a per-frame basis: skin color, face, gender, nudity, and relationship. Similar to VISPR, the labels are binary and specify if an attribute is present or not in the frame.

Training set. We train our models using cross-dataset training on HMDB51 and VISPR datasets. Specifically, we exclude the 515 videos in the PA-HMDB51 dataset from HMDB51 and use the remainder videos to train our action recognition component. On the other hand, we use the VISPR dataset with the same five privacy attributes available in the PA-HMDB51 dataset: skin color, face, gender, nudity, and relationship, to train our adversarial component.

Testing set. We use PA-HMDB51 to test our action recognition and adversarial components. This dataset includes both action and privacy attribute labels.

Training details. In Algorithm 1, we set initial learning rates $\eta_o = 3 \times 10^{-3}$; $\eta_c = \eta_a = 10^{-4}$, and $\alpha_1 = 0.7$; $\alpha_2 = 0.3$ and applied an exponential learning decay with a decay factor of 0.1 that is triggered in the epoch 25. We trained the end-to-end PrivHAR model for 50 epochs, with batch size of 8, and use the Stochastic Gradient Descent (SGD) optimizer to update parameters θ_o ; θ_c ; θ_a . To perform the adversarial attacks during training (adversarial component in Fig. 2), we use the ResNet-50 architecture. Training the PrivHAR for 50 epochs took about 6 hours on 8 Nvidia TITAN RTX GPU with 24 GB of memory.

4.1 Metrics and Evaluation Method

To measure the overall performance of PrivHAR, we evaluate the action recognition task and privacy protection separately. First, to test action recognition, we pass the testing videos through our designed camera lens $O(\cdot)$ to obtain the private videos. Next, we use our learned HAR backbone $C(\cdot)$ to get the predicted actions on each private video. Similarly as C3D [49], and RubiksNet [13], we report the standard average classification accuracy, denoted by A_C .

On the other hand, to evaluate privacy protection, we follow the same evaluation protocol adopted by authors in [56]. Specifically, assuming that an attacker has access to the set of private videos acquired with our $O(\cdot)$ and the corresponding privacy attribute labels, then, the attacker can train different CNNs to try to steal sensitive information from the privacy-protected videos acquired with our camera. To empirically verify that our protection is robust to this kind of attack, we separately train ten different classification networks using the private images acquired with our camera, i.e., these CNNs are different from the selected CNN used during training. To train these networks, we use the same training set defined in the previous section and fix our camera component with the optimal learned parameters θ_o . We use the following architectures: ResNet- $r50$;101g[16], Wide-ResNet- $r50$;101g[58], MobileNet-V2 [43], Inception- $v1$;V3g[47,46], MNASNet- $r0.5$;0.75;1.0g[48]. Among these CNNs, eight randomly selected networks start from ImageNet-pretrained weights. The remaining two models were trained from scratch (random initialization) to eliminate the possibility that the initialization with ImageNet weights affects the correct predictions. After training, we evaluate each model on our defined testing set (videos from PA-HMDB51) and select the model with the highest performance. Similar to previous works [10,56,57,32], we adopt the Class-based Mean Average Precision (C-MAP)[32] to assess the performance of the models. Specifically, we compute the Average Precision (AP) per class, which is the area under the Precision-Recall curve of the privacy-related attribute. Hence, C-MAP corresponds to the average of the AP scores across all the privacy-related attributes. We also denote C-MAP as A_A in our experiments (lower is better).

