

# OPOM: Customized Invisible Cloak towards Face Privacy Protection

Yaoyao Zhong, Weihong Deng

**Abstract**—While convenient in daily life, face recognition technologies also raise privacy concerns for regular users on the social media since they could be used to analyze face images and videos, efficiently and surreptitiously without any security restrictions. In this paper, we investigate the face privacy protection from a technology standpoint based on a new type of customized cloak, which can be applied to all the images of a regular user, to prevent malicious face recognition systems from uncovering their identity. Specifically, we propose a new method, named one person one mask (OPOM), to generate person-specific (class-wise) universal masks by optimizing each training sample in the direction away from the feature subspace of the source identity. To make full use of the limited training images, we investigate several modeling methods, including affine hulls, class centers and convex hulls, to obtain a better description of the feature subspace of source identities. The effectiveness of the proposed method is evaluated on both common and celebrity datasets against black-box face recognition models with different loss functions and network architectures. In addition, we discuss the advantages and potential problems of the proposed method. In particular, we conduct an application study on the privacy protection of a video dataset, Sherlock, to demonstrate the potential practical usage of the proposed method.

**Index Terms**—Privacy protection, adversarial example, class-universal attack

## 1 INTRODUCTION

DEEP learning has achieved considerable success in computer vision [1], [2], [3], [4], significantly improving the state-of-art of face recognition [5], [6], [7], [8], [9], [10], [11], [12], [13]. This ubiquitous technology is now used to create innovative applications for entertainment and commercial services.

However, face recognition can be used for good as well as ill. Due to privacy concerns and fears of an Orwellian surveillance society [14], exploiting and analyzing face images without restriction has promoted considerable controversy. Unfortunately, such concerns and fears are not without merit. Human behavioral and psychological traits on social media and video surveillance, including photographs and videos, are at risk of unauthorized access or can be used inappropriately or with unintended disclosure. Privacy laws can be used to regulate the misuse of face recognition systems [15], [16] to avoid the potential risk of inappropriate usage of private information. However, simply prohibiting face recognition technique due to privacy concerns is not the best way. It may be better to develop deep learning technologies to solve the problems that the technology itself brings [17].

To protect privacy while maintaining practical usage, some studies [17], [18], [19], [20], [21] aim to de-identify face images such that many facial characteristics remain, but the identities of people in the images cannot be reliably recognized. Considering that the generated images may have different visual appearances compared with the original images [17], [18], or behave unnaturally and exhibit undesirable artifacts [19], [20], [21], some recent methods [22], [23], [24], [25], [26], [27] can both hide the identification information and maintain the visual quality of face images,

by generating imperceptible adversarial examples [28] as privacy protection masks. Despite the naturalness and effectiveness, it is incredibly unfriendly for regular users to generate different privacy masks for each photograph or each frame of videos.

Further study is necessary to provide more effective and simple face privacy protections for the general public. In contrast to existing methods, which generate different adversarial masks for the different face images of a person, we aim to generate a type of person-specific (class-wise) universal mask. In this way, a regular user can generate one privacy mask only once, then apply it to all his or her photographs and videos.

Compared with image-specific privacy masks, the benefits of person-specific universal masks are twofold. First, the person-specific mask is generated once, therefore it can dispense with the mask generation time for new images, which may benefit average users and some real-time privacy protection applications in terms of efficiency. Second, compared with image-specific privacy masks which require multiple transmissions of new images between the user and the server, person-specific masks need only one transmission, which can reduce the risk of privacy leakage.

There are two challenges in generating person-specific privacy masks. (1) Individual Universality. Compared with image-specific privacy masks [22], [23], [24], [25], [26], [27], person-specific privacy masks are generated with only a small number of images of an identity, and they can be applied to different unknown images. While face images of the same identity can vary from poses, illuminations, expressions and occlusions [29]. This diversity will undoubtedly increase the difficulty of generating privacy masks for different face images. Therefore, it is crucial to increase the individual universality of adversarial masks for person-specific privacy protection. (2) Model Transferability. The

The authors are with Beijing University of Posts and Telecommunications, Beijing, China (e-mail: zhongyaoyao@bupt.edu.cn; whdeng@bupt.edu.cn). Datasets and code are available at <https://github.com/zhongyy/OPOM>.

privacy masks should be transferable towards different models [30], [31], [32], [33], which means that they are generated from the local surrogate models and applied to different unknown recognition systems. For unknown face recognition models, there are a wide range of options for training databases, training loss functions and network architectures [34]. This variety will undoubtedly increase the difficulty of generating transferable adversarial masks.

In this paper, we propose a method, named one person one mask (OPOM), to provide one privacy mask, for all the face images of one person, similar to a customized invisible cloak, as shown in Figure 1. Specifically, to increase the individual universality, OPOM generates a privacy mask by solving an optimization problem, which maximizes the distance between diverse deep features of the training images and the feature subspace of the identity. We investigate different modeling methods, including affine hulls, class centers and convex hulls, to model the feature subspace of each identity for better description with limited images. As shown in Figure 2, we empirically find that, with better description of the feature subspace in the mask generation process, both individual universality and model transferability can be improved. Furthermore, to increase model transferability, OPOM is also designed to be combined with a variety of model transferability methods, such as the momentum boosting method [32] and DFANet [34]. The main contributions of our paper are as follows:

- We reveal the existence of a new type of person-specific (class-wise) universal adversarial privacy mask, which is generated to protect different face images of the same identity, and therefore can be easier-to-use for regular users.
- We investigate this new type of person-specific privacy mask, and propose an efficient method, OPOM, to generate this type of adversarial mask, which can jointly improve the *Image Universality* and *Model Transferability*, therefore attaining more effective privacy protection.
- The effectiveness of the proposed OPOM method is empirically demonstrated in protecting unconstrained face images towards different black-box deep face recognition models compared with previous universal adversarial perturbations.

The remainder of the paper is organized as follows. Section 2 briefly reviews the literature related to privacy protection and adversarial attacks. In Section 3, we first introduce the proposed OPOM for generating privacy masks; and then describe the potential comparable methods. In Section 4, we first evaluate the effectiveness of the proposed method; and then discuss the proposed method for a deeper understanding. Finally, Section 5 summarizes this paper and provides suggestions for future works.

## 2 RELATED WORKS

In this section, we briefly review the literature related to face privacy protection and adversarial attacks.

### 2.1 Face privacy protection

Privacy protection has been studied for a long time in the literature [35], and has attracted more attention in the

current deep learning era, since multimedia data can be analyzed accurately and efficiently with high-performance deep models [36]. The face image, as an important type of biometric data, can reveal a large amount of identity information.

With the help of deep learning technologies, face recognition has developed with unprecedented success [5], [6], [7], [8], [9], [10], [11], [12], [13]. Face recognition models are trained on large-scale training databases [29], [37], [38], and used as feature extractors to test identities that are usually disjoint from the training set [8]. With these open-set face recognition models [39], images and short videos posted casually on the social media can be collected and analyzed by unauthorized commercial services.

To protect privacy, studies [17], [18], [19], [20], [21] on face de-identification aim to remove identity information from the images. Newton *et al.* [17] replaced the original images with the average of  $k$  face images, and Gross *et al.* [18] reconstructed a de-identified face using a generative multi-factor model, which were early attempts to face de-identification to replace masking, blurring, or pixelation techniques. In recent years, motivated by generative neural networks, face de-identification works [19], [20], [21] have developed the synthesis images to replace the original images for privacy protection. Despite their effectiveness, considering the lack of naturalness of replaced faces in previous methods [17], [18], [19], [20], [21], some researchers [22], [23], [24], [25], [26], [27] are inspired by adversarial examples [28] and made some successful attempts at face adversarial privacy protection.

### 2.2 Adversarial Attacks

Szegedy *et al.* [28] first found that, with elaborate strategies, deep neural networks can be easily fooled by test images with imperceptible noise [28], named adversarial examples. These images could be further classified as image-dependent adversarial examples and image-agnostic adversarial examples (universal adversarial examples).

#### 2.2.1 Image-dependent adversarial examples

Compared with the initial adversarial examples generated by a box-constrained LBFGS method [28], Goodfellow *et al.* [40] proposed a more time-saving method, FGSM, which builds an adversarial example by performing one-step gradient updating along the direction of the sign of the gradient at each pixel. Kurakin *et al.* [41] further developed FGSM to generate iterative attacks, and achieved a higher attack success rate than FGSM in a white-box setting. Similarly, another efficient iterative attack method proposed by Moosavi-Dezfooli *et al.* [42], called DeepFool, generated the minimal perturbation in each step by moving towards the lineared decision boundary [42]. In addition to white-box attacks, transferable black-box attacks [30], [31], [32], [33], [43], [44] are more practical in real-world situations, since this type of attack can be applied in a fully black-box manner without any queries on the target system.

#### 2.2.2 Universal adversarial examples

Moosavi-Dezfooli *et al.* [45] proved the existence of image-agnostic adversarial attacks, which can create a universal

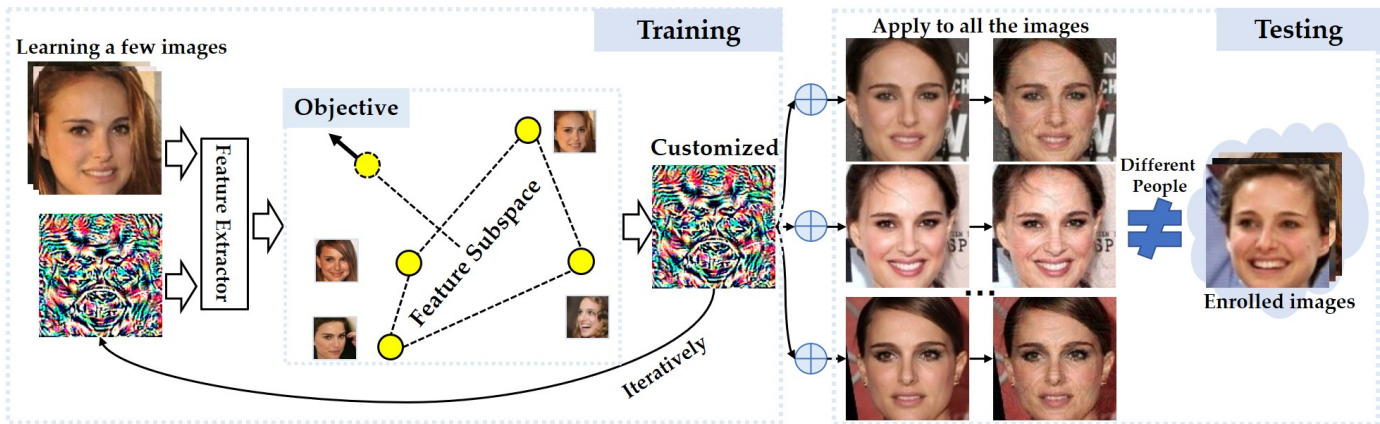


Fig. 1. Illustration of person-specific (class-wise) universal privacy masks. A regular user can generate one privacy mask only once, and then apply it to all his or her photographs and videos. In the training process for mask generation, the proposed OPOM method optimizes each training sample in the direction away from the feature subspace of the source identity. To make full use of the limited training images, several modeling methods, including affine hulls, class centers and convex hulls, have been investigated to obtain a better description of the feature subspace of the source identities. In the testing process, with the customized mask, the protected images will not be recognized as the original identity. Compared with image-specific protection, person-specific masks protect the privacy of regular users in a friendlier way.