To measure image degradation, we use the structural similarity index (SSIM) metric [18]. Large values of SSIM indicate high quality. Thus, in general, we expect to achieve the minimum SSIM values while achieving high A_C and low A_A for HAR and Adversarial accuracy, respectively. Besides, we combine the

| (a) Ablation Study | | | | | (b) Comparisons. | | | | |
|--------------------|-------|-------------|-------------|--------------|--------------------------|-------|-------------|-------------|-------------|
| C3D Backbone | | | | | Methods | | | | |
| Experiment | SSIM# | A_C % | A_A % | P % | | SSIM# | A_C % | A_A % | P % |
| No-Adversarial | 0.603 | 51.1 | 69.1 | 38.6 | No-privacy (C3D) | 1.0 | 71.1 | 76.1 | 35.8 |
| No-TSM | 0.612 | 59.9 | 69.7 | 40.2 | No-privacy (RubiksNet) | 1.0 | 85.2 | 76.1 | 37.3 |
| Zernike-50 | 0.643 | 58.3 | 70.5 | 39.2 | Low-resolution [41] | 0.686 | 48.3 | 70.9 | 36.3 |
| Zernike-100 | 0.629 | 58.8 | 69.3 | 40.4 | Lens in [17]-RubiksNet | 0.608 | 52.4 | 69.4 | 38.6 |
| Zernike-200 | 0.612 | 63.3 | 68.9 | 41.52 | Defocus [37] | 0.688 | 62.1 | 72.5 | 38.1 |
| RubiksNet Backbone | | | | | PDAR-GRL [56] | - | 63.3 | 70.5 | 40.2 |
| No-Adversarial | 0.592 | 57.6 | 68.2 | 40.9 | PDAR-K-Beam [56] | - | 63.5 | 69.3 | 41.4 |
| No-TSM | 0.599 | 72.3 | 67.6 | 44.6 | PDAR-Entropy [56] | - | <u>67.3</u> | 70.3 | 41.2 |
| Zernike-50 | 0.618 | 70.2 | 69.2 | 42.8 | PrivHAR-C3D | 0.612 | 63.3 | 68.9 | 41.7 |
| Zernike-100 | 0.601 | 71.9 | 68.4 | 43.9 | PrivHAR-RubiksNet | 0.588 | 73.8 | 66.5 | 46.1 |
| Zernike-200 | 0.588 | 73.8 | 66.5 | 46.1 | | | | | |

Table 1. Quantitative Results. (a) Multiple ablation studies of our method for two different HAR backbones, C3D and RubiksNet: each component in Fig. 2 is trained separately (No-Adversarial); not using the TSM matrices to preserve temporal information (No-TSM); 50, 100, and 200 Zernike polynomials to design our lens. (b) Comparison of our method (**PrivHAR**) with: three additional privacy-preserving approaches: defocusing, low-resolution cameras, and the lens used in [17]; and the privacy-preserving deep action recognition (PDAR) framework with different learning approaches (GRL, K-Beam, and Entropy) [56]. Accuracy values are reported in percentage.

two accuracy metrics (A_C and A_A) into one using the harmonic mean as:

$$P = \frac{2}{\frac{1}{A_C} + \frac{1}{A_A}} = \frac{2A_C(1 - A_A)}{1 - A_A + A_C}, \quad (8)$$

and we expect to achieve the maximum P value.

4.2 Simulation Experiments

Ablation Studies. We conduct multiple experiments to investigate different configurations for our adversarial approach. We show the quantitative results of our ablations studies in Table 1 (a), for C3D and RubiksNet. We first train the optical and action recognition components to obtain privacy-preserving videos and perform HAR on them. Then, we fix the optical component and train the adversarial CNN to recover the privacy attributes from the videos. We refer to this experiment as 'No-adversarial' in the Table 1 (a). Note that this approach is similar to the prior work [17] but on a different vision task. In our second experiment (No-TSM), we test the performance of our proposed PrivHAR with $q = 200$ Zernike coefficients when not using TSMs to preserve the temporal information. We can observe from the table that, in general, the A_C decreases, which evidences the importance of using TSMs to preserve temporal information. The third experiment consists of training our adversarial framework with a different number of Zernike coefficients. Specifically, we trained our PrivHAR using $q = 50$, $q = 100$, and $q = 200$ Zernike coefficients, see Eq. (3). In general, increasing the number of Zernike coefficients leads to better encoding; hence the A_C value increases while the SSIM decreases. However, memory consumption

Fig. 3. (a) Trade-off between privacy protection and action recognition on PA-HMDB51. Vertical and horizontal, dashed and dotted, purple lines indicate A_A and A_C on the original non-privacy videos, using RubiksNet and C3D backbones for HAR, respectively. The red dashed line indicates where $A_A = A_C$. (b) Face recognition performance on private images (from LFW [19] dataset) acquired with our optimized lens.

also increases since we need to store all the Zernike bases. In general, we use $q = 200$ Zernike coefficients as a default value in all other experiments. The tables show that the best HAR backbone for our proposed PrivHAR network is RubiksNet. We observed that when using RubiksNet, PrivHAR achieves higher distortions (lower SSIM) affecting the performance of the adversarial component while achieving high action recognition accuracy. We empirically verify that RubiksNet is better at preserving the temporal information than C3D; hence it performs better even with high image distortions. Besides, we observed that TSM helps more the C3D backbone, which is more affected by the distortions generated by our lens.

Attribute Estimator Network Performance. The values of A_A reported in the tables corresponds to the C-MAP obtained by the model with highest performance on our testing set, as described in Section 4.1. To analyze the performance of the attribute estimator networks, and hence our privacy protection, we plot the receiver operating characteristic (ROC) and Precision-Recall (PR) curves. In our supplementary document, we show the ROC and PR curves of the attribute estimator network which achieves the best performance on the privacy-preserving images/videos acquired with our camera. Specifically, considering the area under curve (AUC) of the PR curves, we obtain an average precision (AP) of 0.94, 0.72, 0.97, 0.52, 0.18 for skin color, face, gender, nudity, and relationship, respectively. These values of AP are very close to those obtained by a random classifier (null hypothesis), which are 0.95, 0.71, 0.97, 0.58, 0.17. Therefore, based on the Fisher's exact test [51], the best attribute estimator network on our privacy-protected images is not significantly different from the random classifier (p -value < 0.01).

Comparison with other methods. We compare our proposed PrivHAR with two traditional privacy-preserving approaches: low-resolution [42] and defocusing cameras [37]. We simulate both types of cameras and perform a similar training as shown in Fig. 2. To implement the low-resolution approach, we manually downsampled the images with a resolution of 16×16 . In addition, we compare our

proposed PrivHAR with the privacy-preserving deep action recognition (PDAR) framework with different learning approaches (GRL, K-Beam, and Entropy) [56]. We present the quantitative results in Table 1 (b), where all methods use the C3D backbone for HAR if not otherwise specified. We also include our PrivHAR using RubiksNet for comparison. Furthermore, we use the lens designed in [17], which was optimized for human pose estimation and did not consider adversarial attacks during training, for distort the videos and then perform HAR on them. This approach obtains an $A_C = 52.4\%$ using RubiksNet, which is 21.4% lower than our PrivHAR-RubiksNet results. In addition, the trade-off between privacy protection and action recognition is visualized in Fig. 3 (a), which shows PrivHAR obtains the best privacy while maintaining high accuracy.

Face recognition performance. We follow the same face recognition validation on private images acquired by the optimized lens as the prior work in [17]. Specifically, we use an implementation of the face recognition network ArcFace [11], train on Microsoft Celeb (MS-Celeb-1M) [15] and test on LFW [19] datasets. Figure 3 (b) show the ROC curves for each testing approach: "No-privacy Model" uses the pretrained ArcFace model on the original (non-private) images; "Pre-trained model" uses the pretrained ArcFace model on the private version of each dataset; "Trained model" uses an ArcFace model trained from scratch using the private version of the MS-Celeb-1M dataset; "Fine-tuned Model" uses a pre-trained ArcFace model fine-tuned with the private version of the MS-Celeb-1M dataset. From the figure, we can conclude that the ArcFace model does not perform well on the images generated by our designed lens as the best performance is achieved by the fine-tuned model ($AUC = 0.68$), which is still close to random classifier's performance. See results with others datasets in our supplementary.