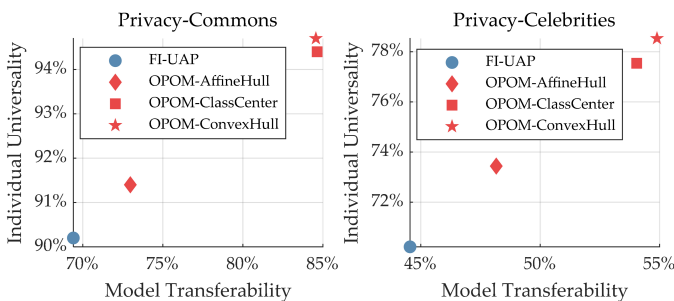


Fig. 2. With better approximation methods (red) for the feature subspace of each identity in the mask generation process, both individual universality and model transferability can be improved, which will lead to more effective privacy protection. Experimental results on two datasets (Privacy-Commons and Privacy-Celebrities) are shown.

adversarial perturbation (UAP) for all the images of the database. Hayes *et al.* [46] investigated the capacity for generative models to craft UAPs and demonstrate that their attack method improves previous methods [45] in terms of crafting more transferable UAPs. Mopuri *et al.* [47] proposed a generative method aiming to model the distribution of adversarial perturbations by generating a wide variety of UAPs. In concurrent work, Poursaeed *et al.* [48] proposed a unifying framework based on generative models for creating universal and image-dependent perturbations and improve the state-of-the-art performance of UAPs. Mopuri *et al.* [49], [50] first proposed a data-free objective, named GD-UAP, to learn perturbations that can adulterate the deep features extracted by multiple layers. GD-UAP, especially GD-UAP with data priors; it demonstrated impressive fooling rates and surprising transferability across various deep models with different architectures, regularizers and underlying tasks.

At present, the most relevant works to ours are the class-wise/discriminative methods [51], [52]. Gupta *et al.* [51] proposed a type of class-wise UAP, using the linearity of the decision boundaries of deep neural networks. Using the absolute accuracy drop as an evaluation metric, Zhang *et al.* [52]

TABLE 1  
Comparison of methods to generate privacy masks with adversarial examples to protect images against malicious face recognition systems.

Methods	Inference	Open-set	Blackbox	Person-specific (Class-wise)
Oh <i>et al.</i> [22]	Yes	No	No	No
P-FGVM [23]	Yes	No	No	No
Fawkes [24]	No	No	Yes	No
TIP-IM [25]	Yes	Yes	Yes	No
APF [26]	Yes	Yes	Yes	No
LowKey [27]	Yes	Yes	Yes	No
This paper	Yes	Yes	Yes	Yes

proposed a class-discriminative universal attack to generate a perturbation that can fool a target network to misclassify only the target classes, while having limited influence on the remaining classes. Note that a class-discriminative universal attack is not suitable for the privacy protection task since it actually sacrifices some protection effects on the target classes in exchange for inhibition effects on the remaining classes. Previous class-wise/discriminative methods [51], [52] are inspiring; however, due to a heavy reliance on class information, they cannot be applied to the open-set face recognition task directly. Even if face recognition models suffer from a lack of class information, the proposed OPOM method models the feature subspace of an identity with a data prior and therefore can generate class-wise universal adversarial perturbations by optimizing the deep features in the direction away from the feature subspace of the source identity, which is exactly the advantage of OPOM.

### 2.3 Adversarial Privacy Protection

Inspired by the adversarial examples, which can mislead deep models to generate incorrect outputs while maintaining the visual quality of images, some works consider gen-

erating adversarial perturbations of face images for privacy protection [22], [23], [24], [25], [26], [27].

Oh *et al.* [22] made an initial attempt to use adversarial perturbations for privacy protection from a game theory perspective. For better visualization effects, Chatzikiriakidis *et al.* [23] proposed P-FGVM to generate white-box privacy masks by considering minimal facial image distortion based on iterative attacks [41]. Fawkes [24] incorporated adversarial perturbations in the training process by poisoning some target training samples; and therefore misled the model on the poisoned targets in the inference period. Fawkes reported good performance on commercial APIs, which drew much attention in the research community. However, compared with the poisoning in the training process method, it is more applicable to design privacy masks that can be directly used [25], [26], [27] in terms of privacy protection for regular users. For a higher privacy protection success rate, Yang *et al.* [25] designed adversarial privacy masks that were directly used on the probe images, with a targeted optimization objective, named TIP-IM, under the assumption of unknown gallery images. To prevent privacy leakage when generating privacy masks, Zhang *et al.* [26] proposed an end-cloud collaborated adversarial attack solution, named APF, where the original face images were only available at the user end. In addition, Valeriia *et al.* [27] applied adversarial privacy masks on the gallery faces so that the probe images could not be recognized correctly.

A comparison of the aforementioned methods [22], [23], [24], [25], [26], [27] and our methods is listed in Table 1, in terms of whether the method can be directly used in the inference process, whether the method is applied in an open-set face recognition setting, whether the privacy mask can confront black-box models, and whether the privacy mask is person-specific (class-wise) universal for convenient usage. Despite the clear motivation and good performance [22], [23], [24], [25], [26], [27], to the best of our knowledge, no works before have considered the person-specific privacy protection for high-efficiency, which we believe could revolutionize the future for face privacy protection.

### 3 PERSON-SPECIFIC PRIVACY MASKS

#### 3.1 Problem Formulation

The objective of a person-specific adversarial mask is to craft a perturbation  $\Delta X$  that can fool a variety of deep face recognition models  $f(\cdot)$  for all the face images  $X^k = \{X_1^k, X_2^k, \dots, X_i^k, \dots\}$  of identity  $k$ , so that any face image  $X_i^k$  can conceal the identity  $X^k$  from deep face recognition models  $f(\cdot)$  with this privacy mask  $\Delta X$ . That is, we find a  $\Delta X$  so that  $X_i^k \in X^k$ ,

$$D(f(X_i^k + \Delta X), f_{X^k}) > t, \|\Delta X\|_\infty < \varepsilon, \quad (1)$$

where  $f(X_i^k) \in \mathbb{R}^d$  is the normalized feature of image  $X_i^k$ ,  $f_{X^k}$  denotes the feature subspace of identity  $k$ , and  $t$  is the distance threshold to decide whether a pair of face images belongs to the same identity.  $D(x_1, x_2)$  denotes the distance between  $x_1$  and  $x_2$ , that is, the shortest distance between point  $x_1$  and subspace  $x_2$ . When  $x_2$  denotes a point, we use  $D(x_1, x_2)$  to represent the normalized Euclidean distance or cosine distance commonly used in face recognition.  $\varepsilon$  limits the maximum deviation of the privacy mask.

#### 3.2 One person one mask (OPOM)

In the above problem formulation as Equation 1, the key point is the formulation of  $f_{X^k}$ . Note that for open-set face recognition models that are more like feature extractors, it is likely that the persons to be protected do not belong to the training databases. That is, we cannot obtain identity information from the face recognition models directly as previous works [51], [52] did in some close-set tasks where the class information of the model can be used to generate class-wise universal perturbations.

Therefore, the only choice is to describe an identity  $f_{X^k}$  from the provided face image set  $X^k$ . For a specific identity, when more face images are provided, the image set  $X^k$  will be more complete, and the feature subspace  $f_{X^k}$  will be more precise. However, incorporating more images for generating adversarial masks is completely at odds with our expectations, since we aim to use as few face images as we can to generate a person-specific adversarial mask for privacy protection. Considering the above function, we explore some approximation methods for  $f_{X^k}$ , with a limited number  $n_k$  of face images  $\tilde{X}^k = \{X_1^k, X_2^k, \dots, X_{n_k}^k\}$  in practice.

The basic intuition is as follows. With more precise approximation of the identity subspace, the individual universality can be enhanced to some extent. In addition, with more diverse gradient information (*i.e.*, gradients of cross-comparisons, images  $X_i^k$  and  $X_j^k$ ) incorporated into the approximation equations, the model transferability can be increased to some degree.

To confirm this intuition, we conduct an analytical experiment as shown in Figure 2. Specifically, we describe individual universality using the protection success rate of the unfamiliar images on the local surrogate model, that is, evaluate the masks with the same model and different images. Model transferability can be described as the average protection success rate of the provided training dataset on different black-box models, that is, evaluating the privacy masks with the same images and different models.

We empirically find that, compared with single points representing the identity (“FI-UAP”, represented with the blue marker), with more precise approximation (represented with the red markers), both individual universality and model transferability can be improved. The empirical experimental results are consistent with the intuition, that is, a better approximation method will lead to more effective privacy protection. Next, we introduce the approximation methods in details.

##### 3.2.1 Approximation methods of the feature subspace

**Affine Hulls.** The affine hull implicitly augments existing face images, by treating any affine combination of normalized deep features as a valid feature for the identity, so that limited face images can realize their full potential. The affine hull of deep features for a specific identity can be expressed as

$$H(f_{\tilde{X}^k}) = \left\{ x = \sum_{i=1}^{n_k} \alpha_i^k f(X_i^k) \mid \sum_{i=1}^{n_k} \alpha_i^k = 1 \right\}, \quad (2)$$

where  $f(X_i^k)$  and  $\alpha_i^k$  ( $i = 1, \dots, n_k$ ) are the features and coefficients to describe the identity  $k$ .

In this way, we can generate the privacy mask  $\Delta X$  with limited face images for identity  $k$  by optimizing the following objective,

$$\Delta X = \arg \max_{\Delta X} \sum_{i=1}^{n_k} D(f(X_i^k + \Delta X), H(f_{\tilde{X}^k})), \|\Delta X\|_\infty < \varepsilon, \quad (3)$$

where  $n_k$  denotes the number of face images,  $H(f_{\tilde{X}^k})$  is the affine hull of the normalized features as Equation 2.

To calculate the distance  $D(f(X_i^k + \Delta X), H(f_{\tilde{X}^k}))$  in Equation 3, we can rewrite  $H(f_{\tilde{X}^k})$  as follows to parametrize the affine hull,

$$H(f_{\tilde{X}^k}) = \left\{ x = U^k V^k + \mu^k \mid V^k \in \mathbb{R}^{n_k} \right\}, \quad (4)$$

where  $\mu^k = \frac{1}{n_k} \sum_{i=1}^{n_k} f(X_i^k)$ ,  $U^k \in \mathbb{R}^{d \times n_k}$  is the orthonormal basis for the directions spanned by the affine hull, obtained by using singular value decomposition (SVD) to  $[f(X_1^k) - \mu^k, \dots, f(X_{n_k}^k) - \mu^k]$ , and  $V^k \in \mathbb{R}^{n_k}$  is a vector of free parameters that provides coordinates for the orthonormal basis. With Equation 4, we can calculate the distance  $D(f(X_i^k + \Delta X), H(f_{\tilde{X}^k}))$  in Equation 3 as

$$\min_{V^k} \|U^k V^k + \mu^k - f(X_i^k + \Delta X)\|, \quad (5)$$

which can be written as a standard least squares problem

$$\min_{V^k} \|U^k V^k - (f(X_i^k + \Delta X) - \mu^k)\|, \quad (6)$$

and the solution is

$$V^k = (U^{kT} U^k)^{-1} U^{kT} (f(X_i^k + \Delta X) - \mu^k). \quad (7)$$

Finally, we can generate the privacy mask  $\Delta X$  by transforming Equation 3 as

$$\begin{aligned} \Delta X &= \arg \max_{\Delta X} J_{AH} \\ &= \arg \max_{\Delta X} \sum_{i=1}^{n_k} \|U^k V_i^k - (f(X_i^k + \Delta X) - \mu^k)\|, \|\Delta X\|_\infty < \varepsilon, \end{aligned} \quad (8)$$

where  $U^k$  and  $\mu^k$  can be calculated as mentioned above, and  $V_i^k$  is the solution of Equation 6, with the value shown in Equation 7.