Qualitative Results. We qualitatively compare our approach with low resolution and defocusing cameras in Fig. 4. We show results on three example videos from the PA-HMDB51 dataset. The first row of the figure shows the non-privacy video acquired using a standard lens and the ground truth (GT) of the actions for reference. As observed, our lens achieves a higher distortion but still performs action recognition. The last video shows a failure case of our method.

Deconvolution Attack. Suppose the attacker has access to the camera or a large collection of acquired images with our proposed camera. In that case, the attacker could use deconvolution methods (blind and non-blind) on our distorted images to recover people's identities. To test the robustness of our designed lens to deconvolution attacks, we assume both scenarios: having access to the camera, we can easily get the PSF (by imaging a point of source light) and hence use a non-blind deconvolution method, e.g. the Wiener deconvolution; on the other hand, not having access to the camera but a large collection of our distorted images then we can train a blind deconvolution network, e.g. DeblurGAN [26]. We describe the training details in our supplementary document. In Fig. 5 we show the results with two video frames from the HMDB51 dataset with people near the camera. We observed that the distortion achieved by PrivHAR-RubiksNet (RBN) is significantly higher than C3D; hence it is more difficult for DeblurGAN and Wiener deconvolution to recover the scene. In both cases, using C3D

Fig. 4. Qualitative Results on PA-HMDB51. Each row shows standard no-privacy videos and ground truth (GT) labels (top); and predictions from our optimized lens (PrivHAR-RubiksNet, bottom) to low-resolution (second) and defocus (third) cameras.

Fig. 5. Deconvolution of private images acquired with our optimized lens using C3D and RubiksNet (RBN) backbones in PrivHAR. The images acquired with our lens are robust to deconvolution, and DeblurGAN cannot recover people's identities.

or RBN, the distortion is sufficient to avoid recovering face details, and the people's identity is protected. However, some attributes are visible in the recovered scene when using PrivHAR-C3D. It is possible to obtain a lens with C3D that provides more distortion; however, the HAR accuracy could be affected.

4.3 Hardware Experiments

To demonstrate the PrivHAR's capability of action recognition, we conduct experimental validations acquiring four human actions: jump, clap, punch, and hair brush in our Lab. We emulate the lens designed with our PrivHAR adversarial framework using a deformable mirror-based 4f system [29,30]. We first train our system using $q = 15$ Zernike coefficients and then load the learned coefficients to the deformable mirror and calibrate the PSF. After calibration we obtained the following learned Zernike coefficients: $f_1 = f_2 = f_3 = 0$; $f_4 = 0.45$; $f_5 = 0.36$; $f_6 = 0.24$; $f_7 = 0.6$; $f_8 = 0.4$; $f_9 = 0.11$; $f_{10} = 0.69$; $f_{11} = 0.31$; $f_{12} =$

Fig. 6. Experimental setup scheme and some results on acquired videos. The deformable mirror configuration and characterized PSF are shown in the upper left corner. The right column shows the non-privacy and private videos obtained with our camera.

0.15; $\gamma_{13} = 0.70$; $\gamma_{14} = 0.85$; $\gamma_{15} = 0.38g$. The resulting PSF and the used phase mask are presented in Fig. 6(Left). Finally, we placed our proof-of-concept system on a movable table to take it out of our Lab and acquire real outdoor images. In Fig. 6(Right), we show the human action recognition for two video sequences recorded by our 4F-based system. The ground truth and the private version were illustrated in the first and second rows, respectively. Outdoor system configuration, additional qualitative and quantitative results, and detailed description of the proof-of-concept system can be found in the supplement.

5 Discussion and Conclusion

We present PrivHAR, a framework for detecting human actions from a privacy-preserving lens. Our framework consists of three components: the hardware component that comprises a camera with a privacy-preserving lens, whose parameters are learned during training and its main function is to obscure sensitive private information; the action recognition component that aims to preserve temporal information using temporal similarity matrices and performs HAR on the degraded video; and the adversarial component, which performs reverse attacks to the private videos seeking to recover the hidden attributes.

Limitations. One limitation of our simulated experiments is that we test our approach on a relatively small set due to the lack of a public dataset containing human actions and privacy attribute labels on the same videos. As future work, we plan to build a video dataset using our proposed optical system, which allows us to acquire both RGB and private videos. In addition, the deformable mirror is the main limitation of the proof-of-concept optical system. This device can only use $q = 15$ Zernike Polynomials, limiting the scene's level of distortion. For now, our small-scale tests show results consistent with our extensive experiments.

Conclusion. We extensively evaluated and experimentally validated our approach in simulations and a hardware prototype. Our qualitative and quantitative results indicate a trade-off between image degradation and HAR accuracy. Our optics modeling can generally be integrated into an embedded chipset or used as a software-only solution by applying the image degradation post-acquisition to deploy a less secure system. However, we show that the learned lens can be deployed as a camera, which provides a higher security layer. One could connect it to an Nvidia Jetson for real-time privacy-preserving HAR.

References

1. Agrawal, P., Narayanan, P.: Person de-identification in videos. *IEEE Transactions on Circuits and Systems for Video Technology* **21**(3), 299{310 (2011) [1](#), [3](#)
2. Ahmad, Z., Illanko, K., Khan, N., Androutsos, D.: Human action recognition using convolutional neural network and depth sensor data. In: *Proceedings of the 2019 International Conference on Information Technology and Computer Communications*. pp. 1{5 (2019) [2](#)
3. Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M.S., Bohg, J., Bosselut, A., Brunskill, E., et al.: On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021) [1](#)
4. Born, M., Wolf, E.: *Principles of optics: electromagnetic theory of propagation, interference and diffraction of light*. Elsevier (2013) [6](#)
5. Brkic, K., Sikiric, I., Hrkac, T., Kalafatic, Z.: I know that person: Generative full body and face de-identification of people in images. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. pp. 1319{1328. IEEE (2017) [3](#)
6. Cao, Z., Hidalgo, G., Simon, T., Wei, S.E., Sheikh, Y.: Openpose: realtime multi-person 2d pose estimation using part affinity fields. *IEEE TPAMI* **43**(1), 172{186 (2019) [2](#)
7. Chen, D., Chang, Y., Yan, R., Yang, J.: Tools for protecting the privacy of specific individuals in video. *EURASIP Journal on Advances in Signal Processing* **2007**, 1{9 (2007) [1](#), [3](#)
8. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1251{1258 (2017) [3](#)
9. Christoph, R., Pinz, F.A.: Spatiotemporal residual networks for video action recognition. *Advances in neural information processing systems* pp. 3468{3476 (2016) [3](#)
10. Dave, I.R., Chen, C., Shah, M.: Spact: Self-supervised privacy preservation for action recognition. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 20164{20173 (2022) [9](#)
11. Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 4690{4699 (2019) [12](#)
12. Dwibedi, D., Aytar, Y., Tompson, J., Sermanet, P., Zisserman, A.: Counting out time: Class agnostic video repetition counting in the wild. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 10387{10396 (2020) [7](#)
13. Fan*, L., Buch*, S., Wang, G., Cao, R., Zhu, Y., Niebles, J.C., Fei-Fei, L.: RubiksNet: Learnable 3D-Shift for Efficient Video Action Recognition. In: *European Conference on Computer Vision*. pp. 505{521. Springer (2020) [3](#), [7](#), [9](#)
14. Goodman, J.W.: *Introduction to Fourier optics*. Macmillan Learning, 4 edition (2017) [5](#), [6](#)
15. Guo, Y., Zhang, L., Hu, Y., He, X., Gao, J.: Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In: *European conference on computer vision*. pp. 87{102. Springer (2016) [12](#)
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770{778 (2016) [9](#)