**Class Centers and Convex Hulls.** Although affine hulls can provide an implicit augmentation for provided face images, this approximation may be too loose since some points in the hull may not be valid and may cause opposite effects when generating the mask. Therefore, we introduce lower and upper bounds  $L$  and  $U$  on the allowable  $\alpha_i^k$  coefficients in Equation 1 to control the looseness,

$$H(f_{\tilde{X}^k}) = \left\{ x = \sum_{i=1}^{n_k} \alpha_i^k f(X_i^k) \mid \sum_{i=1}^{n_k} \alpha_i^k = 1, L \leq \alpha_i^k \leq U \right\}, \quad (9)$$

where  $f(X_i^k)$  and  $\alpha_i^k$  ( $i = 1, \dots, n_k$ ) are the features and coefficients to describe the identity  $k$ .

The approximation will be the affine hull if  $L = -\infty$  and  $U = \infty$ . Otherwise, the incorporation of  $L$  and  $U$  can reduce the over-large region. Another interesting point is  $L = U = 1/n_k$ . In this case, the identity is exactly described as the mean feature

$$H(f_{\tilde{X}^k}) = \frac{1}{n_k} \sum_{i=1}^{n_k} f(X_i^k), \quad (10)$$

which is a popular method [7] to approximate class centers. In this way, the privacy mask  $\Delta X$  can be generated by

$$\begin{aligned} \Delta X &= \arg \max_{\Delta X} J_{CC} \\ &= \arg \max_{\Delta X} \sum_{i=1}^{n_k} \left\| \frac{1}{n_k} \sum_{j=1}^{n_k} f(X_j^k) - f(X_i^k + \Delta X) \right\|, \|\Delta X\|_\infty < \varepsilon. \end{aligned} \quad (11)$$

Note that if  $L = 0$  and  $U = 1$ , then  $H(f_{\tilde{X}^k})$  actually approximates the identity with the convex hull (the smallest convex set) of features  $f(X_i^k)$ , which is also the most effective approximation method we will empirically demonstrate. For the convex hull, *i.e.*, with Equation 9 and  $L = 0, U = 1$ , we can calculate  $D(f(X_i^k + \Delta X), H(f_{\tilde{X}^k}))$  in Equation 3 as a least squares problem with box constraints,

$$\min_{\alpha^k} \|F^k A^k - f(X_i^k + \Delta X)\|, \text{ s.t. } A^k \geq 0, 1^T A^k = 1, \quad (12)$$

where  $F^k \in \mathbb{R}^{d \times n_k}$  is a matrix with columns  $f(X_i^k)$ , and  $A^k \in \mathbb{R}^{n_k}$  is a vector containing the corresponding coefficients ( $\alpha_i^k$  of Equation 9). Then, the privacy mask  $\Delta X$  can be generated by

$$\begin{aligned} \Delta X &= \arg \max_{\Delta X} J_{CH} \\ &= \arg \max_{\Delta X} \sum_{i=1}^{n_k} \|F^k A_i^k - f(X_i^k + \Delta X)\|, \|\Delta X\|_\infty < \varepsilon, \end{aligned} \quad (13)$$

where  $A_i^k$  can be solved by Equation 12.

---

#### Algorithm 1: OPOM-AffineHull

---

**Input:** Face images  $\tilde{X}^k = \{X_1^k, X_2^k, \dots, X_{n_k}^k\}$  of identity  $k$ , deep face model  $f(\cdot)$ , maximum deviation of perturbations  $\varepsilon$ , maximum iterative steps  $N_{max}$ .

- 1 **Initialize:**  $\Delta X_0 \sim U(-\varepsilon, \varepsilon)$ ,  $N = 0$ ,  $g_N = 0$ ;
- 2 **while** *step*  $N < N_{max}$  **do**
- 3     Parametrize the affine hull with  $U^k$  and  $\mu^k$ ;
- 4     Calculate  $V_i^k$  by Equation 7;
- 5      $J_{AH} = \sum_{i=1}^{n_k} \|U^k V_i^k - (f(X_i^k + \Delta X_N) - \mu^k)\|$  (Equation 8);
- 6      $g_{N+1} = \nabla_{X^k + \Delta X_N} J_{AH}$ ;
- 7      $\Delta X_{N+1} = C_\varepsilon(\Delta X_N + \text{sign}(g_{N+1}))$ ,  $N = N + 1$ ;
- 8 **end**

**Output:** Privacy mask  $\Delta X_{N_{max}}$  for identity  $k$ .

---

### 3.2.2 Generation of Privacy Masks

With the approximation methods for modeling the feature subspace of an identity using affine hulls, class centers and convex hulls, we formulate the person-specific privacy mask generation to optimization problems as Equation 8, Equation 11 and Equation 13, named OPOM-AffineHull, OPOM-ClassCenter and OPOM-ConvexHull. The optimization problems can be solved by an iteratively signed gradient ascent as in previous works [34], which is known to be effective for generating adversarial attacks towards deep models [53].

To further detail the generation of person-specific masks, we conclude the proposed OPOM-AffineHull and OPOM-ConvexHull in Algorithm 1 and Algorithm 2, where

**Algorithm 2: OPOM-ConvexHull**


---

**Input:** Face images  $\tilde{X}^k = \{X_1^k, X_2^k, \dots, X_{n_k}^k\}$  of identity  $k$ , deep face model  $f(\cdot)$ , maximum deviation of perturbations  $\varepsilon$ , maximum iterative steps  $N_{max}$ .

- 1 **Initialize:**  $\Delta X_0 \sim U(-\varepsilon, \varepsilon)$ ,  $N = 0$ ,  $g_N = 0$ ;
- 2 **while**  $step\ N < N_{max}$  **do**
- 3     Calculate  $A_i^k$  by solving Equation 12;
- 4      $J_{CH} = \sum_{i=1}^{n_k} \|F^k A_i^k - f(X_i^k + \Delta X_N)\|$  (Equation 13);
- 5      $g_{N+1} = \nabla_{X_i^k + \Delta X_N} J_{CH}$ ;
- 6      $\Delta X_{N+1} = C_\varepsilon(\Delta X_N + sign(g_{N+1}))$ ,  $N = N + 1$ ;
- 7 **end**

**Output:** Privacy mask  $\Delta X_{N_{max}}$  for identity  $k$ .

---

$C_\varepsilon(X) = \min(\varepsilon, \max(-\varepsilon, X))$ . For brevity, we omit OPOM-ClassCenter. To calculate the distance between the protected features and the original feature subspace in each iterative step, the coordinates  $V_i^k$  or coefficients  $A_i^k$  should be calculated at the beginning of each iterative step. We provide a closed-form solution for  $V_i^k$  in Equation 7. For  $A_i^k$ , although there is no closed-form solution, it can be efficiently solved with the convex optimization toolbox CVX [54], [55].

### 3.2.3 Combination with Model Transferability Methods

To generate more transferable person-specific privacy masks from a single source model to protect face images against black-box models, model transferability methods, such as the momentum boosting method [32] and DFANet [34], can be incorporated into OPOM method.

The momentum boosting method [32] integrates the momentum term into the attack process to stabilize the update directions and escape from poor local maxima. Specifically,  $g_{N+1}$  in Algorithm 1 and Algorithm 2 gathers gradients of the first  $N$  iterations with a decay factor  $\mu$ ,

$$g_{N+1} = \mu \cdot g_N + \frac{\nabla_{X_i^k + \Delta X_N} J}{\|\nabla_{X_i^k + \Delta X_N} J\|}, \quad (14)$$

where  $\mu$  is usually set to 1 following the original paper.

DFANet [34] converts the surrogate model  $f(\cdot)$  in the  $N$ -th step of the adversarial example generation to different models  $\tilde{f}^N(x)$  of the  $N$ -th step by incorporating dropout layers. Specifically, for a face recognition model  $f(\cdot)$  composed of convolutional layers, given the output  $o_i \in \mathbb{R}^n$  from the  $i$ -th convolutional layer, a mask  $M_i \in \mathbb{R}^n$  is generated with each element  $m_i$  independently sampled from a Bernoulli distribution with probability  $p_d$ :

$$m_i \sim \text{Bernoulli}(p_d), \quad m_i \in M_i. \quad (15)$$

Then, this mask is applied to modify the output as  $o_i = M_i \times o_i$ , where  $\times$  denotes the Hadamard product.

## 3.3 Comparison methods

In addition to OPOM, in this paper we explore other methods that can generate universal adversarial perturbations or class-universal adversarial perturbations, for the potential capabilities in the person-specific privacy protection task.

### 3.3.1 GD-UAP

To generate universal perturbations, GD-UAP seeks  $\Delta X$  to produce maximal spurious activations at each layer of a model. Although GD-UAP is a data-free optimization for crafting image-agnostic perturbations, according to the original paper [50], for the best performance of GD-UAP, it can be utilized with simple data priors, such as target data samples. Therefore, for a fair comparison with OPOM,  $\tilde{X}^k = \{X_1^k, X_2^k, \dots, X_{n_k}^k\}$  is used as the data priors to generate a privacy mask for identity  $k$ ,

$$\begin{aligned} \Delta X &= \arg \min_{\Delta X} J_{GD-UAP} \\ &= \arg \min_{\Delta X} - \sum_i \log \left( \prod_{j=1}^K \|l_j(X_i^k + \Delta X)\| \right), \|\Delta X\|_\infty < \varepsilon, \end{aligned} \quad (16)$$

where  $l_j(X_i^k + \Delta X)$  is the activation in the output tensor (after the non-linearity) at layer  $j$  when  $X_i^k + \Delta X$  is fed to the network  $f(\cdot)$ , and  $K$  is the number of layers.

### 3.3.2 FI-UAP, FI-UAP+ and FI-UAP-all

One straightforward idea to generate class-wise universal adversarial perturbations towards open-set face recognition models is actually a similar way in the spirit of UAP [45], which leverages a set of images to seek the minimal perturbation iteratively by DeepFool [42] and aggregates them to the universal perturbations. The feature iterative attack method (FIM) [34], [56], which has been demonstrated to be an applicable method for face recognition models, can also be used in this way. Specifically, we aggregate gradients of  $\tilde{X}^k = \{X_1^k, X_2^k, \dots, X_{n_k}^k\}$  to generate person-specific masks,

$$\begin{aligned} \Delta X &= \arg \max_{\Delta X} J_{FI-UAP} \\ &= \arg \max_{\Delta X} \sum_{i=1}^{n_k} \|f(X_i^k) - f(X_i^k + \Delta X)\|, \|\Delta X\|_\infty < \varepsilon, \end{aligned} \quad (17)$$

and name this method FI-UAP, which can be considered one approach for using the single point to model the identity information.

FI-UAP can be enhanced by incorporating intra-class interactions, which we refer to as FI-UAP+.

$$\begin{aligned} \Delta X &= \arg \max_{\Delta X} J_{FI-UAP+} \\ &= \arg \max_{\Delta X} \sum_{i=1}^{n_k} \sum_{j=1}^{n_k} \|f(X_i^k) - f(X_j^k + \Delta X)\|, \|\Delta X\|_\infty < \varepsilon. \end{aligned} \quad (18)$$

Furthermore, FI-UAP can be extended to generate a universal mask if all the training images are used, which we refer to as FI-UAP-all.

### 3.3.3 GAP

Generative adversarial perturbations (GAP) [48] pass a fixed pattern sampled from a uniform distribution through the generator to generate universal adversarial perturbations, which can be added to a set of images to mislead models. We strictly follow the original paper using the ResNet Generator but replace the label-level adversarial losses as feature-level ones, as Equation 17, to make it more suitable for face recognition models [34].