17. Hinojosa, C., Niebles, J.C., Arguello, H.: Learning privacy-preserving optics for human pose estimation. In: ICCV. pp. 2573{2582 (October 2021) [2](#), [5](#), [6](#), [8](#), [10](#), [12](#)
18. Hore, A., Ziou, D.: Image quality metrics: Psnr vs. ssim. In: 2010 20th international conference on pattern recognition. pp. 2366{2369. IEEE (2010) [2](#), [9](#)
19. Huang, G.B., Mattar, M., Lee, H., Learned-Miller, E.: Learning to align from scratch. In: NIPS (2012) [11](#), [12](#)
20. Ji, X., Cheng, J., Feng, W., Tao, D.: Skeleton embedded motion body partition for human action recognition using depth sequences. *Signal Processing* **143**, 56{68 (2018) [2](#)
21. Junejo, I.N., Dexter, E., Laptev, I., Perez, P.: View-independent action recognition from temporal self-similarities. *IEEE transactions on pattern analysis and machine intelligence* **33**(1), 172{185 (2010) [7](#)
22. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 1725{1732 (2014) [3](#)
23. Kopuklu, O., Kose, N., Gunduz, A., Rigoll, G.: Resource efficient 3d convolutional neural networks. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (2019) [3](#)
24. Krishna, R., Gordon, M., Fei-Fei, L., Bernstein, M.: Visual intelligence through human interaction. In: Artificial Intelligence for Human Computer Interaction: A Modern Approach, pp. 257{314. Springer (2021) [1](#)
25. Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: Hmdb: a large video database for human motion recognition. In: 2011 International conference on computer vision. pp. 2556{2563. IEEE (2011) [8](#)
26. Kupyn, O., Martyniuk, T., Wu, J., Wang, Z.: Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8878{8887 (2019) [12](#)
27. Lakshminarayanan, V., Fleck, A.: Zernike polynomials: a guide. *Journal of Modern Optics* **58**(7), 545{561 (2011) [6](#)
28. Liu, B., Adeli, E., Cao, Z., Lee, K.H., Shenoi, A., Gaidon, A., Niebles, J.C.: Spatiotemporal relationship reasoning for pedestrian intent prediction. *IEEE Robotics and Automation Letters* **5**(2), 3485{3492 (2020) [1](#)
29. Marquez, M., Meza, P., Arguello, H., Vera, E.: Compressive spectral imaging via deformable mirror and colored-mosaic detector. *Optics express* **27**(13), 17795{17808 (2019) [13](#)
30. Marquez, M., Meza, P., Rojas, F., Arguello, H., Vera, E.: Snapshot compressive spectral depth imaging from coded aberrations. *Optics Express* **29**(6), 8142{8159 (2021) [13](#)
31. Metzler, C.A., Ikoma, H., Peng, Y., Wetzstein, G.: Deep optics for single-shot high-dynamic-range imaging. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020) [2](#)
32. Orekondy, T., Schiele, B., Fritz, M.: Towards a visual privacy advisor: Understanding and predicting privacy risks in images. In: Proceedings of the IEEE international conference on computer vision. pp. 3686{3695 (2017) [8](#), [9](#)
33. Padilla-Lopez, J.R., Chaaraoui, A.A., Florez-Revuelta, F.: Visual privacy protection methods: A survey. *Expert Systems with Applications* **42**(9), 4177{4195 (2015) [1](#), [3](#)
34. Panagiotakis, C., Karvounas, G., Argyros, A.: Unsupervised detection of periodic segments in videos. In: 2018 25th IEEE International Conference on Image Processing (ICIP). pp. 923{927. IEEE (2018) [7](#)