## 4 EXPERIMENTS

In this section, we first introduce the experimental settings. Next, we report the protection performance of the comparison methods and the proposed OPOM method, explore the performance of OPOM combined with the transferability methods, and show the protection effects against commercial APIs. Finally, we discuss the proposed method in the following aspects: analyses of the failure cases, the strengths and weaknesses of the person-specific (class-wise) universal masks compared with image-specific masks and universal masks, the practicality of OPOM in the privacy protection of a video dataset, and the diversity of privacy masks for the potential leakage of the used privacy masks.

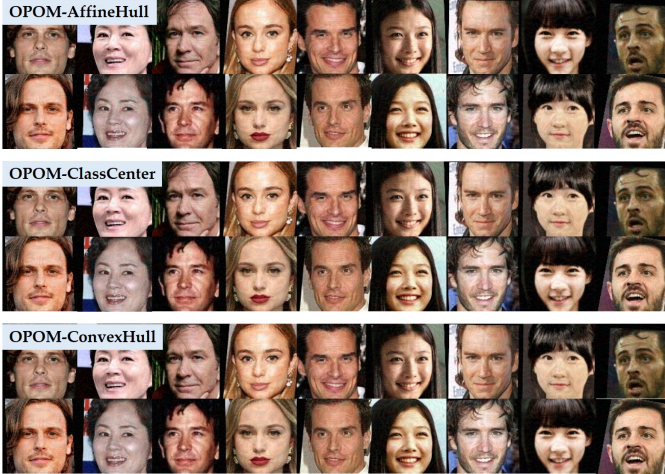


Fig. 3. Some protected images in the Privacy-Celebrities dataset with masks ( $\epsilon = 8$ ), generated by OPOM-AffineHull, OPOM-ClassCenter, and OPOM-ConvexHull respectively.

### 4.1 Experimental Settings

#### 4.1.1 Datasets and Evaluation Metrics

In this paper, we simulate a real-world scenario where regular users would like to provide a few images to generate person-specific privacy protection masks. Person-specific masks are generated on the source model with the provided training images. Then, the generated privacy masks can be applied to other test images to protect them against different black-box face recognition models.

We expect to investigate person-specific adversarial privacy masks, and therefore need identities with relatively enough images to evaluate our method. We build two datasets using common people and celebrities to simulate the real-world situation, and mainly report the 1:N identification performance for measuring the privacy protection rate. Specifically, the 1:N identification aims to identify the image with the same identity as the probe image in the gallery set. Given a probe image and a gallery containing one photo of the same person, the algorithm rank-orders all the photos in the gallery based on the feature distances to the probe. We test each of the  $M$  images per person by adding it to the gallery of distractors and use each of the other  $M-1$  images as a probe. Note that we apply the privacy masks to the probe images, not the gallery images, since we assume that some unauthorized face recognition

services may have already obtained the original images of the regular users. In this way, new images cannot be recognized and analyzed, since a wrong image in the gallery set may be matched with the mask-protected probe image.

We report the Top-1 and Top-5 protection success rate (100% - Top1 or Top 5 accuracy), which means that the Top-1 and Top-5 images do not have the same identity as the probe. Briefly, higher protection success rate is better. Note that, the more easily a face image is recognized, the more difficult it is to be protected. If a person cannot be recognized, there is no need for us to use a privacy mask to protect it anymore. Therefore, face images in the test dataset should be recognized normally by face recognition models; that is, the initial protection success rate should be almost zero without privacy masks.

**Privacy-Commons dataset.** We select 500 persons in the MegaFace challenge2 [57] database, who have at least 15 clean images (without label noise). For each person, we use at most 10 images as training images to generate the adversarial privacy mask; while another 5 images are used as the testing images to evaluate the protection performance. For the 1:N identification, the selected 500 persons with 5 test images are used as the probe set. Correspondingly, we conduct  $500 \times 5 \times 4 = 10,000$  tests. 10,000 distractors, *i.e.*, individuals who are not in the probe set, serve as other persons in the gallery set. Note that distractors belong to different persons from the training databases (therefore, different from the 500 persons in the probe set).

**Privacy-Celebrities dataset.** We select 500 persons in the one-million celebrity list, avoiding repetition of the MS-Celeb-1M [38] and LFW database [58]. For each person, there are at least 20 clean images (without label noise). We use at most 10 images as training images to generate the adversarial privacy mask, while another 10 images are used as the testing images to evaluate the protection performance. For the 1:N identification, the selected 500 persons with 10 test images are used as the probe set. Correspondingly, we conduct  $500 \times 10 \times 9 = 45,000$  tests. To reduce the domain gap between probe and gallery set, we use a celebrity database, 13,233 images of LFW database [58], as distractors, which serve as other persons in the gallery set.

#### 4.1.2 Face Recognition Models

To simulate potential scenarios, we increase the difference between the source (training) and target (testing) models by attacking deep face recognition models with different training loss functions and network architectures.

Three source (training) models to generate privacy masks are the modified version [11] of ResNet-50 [3] trained on the CAISA-WebFace database [37], supervised by Softmax loss, CosFace [10], and ArcFace [11]. The aim of privacy protection masks is to protect the original face images from being recognized by black-box models. Here, we use six target (testing) models different from the source models. Three of the black-box models differ in loss functions, *i.e.*, CosFace [10], ArcFace [11] and SFace [13]. The other three of the black-box models differ in network architectures, *i.e.*, the modified version [11] of the squeeze-and-excitation network (SENet) [59], MobileNet [60], and Inception-ResNet [61]. The recognition performance of the source and black-box models can be found in the Appendix.

TABLE 2

Comparison of different methods to generate person-specific privacy masks ( $\varepsilon = 8$ ) from a single source model to protect face images against black-box models. We report Top-1 and Top-5 protection success rate (%) under 1:N identification setting of the Privacy-Commons dataset. The higher protection success rate is better.

Source	Method	Target											
		ArcFace		CosFace		SFace		MobileNet		SENet		Inception-ResNet	
		Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Softmax	GD-UAP	6.3	2.7	3.1	1.6	3.0	1.5	8.7	3.5	5.8	2.7	3.0	1.5
	GAP	30.7	20.7	19.7	13.9	23.9	15.9	33.6	20.8	31.0	20.1	13.1	7.4
	FI-UAP	72.3	62.4	63.5	53.3	70.3	61.9	73.9	61.8	77.4	67.6	52.4	40.7
	FI-UAP+	76.9	67.8	69.2	60.3	75.1	67.3	78.3	67.8	82.1	73.3	57.2	45.6
	FI-UAP-all	53.4	43.2	38.5	30.3	39.8	31.0	54.0	39.4	51.3	38.9	24.2	17.1
	OPOM-AffineHull	73.0	63.1	63.3	53.8	71.1	62.0	74.7	63.0	77.8	68.6	52.7	40.8
	OPOM-ClassCenter	76.9	67.8	69.2	60.3	75.0	67.3	78.4	67.9	82.1	73.4	57.1	45.6
OPOM-ConvexHull	<b>78.0</b>	<b>69.4</b>	<b>70.2</b>	<b>61.4</b>	<b>76.1</b>	<b>68.7</b>	<b>79.2</b>	<b>69.1</b>	<b>82.9</b>	<b>74.2</b>	<b>58.7</b>	<b>47.2</b>	
ArcFace	GD-UAP	2.6	1.0	1.1	0.6	1.1	0.4	2.9	1.0	1.9	0.8	1.5	0.7
	GAP	52.7	36.7	36.7	26.1	41.8	30.6	51.8	36.6	49.5	36.1	28.1	18.9
	FI-UAP	82.3	75.0	71.2	63.6	77.3	70.0	65.2	50.4	73.9	63.1	56.6	45.3
	FI-UAP+	85.9	79.6	76.2	69.5	82.0	75.7	69.8	56.7	78.7	69.8	62.0	50.9
	FI-UAP-all	60.5	50.9	45.1	35.6	49.1	39.8	46.0	33.1	48.4	36.1	26.6	19.0
	OPOM-AffineHull	82.9	75.8	72.3	64.8	78.1	70.9	66.4	52.2	75.4	64.5	58.4	46.8
	OPOM-ClassCenter	86.0	79.5	76.1	69.6	82.0	75.8	69.6	56.7	78.6	69.7	61.8	51.0
OPOM-ConvexHull	<b>86.5</b>	<b>80.1</b>	<b>76.8</b>	<b>70.0</b>	<b>82.7</b>	<b>76.5</b>	<b>70.5</b>	<b>57.3</b>	<b>79.3</b>	<b>70.2</b>	<b>63.2</b>	<b>52.2</b>	
CosFace	GD-UAP	3.7	1.4	1.1	0.4	1.3	0.4	4.0	1.5	3.2	1.0	1.8	0.6
	GAP	59.7	46.7	40.2	28.4	38.8	27.0	41.5	28.4	47.2	34.0	27.5	18.0
	FI-UAP	80.4	72.3	72.8	65.1	76.7	68.9	55.6	41.0	66.9	54.6	49.9	37.1
	FI-UAP+	85.4	78.5	78.6	71.4	82.4	75.4	61.4	47.6	73.0	61.7	55.3	43.0
	FI-UAP-all	64.2	51.7	48.6	38.4	52.1	40.8	49.6	36.6	56.4	43.3	29.3	20.5
	OPOM-AffineHull	82.4	74.5	74.3	67.0	78.4	71.3	58.6	43.8	69.0	56.8	51.8	39.4
	OPOM-ClassCenter	85.4	78.3	78.8	71.6	82.4	75.5	61.4	47.5	73.0	61.5	55.2	43.0
OPOM-ConvexHull	<b>86.6</b>	<b>79.3</b>	<b>79.5</b>	<b>72.7</b>	<b>83.0</b>	<b>76.3</b>	<b>62.3</b>	<b>48.2</b>	<b>74.0</b>	<b>63.1</b>	<b>56.8</b>	<b>44.5</b>	

TABLE 3

Comparison of different methods to generate person-specific privacy masks ( $\varepsilon = 8$ ) from a single source model to protect face images against black-box models. We report Top-1 and Top-5 protection success rate under 1:N identification setting of the Privacy-Celebrities dataset. The higher protection success rate is better.

Source	Method	Target											
		ArcFace		CosFace		SFace		MobileNet		SENet		Inception-ResNet	
		Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Softmax	GD-UAP	6.5	2.4	3.4	1.2	3.6	1.3	10.7	3.8	5.7	2.1	3.1	0.9
	GAP	41.5	31.3	31.7	22.2	35.0	25.2	53.2	40.7	44.5	33.5	20.2	12.1
	FI-UAP	56.5	45.5	45.2	34.9	52.9	42.4	60.2	46.9	61.6	49.8	37.4	26.5
	FI-UAP+	62.4	51.6	51.8	42.3	59.8	49.7	66.2	54.7	67.8	57.5	42.5	31.4
	FI-UAP-all	47.6	37.9	38.4	28.9	42.5	33.3	59.0	47.8	53.9	43.3	24.4	16.2
	OPOM-AffineHull	58.8	47.9	47.0	37.2	54.8	44.9	62.4	49.3	63.5	52.1	39.0	27.8
	OPOM-ClassCenter	62.3	51.5	51.8	42.5	59.8	49.8	66.2	54.8	67.7	57.7	42.4	31.4
OPOM-ConvexHull	<b>63.8</b>	<b>53.5</b>	<b>53.6</b>	<b>44.0</b>	<b>61.3</b>	<b>51.7</b>	<b>67.0</b>	<b>55.7</b>	<b>68.9</b>	<b>58.9</b>	<b>44.2</b>	<b>33.2</b>	
ArcFace	GD-UAP	6.5	2.4	3.4	1.2	4.0	1.3	10.0	3.8	6.4	2.5	3.8	1.3
	GAP	52.4	41.5	40.1	30.2	45.4	36.1	55.9	43.4	50.8	39.5	36.5	26.3
	FI-UAP	62.5	51.5	47.9	37.8	55.6	45.8	50.0	35.6	54.3	41.0	37.7	26.3
	FI-UAP+	68.7	59.3	55.2	46.1	63.0	54.2	56.4	42.5	60.9	49.0	43.6	31.8
	FI-UAP-all	57.2	46.6	43.4	33.4	50.6	40.9	<b>58.8</b>	<b>45.7</b>	57.3	46.0	35.1	25.1
	OPOM-AffineHull	65.8	55.2	51.3	41.3	59.1	49.6	53.1	39.0	57.7	45.1	40.5	28.4
	OPOM-ClassCenter	68.9	59.3	55.2	46.1	63.0	54.3	56.6	42.6	61.1	49.0	43.8	32.0
OPOM-ConvexHull	<b>69.9</b>	<b>60.4</b>	<b>56.2</b>	<b>47.3</b>	<b>64.1</b>	<b>55.1</b>	57.5	43.1	<b>62.3</b>	<b>50.0</b>	<b>44.9</b>	<b>33.2</b>	
CosFace	GD-UAP	7.4	3.0	3.7	1.3	4.2	1.4	10.2	4.2	6.3	2.4	3.8	1.3
	GAP	53.3	42.7	40.4	30.5	41.6	31.7	48.9	35.9	48.1	36.9	24.8	15.0
	FI-UAP	63.7	52.1	51.5	41.3	57.6	47.2	46.6	31.7	50.8	37.1	33.0	21.9
	FI-UAP+	68.6	58.6	58.4	48.9	63.8	54.7	51.2	36.4	57.2	43.9	38.6	26.5
	FI-UAP-all	58.9	48.2	41.4	31.7	49.2	39.5	<b>55.7</b>	<b>43.5</b>	57.1	<b>46.6</b>	33.2	23.8
	OPOM-AffineHull	66.5	55.6	55.4	45.3	61.1	51.4	49.5	34.7	54.7	40.7	35.7	23.9
	OPOM-ClassCenter	68.3	58.4	58.5	48.9	63.8	54.5	51.2	36.3	57.0	43.9	38.5	26.5
OPOM-ConvexHull	<b>69.9</b>	<b>60.6</b>	<b>59.8</b>	<b>50.5</b>	<b>65.6</b>	<b>56.2</b>	53.0	38.1	<b>58.6</b>	45.5	<b>40.2</b>	<b>28.0</b>	



## 4.2 Protection Performance towards Black-box Face Recognition Models

### 4.2.1 Effectiveness of OPOM

We evaluate the proposed OPOM method, including OPOM-AffineHull, OPOM-ClassCenter, and OPOM-ConvexHull, on the Privacy-Commons dataset and Privacy-Celebrities dataset. In addition, we also implement the comparison methods, *i.e.*, GD-UAP [50], FI-UAP, FI-UAP+, FI-UAP-all, and GAP [48] as described in Section 3.3, to explore their potential capabilities in this person-specific privacy protection task.

Some details are as follows. The maximum deviation of perturbations  $\varepsilon$  is set to 8 under the  $L_\infty$  constraint. For GD-UAP, we optimize all the residual blocks at the last layers and the independent convolutional layers following the original paper [50]. Correspondingly, the number of iterations for the FI-UAP, FI-UAP+, FI-UAP-all and OPOM methods are chosen to be 16, the maximum number of iterations for GD-UAP method is set as 10,000. For GAP, all the images of training datasets are applied to train the Generator for 10 epochs.

The Top-1 and Top-5 protection success rates (%) for the 1:N identification on the Privacy-Commons dataset are reported in Table 2. The protection success rates on the Privacy-Celebrities dataset are shown in Table 3. We can see that OPOM is a more appropriate method for the person-specific privacy protection task. Some protected images with masks generated by OPOM are shown in Figure 3.

Note that FI-UAP can also be understood as a special case for Equation 13, where  $\alpha_i^k = 1$  and  $\alpha_j^k = 0, j \neq i$  for  $A_i^k$ . FI-UAP+ can be considered as a similar method as OPOM-ClassCenter. From the comparison of FI-UAP, FI-UAP+, OPOM-AffineHull, OPOM-ClassCenter and OPOM-ConvexHull, we can conclude that, the approximation method for describing the feature subspace has influence on the person-specific privacy protection task. FI-UAP only uses the single training point, which cannot make use of other images of this identity; OPOM-AffineHull increases the source space, while it may lead to too loose description; OPOM-ClassCenter (similar to FI-UAP+) indeed relies evenly on all the points of this identity, while it may neglect the differences of the training points  $f(X_i^k + \Delta X)$ ; OPOM-ConvexHull increases the source space to an appropriate degree, which will rely on support points differently for a specific training point  $f(X_i^k + \Delta X)$ , and therefore can effectively adapt to different face images better.

### 4.2.2 Combination with Transferability Methods

As privacy protection masks should protect face images from different kinds of unknown models, it is vital to generate transferable masks. We investigate the performance of the proposed OPOM combined with transferability enhancement methods including the momentum boosting method [32] and DFANet [34], as mentioned in Section 3.2.2. Experimental results on the Privacy-Commons dataset and Privacy-Celebrities dataset are listed in Table 4 and Table 5 respectively. With the enhancement of the momentum boosting method and DFANet, the privacy protection success rate of OPOM can increase further, while there still exists some scope for improvement.

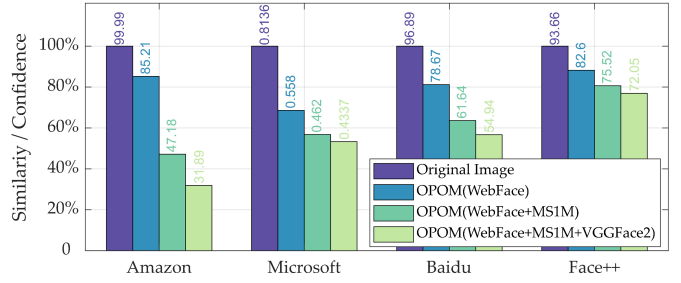


Fig. 4. Protection against Commercial APIs (Amazon [62], Microsoft [63], Baidu [64] and Face++ [65]). Fifty identities in the Privacy-Commons dataset, each with 5 test images are used for the face verification test. The normalized average similarity/confidence scores are shown (lower is better). The original scores are listed above the bar.

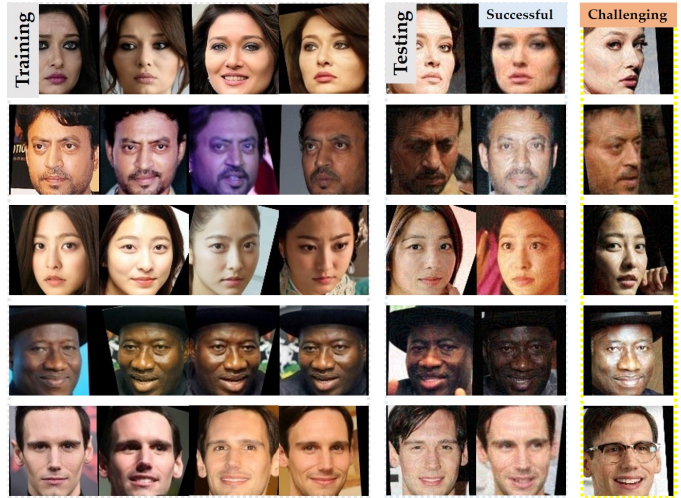


Fig. 5. Some failure cases, as well as the corresponding training samples and successful easily protected samples for analysis. Each row represents an identity. The privacy masks generated with OPOM can generalize to different testing images to some degree. However, if there are obvious differences between the testing images and the training samples, such as, large poses, different illuminations, and occlusions, the mask protection tends to break down.

### 4.2.3 Protection against APIs

In Figure 4, we conduct protection experiments against Commercial APIs (Amazon [62], Microsoft [63], Baidu [64] and Face++ [65]). Fifty identities in the Privacy-Commons dataset are randomly chosen for the face verification test, each with 5 test images. Since commercial APIs are based on extremely large datasets, a single model trained on CASIA-WebFace cannot obtain a better performance. With privacy masks generated with appropriate source models trained on CASIA-WebFace, VGGFace2 [29], and MS-Celeb-1M [38], the average similarity score of the same person decreases significantly, which indicates better privacy protection.

## 4.3 Discussion

### 4.3.1 Failure Case Analysis

Considering the existing performance of the person-specific universal privacy masks, it is pertinent to discuss the current challenges for privacy protection. We select some representative failure cases for analysis, as well as the corresponding training samples and successful easily protected samples, as

TABLE 4

Comparison of different methods combined with the momentum boosting method [32] and DFANet [34] to generate more transferable person-specific privacy masks ( $\varepsilon = 8$ ) from a single source model to protect face images against black-box models. We report Top-1 and Top-5 protection success rate (%) under 1:N identification setting of the Privacy-Commons dataset. The increment compared with TABLE 2 is indicated by symbol  $\uparrow$ .