35. Pareek, P., Thakkar, A.: A survey on video-based human action recognition: recent updates, datasets, challenges, and applications. *Artificial Intelligence Review* **54**(3), 2259{2322 (2021) [3](#)
36. Pittaluga, F., Koppal, S., Chakrabarti, A.: Learning privacy preserving encodings through adversarial training. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 791{799. IEEE (2019) [3](#)
37. Pittaluga, F., Koppal, S.J.: Privacy preserving optics for miniature vision sensors. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 314{324 (2015) [2](#), [4](#), [10](#), [11](#)
38. Pittaluga, F., Koppal, S.J.: Pre-capture privacy for small vision sensors. *IEEE transactions on pattern analysis and machine intelligence* **39**(11), 2215{2226 (2016) [4](#)
39. Purwanto, D., Renanda Adhi Pramono, R., Chen, Y.T., Fang, W.H.: Extreme low resolution action recognition with spatial-temporal multi-head self-attention and knowledge distillation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops. pp. 0{0 (2019) [4](#)
40. Ren, Z., Lee, Y.J., Ryoo, M.S.: Learning to anonymize faces for privacy preserving action detection. In: ECCV. pp. 620{636 (2018) [1](#), [3](#)
41. Ryoo, M.S., Kim, K., Yang, H.J.: Extreme low resolution activity recognition with multi-siamese embedding learning. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018) [2](#), [10](#)
42. Ryoo, M.S., Rothrock, B., Fleming, C., Yang, H.J.: Privacy-preserving human activity recognition from extreme low resolution. In: AAAI (2017) [2](#), [4](#), [11](#)
43. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4510{4520 (2018) [9](#)
44. Sitzmann, V., Diamond, S., Peng, Y., Dun, X., Boyd, S., Heidrich, W., Heide, F., Wetzstein, G.: End-to-end optimization of optics and image processing for achromatic extended depth of field and super-resolution imaging. *ACM TOG* (2018) [2](#), [5](#)
45. Sun, C., Junejo, I.N., Tappen, M., Foroosh, H.: Exploring sparseness and self-similarity for action recognition. *IEEE Transactions on Image Processing* **24**(8), 2488{2501 (2015) [7](#)
46. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1{9 (2015) [9](#)
47. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2818{2826 (2016) [9](#)
48. Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: Mnasnet: Platform-aware neural architecture search for mobile. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2820{2828 (2019) [9](#)
49. Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M.: Learning spatiotemporal features with 3d convolutional networks. In: Proceedings of the IEEE international conference on computer vision. pp. 4489{4497 (2015) [3](#), [6](#), [9](#)
50. Tran, D., Wang, H., Torresani, L., Feiszli, M.: Video classification with channel-separated convolutional networks. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5552{5561 (2019) [3](#)

51. Upton, G.J.: Fisher's exact test. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **155**(3), 395{402 (1992) [11](#)
52. Van Der Maaten, L., Postma, E., Van den Herik, J., et al.: Dimensionality reduction: a comparative. *J Mach Learn Res* **10**(66-71), 13 (2009) [1](#)
53. Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Van Gool, L.: Temporal segment networks: Towards good practices for deep action recognition. In: *European conference on computer vision*. pp. 20{36. Springer (2016) [3](#)
54. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600{612 (2004) [6](#)
55. Wang, Z.W., Vineet, V., Pittaluga, F., Sinha, S.N., Cossairt, O., Bing Kang, S.: Privacy-preserving action recognition using coded aperture videos. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. pp. 0{0 (2019) [4](#)
56. Wu, Z., Wang, H., Wang, Z., Jin, H., Wang, Z.: Privacy-preserving deep action recognition: An adversarial learning framework and a new dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020) [1](#), [3](#), [4](#), [8](#), [9](#), [10](#), [12](#)
57. Wu, Z., Wang, Z., Wang, Z., Jin, H.: Towards privacy-preserving visual recognition via adversarial training: A pilot study. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 606{624 (2018) [1](#), [4](#), [9](#)
58. Zagoruyko, S., Komodakis, N.: Wide residual networks. *arXiv preprint arXiv:1605.07146* (2016) [9](#)