Source	Method	Target											
		ArcFace		CosFace		SFace		MobileNet		SENet		Inception-ResNet	
		Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Softmax	FI-UAP	80.2 (7.9 $\uparrow$ )	71.2 (8.9 $\uparrow$ )	72.0 (8.6 $\uparrow$ )	63.9 (10.5 $\uparrow$ )	77.9 (7.6 $\uparrow$ )	70.9 (9.1 $\uparrow$ )	82.7 (8.8 $\uparrow$ )	72.9 (11.2 $\uparrow$ )	84.5 (7.1 $\uparrow$ )	77.2 (9.6 $\uparrow$ )	61.5 (9.1 $\uparrow$ )	50.9 (10.3 $\uparrow$ )
	OPOM-AffineHull	81.0 (7.9 $\uparrow$ )	72.7 (9.6 $\uparrow$ )	73.8 (10.6 $\uparrow$ )	65.3 (11.4 $\uparrow$ )	79.1 (8.0 $\uparrow$ )	72.5 (10.5 $\uparrow$ )	83.0 (8.2 $\uparrow$ )	73.7 (10.6 $\uparrow$ )	85.0 (7.2 $\uparrow$ )	77.9 (9.3 $\uparrow$ )	61.7 (8.9 $\uparrow$ )	51.6 (10.8 $\uparrow$ )
	OPOM-ClassCenter	80.6 (3.7 $\uparrow$ )	72.2 (4.4 $\uparrow$ )	73.9 (4.7 $\uparrow$ )	65.6 (5.3 $\uparrow$ )	78.8 (3.9 $\uparrow$ )	72.1 (4.8 $\uparrow$ )	83.1 (4.6 $\uparrow$ )	74.0 (6.1 $\uparrow$ )	85.0 (3.0 $\uparrow$ )	78.2 (4.9 $\uparrow$ )	62.5 (5.4 $\uparrow$ )	51.9 (6.4 $\uparrow$ )
	OPOM-ConvexHull	81.2 (3.3 $\uparrow$ )	72.6 (3.3 $\uparrow$ )	73.6 (3.4 $\uparrow$ )	65.7 (4.3 $\uparrow$ )	79.1 (3.0 $\uparrow$ )	72.2 (3.5 $\uparrow$ )	83.6 (4.4 $\uparrow$ )	74.6 (5.5 $\uparrow$ )	85.0 (2.1 $\uparrow$ )	78.0 (3.8 $\uparrow$ )	62.4 (3.7 $\uparrow$ )	52.1 (4.8 $\uparrow$ )
ArcFace	FI-UAP	89.8 (7.5 $\uparrow$ )	84.9 (9.9 $\uparrow$ )	81.1 (9.9 $\uparrow$ )	75.3 (11.7 $\uparrow$ )	86.2 (8.9 $\uparrow$ )	81.0 (11.0 $\uparrow$ )	79.3 (14.1 $\uparrow$ )	68.5 (18.2 $\uparrow$ )	85.6 (11.7 $\uparrow$ )	78.8 (15.7 $\uparrow$ )	70.6 (14.0 $\uparrow$ )	60.9 (15.5 $\uparrow$ )
	OPOM-AffineHull	91.0 (8.1 $\uparrow$ )	86.1 (10.3 $\uparrow$ )	82.9 (10.6 $\uparrow$ )	77.3 (12.6 $\uparrow$ )	87.5 (9.4 $\uparrow$ )	82.7 (11.8 $\uparrow$ )	80.9 (14.6 $\uparrow$ )	70.9 (18.7 $\uparrow$ )	86.4 (11.0 $\uparrow$ )	80.5 (16.1 $\uparrow$ )	72.1 (13.7 $\uparrow$ )	63.4 (16.6 $\uparrow$ )
	OPOM-ClassCenter	90.7 (4.8 $\uparrow$ )	86.1 (6.6 $\uparrow$ )	82.9 (6.8 $\uparrow$ )	77.1 (7.5 $\uparrow$ )	87.6 (5.6 $\uparrow$ )	82.5 (6.7 $\uparrow$ )	80.6 (10.9 $\uparrow$ )	70.5 (13.8 $\uparrow$ )	86.5 (8.0 $\uparrow$ )	80.1 (10.4 $\uparrow$ )	72.1 (10.3 $\uparrow$ )	62.4 (11.4 $\uparrow$ )
	OPOM-ConvexHull	90.7 (4.3 $\uparrow$ )	86.2 (6.0 $\uparrow$ )	82.7 (5.9 $\uparrow$ )	77.0 (7.1 $\uparrow$ )	87.6 (4.9 $\uparrow$ )	82.7 (6.2 $\uparrow$ )	81.2 (10.7 $\uparrow$ )	71.3 (14.0 $\uparrow$ )	86.5 (7.2 $\uparrow$ )	80.6 (10.4 $\uparrow$ )	72.0 (8.8 $\uparrow$ )	63.0 (10.8 $\uparrow$ )
CosFace	FI-UAP	88.6 (8.2 $\uparrow$ )	82.8 (10.5 $\uparrow$ )	82.5 (9.7 $\uparrow$ )	76.4 (11.2 $\uparrow$ )	85.5 (8.9 $\uparrow$ )	79.7 (10.8 $\uparrow$ )	69.7 (14.1 $\uparrow$ )	56.0 (15.0 $\uparrow$ )	79.2 (12.3 $\uparrow$ )	69.6 (15.1 $\uparrow$ )	62.0 (12.0 $\uparrow$ )	51.3 (14.2 $\uparrow$ )
	OPOM-AffineHull	89.4 (7.1 $\uparrow$ )	84.3 (9.8 $\uparrow$ )	84.0 (9.7 $\uparrow$ )	78.8 (11.8 $\uparrow$ )	86.9 (8.4 $\uparrow$ )	81.8 (10.5 $\uparrow$ )	72.9 (14.3 $\uparrow$ )	60.7 (16.9 $\uparrow$ )	81.8 (12.8 $\uparrow$ )	73.4 (16.6 $\uparrow$ )	64.5 (12.7 $\uparrow$ )	53.8 (14.4 $\uparrow$ )
	OPOM-ClassCenter	89.4 (4.0 $\uparrow$ )	84.1 (5.7 $\uparrow$ )	83.8 (5.0 $\uparrow$ )	78.4 (6.8 $\uparrow$ )	86.8 (4.4 $\uparrow$ )	81.5 (6.0 $\uparrow$ )	72.4 (11.0 $\uparrow$ )	60.4 (12.9 $\uparrow$ )	81.1 (8.1 $\uparrow$ )	72.1 (10.6 $\uparrow$ )	64.4 (9.2 $\uparrow$ )	53.7 (10.8 $\uparrow$ )
	OPOM-ConvexHull	89.6 (3.0 $\uparrow$ )	84.1 (4.7 $\uparrow$ )	84.2 (4.7 $\uparrow$ )	78.8 (6.1 $\uparrow$ )	87.8 (4.8 $\uparrow$ )	82.3 (6.0 $\uparrow$ )	73.0 (10.7 $\uparrow$ )	60.8 (12.6 $\uparrow$ )	82.0 (8.1 $\uparrow$ )	73.4 (10.2 $\uparrow$ )	64.8 (8.0 $\uparrow$ )	53.6 (9.1 $\uparrow$ )

TABLE 5

Comparison of different methods combined with the momentum boosting method [32] and DFANet [34] to generate more transferable person-specific privacy masks ( $\varepsilon = 8$ ) from a single source model to protect face images against black-box models. We report Top-1 and Top-5 protection success rate (%) under 1:N identification setting of the Privacy-Celebrities dataset. The increment compared with TABLE 3 is indicated by symbol  $\uparrow$ .

Source	Method	Target											
		ArcFace		CosFace		SFace		MobileNet		SENet		Inception-ResNet	
		Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Softmax	FI-UAP	62.7 (6.2 $\uparrow$ )	52.2 (6.6 $\uparrow$ )	52.9 (7.7 $\uparrow$ )	43.2 (8.3 $\uparrow$ )	60.3 (7.4 $\uparrow$ )	50.8 (8.3 $\uparrow$ )	67.8 (7.6 $\uparrow$ )	56.2 (9.3 $\uparrow$ )	68.6 (7.0 $\uparrow$ )	58.6 (8.8 $\uparrow$ )	43.6 (6.2 $\uparrow$ )	32.1 (5.6 $\uparrow$ )
	OPOM-AffineHull	65.1 (6.3 $\uparrow$ )	54.8 (7.0 $\uparrow$ )	54.9 (8.0 $\uparrow$ )	45.1 (7.9 $\uparrow$ )	62.4 (7.6 $\uparrow$ )	53.2 (8.3 $\uparrow$ )	69.9 (7.4 $\uparrow$ )	59.0 (9.6 $\uparrow$ )	70.5 (7.0 $\uparrow$ )	60.8 (8.6 $\uparrow$ )	45.0 (6.0 $\uparrow$ )	33.4 (5.5 $\uparrow$ )
	OPOM-ClassCenter	64.3 (1.9 $\uparrow$ )	54.0 (2.5 $\uparrow$ )	54.2 (2.36 $\uparrow$ )	45.4 (3.0 $\uparrow$ )	61.7 (1.9 $\uparrow$ )	52.5 (2.8 $\uparrow$ )	69.7 (3.5 $\uparrow$ )	58.7 (3.9 $\uparrow$ )	69.8 (2.1 $\uparrow$ )	60.1 (2.4 $\uparrow$ )	45.2 (2.8 $\uparrow$ )	33.9 (2.5 $\uparrow$ )
	OPOM-ConvexHull	65.7 (1.9 $\uparrow$ )	55.5 (2.0 $\uparrow$ )	55.5 (1.9 $\uparrow$ )	46.0 (2.1 $\uparrow$ )	63.0 (1.7 $\uparrow$ )	53.7 (2.1 $\uparrow$ )	70.1 (3.1 $\uparrow$ )	59.1 (3.4 $\uparrow$ )	70.9 (2.0 $\uparrow$ )	61.3 (2.4 $\uparrow$ )	46.5 (2.3 $\uparrow$ )	35.3 (2.1 $\uparrow$ )
ArcFace	FI-UAP	73.0 (10.5 $\uparrow$ )	64.5 (13.1 $\uparrow$ )	60.3 (12.4 $\uparrow$ )	50.9 (13.1 $\uparrow$ )	67.4 (11.8 $\uparrow$ )	58.9 (13.0 $\uparrow$ )	63.9 (13.9 $\uparrow$ )	50.4 (14.8 $\uparrow$ )	68.0 (13.7 $\uparrow$ )	56.8 (15.8 $\uparrow$ )	49.2 (11.5 $\uparrow$ )	37.6 (11.3 $\uparrow$ )
	OPOM-AffineHull	75.9 (10.1 $\uparrow$ )	67.7 (12.5 $\uparrow$ )	64.0 (12.7 $\uparrow$ )	55.4 (14.1 $\uparrow$ )	71.0 (11.9 $\uparrow$ )	63.4 (13.8 $\uparrow$ )	67.0 (13.9 $\uparrow$ )	54.8 (15.7 $\uparrow$ )	71.5 (13.8 $\uparrow$ )	61.1 (16.0 $\uparrow$ )	53.1 (12.6 $\uparrow$ )	41.6 (13.2 $\uparrow$ )
	OPOM-ClassCenter	75.3 (6.5 $\uparrow$ )	67.2 (7.9 $\uparrow$ )	63.4 (8.2 $\uparrow$ )	54.9 (8.8 $\uparrow$ )	70.6 (7.6 $\uparrow$ )	62.8 (8.5 $\uparrow$ )	66.0 (9.4 $\uparrow$ )	53.5 (10.9 $\uparrow$ )	70.8 (9.7 $\uparrow$ )	60.4 (11.4 $\uparrow$ )	52.4 (8.5 $\uparrow$ )	40.9 (8.9 $\uparrow$ )
	OPOM-ConvexHull	76.3 (6.5 $\uparrow$ )	68.0 (7.6 $\uparrow$ )	64.4 (8.2 $\uparrow$ )	55.6 (8.3 $\uparrow$ )	71.2 (7.1 $\uparrow$ )	63.4 (8.2 $\uparrow$ )	67.0 (9.5 $\uparrow$ )	54.5 (11.4 $\uparrow$ )	72.0 (9.7 $\uparrow$ )	61.4 (11.4 $\uparrow$ )	53.9 (9.0 $\uparrow$ )	42.3 (9.1 $\uparrow$ )
CosFace	FI-UAP	71.7 (8.0 $\uparrow$ )	62.1 (10.1 $\uparrow$ )	60.7 (9.1 $\uparrow$ )	51.1 (9.7 $\uparrow$ )	66.5 (8.8 $\uparrow$ )	57.4 (10.2 $\uparrow$ )	56.0 (9.5 $\uparrow$ )	41.7 (10.0 $\uparrow$ )	61.4 (10.6 $\uparrow$ )	48.9 (11.9 $\uparrow$ )	42.4 (9.4 $\uparrow$ )	30.3 (8.5 $\uparrow$ )
	OPOM-AffineHull	74.9 (8.4 $\uparrow$ )	66.1 (10.5 $\uparrow$ )	65.1 (9.6 $\uparrow$ )	56.8 (11.4 $\uparrow$ )	70.6 (9.5 $\uparrow$ )	62.7 (11.2 $\uparrow$ )	60.7 (11.2 $\uparrow$ )	46.5 (11.8 $\uparrow$ )	66.2 (11.5 $\uparrow$ )	54.9 (14.1 $\uparrow$ )	46.2 (10.5 $\uparrow$ )	34.4 (10.5 $\uparrow$ )
	OPOM-ClassCenter	73.8 (5.5 $\uparrow$ )	64.8 (6.4 $\uparrow$ )	64.3 (5.8 $\uparrow$ )	55.9 (7.1 $\uparrow$ )	69.8 (6.0 $\uparrow$ )	61.6 (7.1 $\uparrow$ )	59.9 (8.7 $\uparrow$ )	45.3 (9.0 $\uparrow$ )	65.1 (8.1 $\uparrow$ )	53.4 (9.5 $\uparrow$ )	45.7 (7.2 $\uparrow$ )	33.6 (7.1 $\uparrow$ )
	OPOM-ConvexHull	75.2 (5.3 $\uparrow$ )	66.6 (6.0 $\uparrow$ )	65.1 (5.3 $\uparrow$ )	57.0 (6.5 $\uparrow$ )	70.9 (5.3 $\uparrow$ )	63.0 (6.8 $\uparrow$ )	60.6 (7.6 $\uparrow$ )	46.4 (8.3 $\uparrow$ )	66.2 (7.5 $\uparrow$ )	54.9 (9.4 $\uparrow$ )	47.0 (6.8 $\uparrow$ )	34.8 (6.7 $\uparrow$ )

shown in figure 5, where each row represents an identity. We can see that, the generated privacy masks can generalize to different testing images to some degree. However, if there are obvious differences between the testing images and the training samples, such as large poses, different illuminations and occlusions, the mask protection tends to break down.

#### 4.3.2 Why person-specific? Effectiveness and Efficiency

One might naturally wonder what are the strengths and weaknesses of person-specific (class-wise) universal masks compared with image-specific and universal masks. Here, we give an analysis in terms of effectiveness and efficiency. Specifically, universal masks (GD-UAP, GAP [48] and FI-UAP-all), person-specific universal masks (FI-UAP, OPOM-AffineHull and OPOM-ConvexHull), and image-specific masks (FIM [34], LowKey [27] and M-DI<sup>2</sup>-APF [26], [32], [33]) are compared in Figure 6. For brevity, we use the average protection success rate of six black-box models to represent the effectiveness, and use the average generation time of 100 images based on the ResNet-50 model to represent the efficiency.

The experimental results show that the protection success rate of image-specific masks (FIM, LowKey and M-DI<sup>2</sup>-APF) is the best, and the effectiveness of person-specific universal masks (FI-UAP, OPOM-AffineHull and OPOM-ConvexHull) is better than that of universal masks (GD-UAP, GAP and FI-UAP-all). While considering the efficiency, it takes 1.65s, 2.35s or 3.50s to generate a mask of a new image for image-specific privacy protection, even with a GPU (6.69s, 10.88s, and 41.20s with CPU). In contrast, person-specific universal masks and universal masks can dispense with the mask generation time for new images, that is, the generation time of new images is 0s. Compared with universal masks and image-specific masks, person-specific (class-wise) universal masks show a tradeoff between effectiveness and efficiency, which may benefit average users and some real-time video applications.

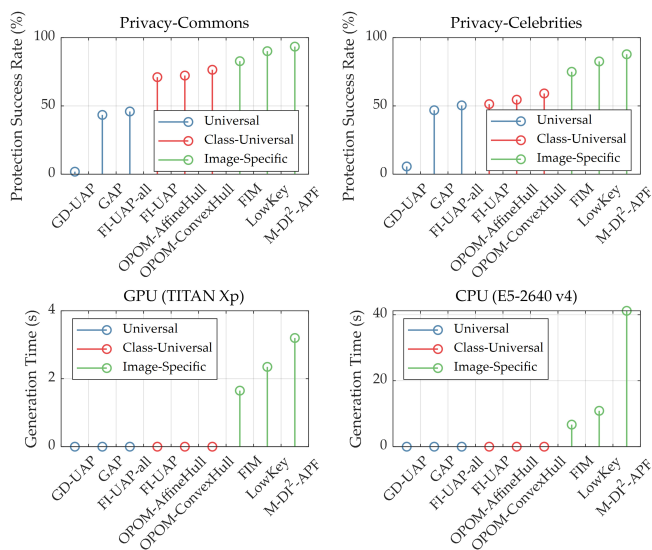


Fig. 6. Comparison of universal, person-specific (class-wise) and image-specific masks in terms of effectiveness and efficiency.

#### 4.3.3 Privacy Protection in Videos

The proposed OPOM method can be applied to the privacy protection of photographs to be posted on the social media platform. The still image, of course, is not the only application. We demonstrate that OPOM is also applicable to the privacy protection in videos.

Specifically, we use the series Sherlock database [66] for experiments. Following the original work [66], we also collected images of actors that were the main characters by searching the names Benedict Cumberbatch and Martin Freeman. We manually found 10 high quality web-collected images of the two characters and use the mean value of their deep features to represent their feature subspace. Then, some other web-collected images, which can differ from the target images (from the TV material) in hairstyle, makeup, lighting, and viewpoint, are used as the training samples. Finally, we apply the generated privacy masks to target images of the videos. Some frames are shown in Figure 7, where the first row shows the original frames of the video, and the second row shows the modified frames with the privacy masks of Sherlock generated from other face images of actors (Benedict Cumberbatch).

The effectiveness of the privacy protection is evaluated using the cosine similarity between the deep features of images from videos and the feature prototype. As shown in Figure 8, the average value and the standard deviation of cosine similarity between the face features of characters (Sherlock Holmes and Doctor John Watson), tested on six black-box models, are used to demonstrate the effectiveness of OPOM. With no privacy masks (no training images) provided, it is easy to recognize the character. With privacy masks, privacy can be protected against black-box deep face recognition models. We gradually increase the number of web images for training and found that 15 training images can obtain good protection performance.

#### 4.3.4 Diversity of Privacy Masks

In this paper, we apply privacy masks to the probe images, not the gallery images, since we assume that some unauthorized face recognition services may have already obtained the original images of regular users. Despite the effectiveness of person-specific (class-wise) universal masks in one-shot usage, we have to consider a more difficult situation since new images on the internet are being continuously collected. If masked face images have been labeled and recorded in the gallery set, then new images with the same mask (probe) may no longer be safe. Therefore, we explore whether a diversity of person-specific masks can be generated for an identity to advance the possibility of potential solutions for this problem.

That is, the objective is to generate a set of diverse person-specific (class-wise) universal masks to fool a variety of deep face recognition models  $f(\cdot)$  for all the face images  $X^k = \{X_1^k, X_2^k, \dots, X_i^k, \dots\}$  of identity  $k$ , so that any face image  $X_i^k$  can conceal the identity with any mask of  $\{\Delta X_1, \Delta X_2, \dots, \Delta X_j, \dots\}$ :

$$\begin{aligned}
 D(f(X_i^k + \Delta X_{j_1}), f_{X^k}) &> t, \\
 D(f(X_i^k + \Delta X_{j_2}), f_{X^k}) &> t, \\
 D(f(X_i^k + \Delta X_{j_1}), f(X_i^k + \Delta X_{j_2})) &> t, \\
 \|\Delta X_{j_1}\|_\infty &< \varepsilon, \|\Delta X_{j_2}\|_\infty &< \varepsilon.
 \end{aligned} \tag{19}$$

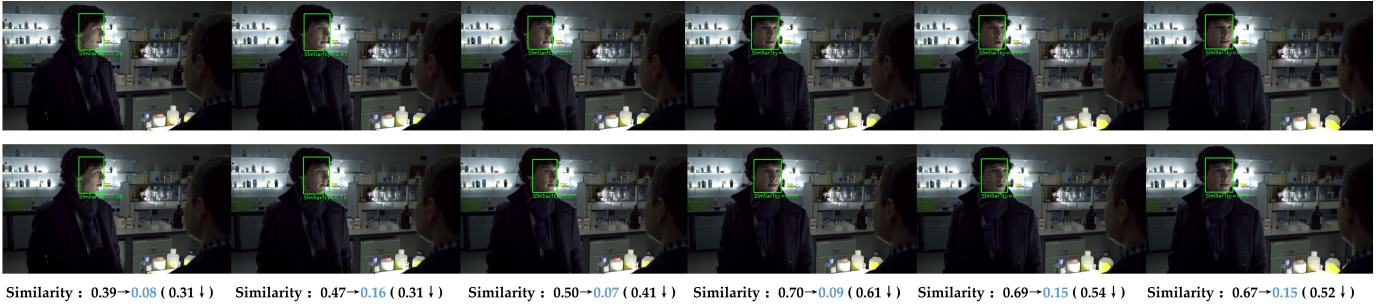


Fig. 7. Application of OPOM in video privacy protection, Sherlock [66]. The first row shows the original frames. The second row shows the modified frames with the privacy masks of Sherlock generated from other images of actors (Benedict Cumberbatch). The cosine similarity between the deep features of the detected face in the video and the deep features of the corresponding character (Sherlock Holmes) is used to show the effectiveness.

TABLE 6

Diverse person-specific privacy masks ( $\varepsilon = 8$ ) from a single source model to protect face images against black-box models. Here, “M-O” denotes that the masked image is used as the probe, while the original image is recorded in the gallery set. “M-M” represents masked image has been collected and applied in the gallery set, and new masked image is used as the probe. (The same mask is used for  $n_M=1$ .) We report the Top-1 protection success rate (%) under 1:N identification setting of Privacy-Commons dataset. The higher protection success rate is better.

Source	$n_M$	Target											
		ArcFace		CosFace		SFace		MobileNet		SENet		Inception-ResNet	
		M-O	M-M	M-O	M-M	M-O	M-M	M-O	M-M	M-O	M-M	M-O	M-M
Softmax	1	78.0	20.1	70.2	9.3	76.1	11.0	79.2	18.7	82.9	20.8	58.7	14.0
	5	65.8±1.9	82.8±4.2	56.2±3.4	72.3±6.5	63.0±3.1	79.4±5.9	66.6±2.4	77.7±4.0	71.0±2.3	85.9±3.3	43.6±2.8	66.6±4.1
	10	54.7±8.3	81.4±3.8	43.8±9.5	71.5±5.5	50.4±9.7	77.7±5.0	55.4±8.7	76.6±3.7	59.3±9.0	83.8±3.3	33.2±7.5	64.3±4.2
ArcFace	1	86.5	24.5	76.8	11.5	82.7	14.2	70.5	18.0	79.3	20.7	63.2	16.3
	5	78.2±2.3	89.9±3.2	65.0±3.1	76.6±6.3	71.6±2.9	84.3±4.6	58.4±3.1	70.3±4.0	68.4±3.4	84.0±3.2	49.4±3.4	71.1±3.9
	10	69.0±7.6	88.3±2.4	53.5±9.3	76.2±4.5	60.6±9.2	82.8±3.5	47.4±8.3	69.6±3.2	57.4±9.2	82.3±2.8	38.4±8.1	68.8±3.5
CosFace	1	86.6	22.3	79.5	11.9	83.0	13.3	62.3	15.7	74.0	18.7	56.8	13.4
	5	76.3±2.9	90.2±3.0	66.9±3.5	82.1±5.3	71.0±3.6	85.7±4.8	50.2±3.6	63.6±3.9	62.0±4.0	79.6±3.5	42.3±3.9	64.6±4.0
	10	66.3±8.2	88.4±2.5	54.6±10.0	80.3±3.9	58.9±9.8	83.8±3.7	39.9±8.0	62.7±3.5	49.9±9.2	77.8±3.3	31.7±7.8	61.9±3.7

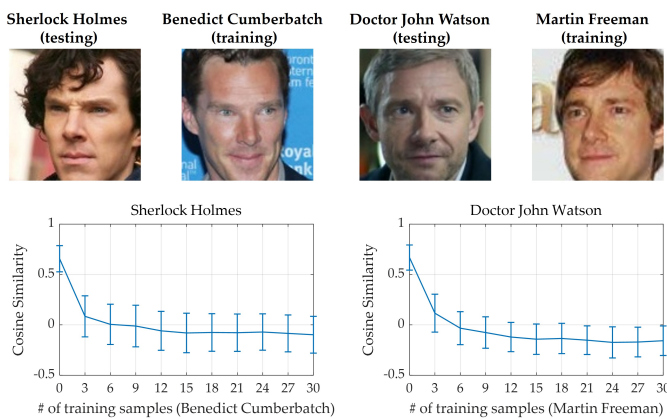


Fig. 8. Application of OPOM in video privacy protection, Sherlock [66]. The privacy masks are trained with other face images of actors (Benedict Cumberbatch and Martin Freeman). The average cosine similarity between the deep features of the detected face in the video and the deep features of the corresponding characters (Sherlock Holmes and Doctor John Watson) is used to demonstrate the effectiveness.

In the above,  $f(X_i^k) \in \mathbb{R}^d$  is the normalized feature of image  $X_i^k$ .  $f_{X^k}$  denotes the feature subspace of identity  $k$ .  $\Delta X_{j_1}$  and  $\Delta X_{j_2}$  are any two generated privacy masks of the set.  $t$  is the distance threshold to decide whether a pair of face images belongs to the same identity.  $D(x_1, x_2)$  denotes the distance between  $x_1$  and  $x_2$ , that is, the shortest distance between point  $x_1$  and subspace  $x_2$ . When  $x_2$  denotes a point, we use

$D(x_1, x_2)$  to represent the normalized Euclidean distance or cosine distance commonly used in face recognition.  $\varepsilon$  limits the maximum deviation of the privacy mask.

The objective not only optimizes each masked sample in the direction away from the feature subspace of the source identity, but also optimizes samples with different masks away from each other.

Considering the above function, we still use the proposed approximation methods for  $f_{X^k}$  with a limited number  $n_k$  of face images  $\tilde{X}^k = \{X_1^k, X_2^k, \dots, X_{n_k}^k\}$  and generate  $n_M$  masks  $\{\Delta X_1, \Delta X_2, \dots, \Delta X_{n_M}\}$  with previous source models. For experimental settings, in addition to the original image as gallery set, we also consider that different masked images are used in the gallery set.

Experimental results in Table 6 indicate that person-specific privacy masks can be diverse and varied. In this way, the average user can apply different masks for privacy protection each time to cope with the situation that some masks have been ineffective. However, it seems that the protection rate will be lessened as the number of privacy masks increases. Note that we are only presenting a preliminary exploration; it will be for future work to study more elaborate solutions for this purpose.

## 5 CONCLUSION AND FUTURE WORKS

In this paper, we have presented a type of class-wise universal adversarial perturbation for a new customized privacy protection task, by generating a person-specific masks that

can be applied universally for all the images of an identity. Experimental results on two test benchmarks have demonstrated the effectiveness and superiority of the proposed method on this customized privacy protection task. We have demonstrated that the proposed method can be used in the privacy protection of videos.

Great progress has been made, yet much still remains to be done. As we stated before, it is challenging to offer protection if there are obvious differences between the testing images and the training samples. Therefore, it is imperative to further improve the image universality to cover the potential diverse variances of face images, such as large poses, different illuminations, and occlusions. As analyzed before, we considered a diverse of masks to address the potential problems of privacy mask leakage. However, the proposed method has challenges in terms of the diversity and effects of the privacy masks, and therefore remains a preliminary exploration. In addition, we are currently at the implementation stage of generating digital privacy cloaks that can only be used for digital images and videos on social media. They will have more impact if they can be applied in practical-world video surveillance.

## ACKNOWLEDGMENT

We would like to thank Jinglin Zhang and Xuannan Liu for useful feedback, and anonymous referees for their valuable comments. This work was partially supported by the National Natural Science Foundation of China under Grants No. 61871052 and 62192784.

## REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representation*, 2015.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [4] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [5] Y. Sun, Y. Chen, X. Wang, and X. Tang, "Deep learning face representation by joint identification-verification," in *Advances in neural information processing systems*, 2014, pp. 1988–1996.
- [6] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [7] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *European conference on computer vision*. Springer, 2016, pp. 499–515.
- [8] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 212–220.
- [9] B. Chen, W. Deng, and J. Du, "Noisy softmax: Improving the generalization ability of dcnn via postponing the early softmax saturation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5372–5381.
- [10] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5265–5274.
- [11] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [12] Y. Huang, Y. Wang, Y. Tai, X. Liu, P. Shen, S. Li, J. Li, and F. Huang, "Curriculumface: adaptive curriculum learning loss for deep face recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5901–5910.
- [13] Y. Zhong, W. Deng, J. Hu, D. Zhao, X. Li, and D. Wen, "Sface: sigmoid-constrained hypersphere loss for robust face recognition," *IEEE Transactions on Image Processing*, vol. 30, pp. 2587–2598, 2021.
- [14] K. W. Bowyer, "Face recognition technology: security versus privacy," *IEEE Technology and society magazine*, vol. 23, no. 1, pp. 9–19, 2004.
- [15] "Facial Recognition Is Here But We Have No Laws," <https://www.nextgov.com/ideas/2020/07/facial-recognition-here-we-have-no-laws/166711/>.
- [16] "Facial Recognition: When Convenience and Privacy Collide," <https://www.securitymagazine.com/articles/90533-facial-recognition-when-convenience-and-privacy-collide>.
- [17] E. M. Newton, L. Sweeney, and B. Malin, "Preserving privacy by de-identifying face images," *IEEE transactions on Knowledge and Data Engineering*, vol. 17, no. 2, pp. 232–243, 2005.
- [18] R. Gross, L. Sweeney, J. Cohn, F. De la Torre, and S. Baker, "Face de-identification," in *Protecting privacy in video surveillance*. Springer, 2009, pp. 129–146.
- [19] B. Meden, R. C. Malli, S. Fabijan, H. K. Ekenel, V. Štruc, and P. Peer, "Face deidentification with generative deep neural networks," *IET Signal Processing*, vol. 11, no. 9, pp. 1046–1054, 2017.
- [20] O. Gafni, L. Wolf, and Y. Taigman, "Live face de-identification in video," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 9378–9387.
- [21] T. Li and L. Lin, "Anonymousnet: Natural face de-identification with measurable privacy," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [22] S. Joon Oh, M. Fritz, and B. Schiele, "Adversarial image perturbation for privacy protection—a game theory perspective," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1482–1491.
- [23] E. Chatzikyriakidis, C. Papaioannidis, and I. Pitas, "Adversarial face de-identification," in *2019 IEEE International conference on image processing (ICIP)*. IEEE, 2019, pp. 684–688.
- [24] S. Shan, E. Wenger, J. Zhang, H. Li, H. Zheng, and B. Y. Zhao, "Fawkes: Protecting privacy against unauthorized deep learning models," in *29th {USENIX} Security Symposium ({USENIX} Security 20)*, 2020, pp. 1589–1604.
- [25] X. Yang, Y. Dong, T. Pang, J. Zhu, and H. Su, "Towards privacy protection by generating adversarial identity masks," *arXiv preprint arXiv:2003.06814*, 2020.
- [26] J. Zhang, J. Sang, X. Zhao, X. Huang, Y. Sun, and Y. Hu, "Adversarial privacy-preserving filter," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 1423–1431.
- [27] V. Cherepanova, M. Goldblum, H. Foley, S. Duan, J. P. Dickerson, G. Taylor, and T. Goldstein, "Lowkey: Leveraging adversarial attacks to protect social media users from facial recognition," in *International Conference on Learning Representations*, 2021.
- [28] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, and I. Goodfellow, "Intriguing properties of neural networks," in *International Conference on Learning Representation*, 2014.
- [29] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vg-gface2: A dataset for recognising faces across pose and age," in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 67–74.
- [30] Y. Liu, X. Chen, C. Liu, and D. Song, "Delving into transferable adversarial examples and black-box attacks," in *International Conference on Learning Representation*, 2017.
- [31] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical black-box attacks against machine learning," in *Proceedings of the 2017 ACM on Asia conference on computer and communications security*. ACM, 2017, pp. 506–519.
- [32] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 9185–9193.
- [33] C. Xie, Z. Zhang, Y. Zhou, S. Bai, J. Wang, Z. Ren, and A. L. Yuille, "Improving transferability of adversarial examples with

- input diversity," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2730–2739.
- [34] Y. Zhong and W. Deng, "Towards transferable adversarial attack against deep face recognition," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 1452–1466, 2020.
- [35] H. Cho and R. LaRose, "Privacy issues in internet surveys," *Social science computer review*, vol. 17, no. 4, pp. 421–434, 1999.
- [36] B. Liu, M. Ding, S. Shaham, W. Rahayu, F. Farokhi, and Z. Lin, "When machine learning meets privacy: A survey and outlook," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–36, 2021.
- [37] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Learning face representation from scratch," *arXiv preprint arXiv:1411.7923*, 2014.
- [38] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, "Ms-celeb-1m: A dataset and benchmark for large-scale face recognition," in *European Conference on Computer Vision*. Springer, 2016, pp. 87–102.
- [39] M. Wang and W. Deng, "Deep face recognition: A survey," *Neuro-computing*, vol. 429, pp. 215–244, 2021.
- [40] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *International Conference on Learning Representation*, 2015.
- [41] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *International Conference on Learning Representation Workshop*, 2017.
- [42] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2574–2582.
- [43] Y. Dong, T. Pang, H. Su, and J. Zhu, "Evading defenses to transferable adversarial examples by translation-invariant attacks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4312–4321.
- [44] D. Wu, Y. Wang, S.-T. Xia, J. Bailey, and X. Ma, "Skip connections matter: On the transferability of adversarial examples generated with resnets," in *International Conference on Learning Representation*, 2020.
- [45] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1765–1773.
- [46] J. Hayes and G. Danezis, "Learning universal adversarial perturbations with generative models," in *2018 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2018, pp. 43–49.
- [47] K. Reddy Mopuri, U. Ojha, U. Garg, and R. Venkatesh Babu, "Nag: Network for adversary generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 742–751.
- [48] O. Poursaeed, I. Katsman, B. Gao, and S. Belongie, "Generative adversarial perturbations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4422–4431.
- [49] K. R. Mopuri, U. Garg, and R. V. Babu, "Fast feature fool: A data independent approach to universal adversarial perturbations," *arXiv preprint arXiv:1707.05572*, 2017.
- [50] K. R. Mopuri, A. Ganeshan, and R. V. Babu, "Generalizable data-free objective for crafting universal adversarial perturbations," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 10, pp. 2452–2465, 2018.
- [51] T. Gupta, A. Sinha, N. Kumari, M. Singh, and B. Krishnamurthy, "A method for computing class-wise universal adversarial perturbations," *arXiv preprint arXiv:1912.00466*, 2019.
- [52] C. Zhang, P. Benz, T. Imtiaz, and I.-S. Kweon, "Cd-uap: Class discriminative universal adversarial perturbation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 6754–6761.
- [53] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.
- [54] S. Diamond and S. Boyd, "CVXPY: A Python-embedded modeling language for convex optimization," *Journal of Machine Learning Research*, vol. 17, no. 83, pp. 1–5, 2016.
- [55] A. Agrawal, R. Verschueren, S. Diamond, and S. Boyd, "A rewriting system for convex optimization problems," *Journal of Control and Decision*, vol. 5, no. 1, pp. 42–60, 2018.
- [56] Y. Zhong and W. Deng, "Adversarial learning with margin-based triplet embedding regularization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6549–6558.
- [57] A. Nech and I. Kemelmacher-Shlizerman, "Level playing field for million scale face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7044–7053.
- [58] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," University of Massachusetts, Amherst, Tech. Rep. 07-49, October 2007.
- [59] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [60] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [61] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [62] "Amazon's Rekognition Tool," <https://aws.amazon.com/rekognition/>.
- [63] "Microsoft Azure," <https://www.azure.cn>.
- [64] "Baidu Cloud Vision Api," <http://ai.baidu.com>.
- [65] "Face++ Research Toolkit," <https://www.faceplusplus.com.cn/>.
- [66] A. Nagrani and A. Zisserman, "From benedict cumberbatch to sherlock holmes: Character identification in tv series without a script," in *The British Machine Vision Conference (BMVC)*, 2017.