# Face Deidentification with Generative Deep Neural Networks

Blaž Meden[1*], Refik Can Mallı[2], Sebastjan Fabijan[3], Hazım Kemal Ekenel[2], Vitomir Štruc[3], Peter Peer[1]

[1]Faculty of Computer and Information Science, University of Ljubljana, Večna pot 113, SI-1001 Ljubljana, Slovenia
[2]Department of Computer Engineering, Istanbul Technical University, 34469 Maslak, Istanbul, Turkey
[3]Faculty of Electrical Engineering, University of Ljubljana, Tržaška cesta 25, SI-1000 Ljubljana, Slovenia
[*]blaz.meden@fri.uni-lj.si

**Abstract:** Face deidentification is an active topic amongst privacy and security researchers. Early deidentification methods relying on image blurring or pixelization were replaced in recent years with techniques based on formal anonymity models that provide privacy guaranties and at the same time aim at retaining certain characteristics of the data even after deidentification. The latter aspect is particularly important, as it allows to exploit the deidentified data in applications for which identity information is irrelevant. In this work we present a novel face deidentification pipeline, which ensures anonymity by synthesizing artificial surrogate faces using generative neural networks (GNNs). The generated faces are used to deidentify subjects in images or video, while preserving non-identity-related aspects of the data and consequently enabling data utilization. Since generative networks are very adaptive and can utilize a diverse set of parameters (pertaining to the appearance of the generated output in terms of facial expressions, gender, race, etc.), they represent a natural choice for the problem of face deidentification. To demonstrate the feasibility of our approach, we perform experiments using automated recognition tools and human annotators. Our results show that the recognition performance on deidentified images is close to chance, suggesting that the deidentification process based on GNNs is highly effective.

## 1. Introduction

Over the last decades an extensive amount of video data is being recorded and stored. Since access to and exploitation of such data is difficult to monitor let alone prevent, appropriate measures need to be taken to ensure that such data is not misused and the privacy of people is adequately protected.

A popular approach towards privacy protection in image and video data is the use of deidentification. Ribarić et al. [1] define deidentification as the process of concealing or removing personal identifiers from source content in order to prevent disclosure and use of data for unauthorized purposes. For video data, for example, this may translate to "blurring" or "pixelation" of the facial areas [2], both of which represent early deidentification examples. These naive methods are typi-

cally useful for preventing humans from recognizing subjects in videos, but are far less successful with automated recognition techniques, where repeating the (naive) deidentification process on the test data still enables automated recognition, i.e., parrot attack [3]. Thus, for successful deidentification of images and videos, more advanced techniques are needed.

Another shortcoming of naive deidentification techniques is the fact that all information contained in the data is typically removed even if the information is not related to identity. This raises the question of data utility. If the deidentified data is to be useful for purposes that do not require identity information, but, for example, rely on gender or age information (e.g., customer-profiling applications in shopping malls), this information needs to be preserved even after deidentification. Recent deidentification approaches, therefore, focus on ways of removing identity information from images and videos, while still retaining other non-identity related information [4], [5].

In this paper we follow these recent trends and present a new deidentification approach exploiting generative neural networks (GNNs), which represent contemporary generative models capable of synthesizing photo-realistic artificial images of any object (see, e.g., [6], [7], [8]) based on supplied high-level information. Similarly to existing deidentification techniques, we replace the original faces in the input data with surrogates generated from a small number of identities. However, instead of synthesizing the surrogate faces through pixel averaging as in prior work, we use a GNN to combine identities and generate artificial surrogates for deidentification. The flexibility of the GNN also allows us to parameterize the generation process with respect to various appearance-related characteristics and synthesize faces under different appearances (under varying pose, with different facial expressions, etc.). This property ensures that our deidentification approach is able to conceal the identity of individuals, but also to preserve the utility of the data.

We demonstrate the feasibility of the proposed deidentification pipeline through extensive experiments on the ChokePoint dataset [9]. Our experimental results show that GNNs are a viable solution for the problem of face deidentification and are able to generate realistic, visually convincing deidentification results. Furthermore, the deidentified faces offer a suitable level of privacy protection as evidenced by experiments with a number of contemporary recognition models as well as humans. In summary, we make the following contributions:

- We introduce a face deidentification pipeline that exploits GNNs to produce artificial surrogate faces for deidentification and offers a level of flexibility in the generation process that is not available with existing deidentification approaches.

- We present a qualitative evaluation of the proposed pipeline with challenging data captured in a real surveillance scenario and discuss the advantages and limitations of our deidentification approach.

- We demonstrate the efficacy of the proposed pipeline in comprehensive quantitative experiments with several state-of-the-art recognition techniques from the literature and human annotators.

## 2. Related work

In this section we review the most important work related to our deidentification pipeline. For a more comprehensive review please refer to the surveys by Ribarić et al. [1], [10].

Existing approaches to deidentification often implement formal privacy protection models such as $k$-anonymity [11], $l$-diversity [12], or $t$-closeness [13]. Among these, the $k$-anonymity models

have likely received the most attention in the area of face deidentification and resulted in the so-called $k$-same family of algorithms [3], [4], [14]. These algorithms operate on a closed set of static facial images and substitute each image in the set with the average of the closest $k$ identities computed from the same closed set of images. Because several images are replaced with the same average face, data anonymity of a certain level is guaranteed. A number of $k$-same variants was presented in the literature, including the original $k$-same algorithm [3], $k$-same-select [4], and $k$-same-model [15] to name a few. The majority of these techniques is implemented using Active Appearance Models (AAMs).

Another example of a deidentification technique using AAMs was recently presented by Jourabloo et al. in [16]. Here, the authors combine facial-attribute and face-verification classifiers in a joint objective function. By optimizing the objective function, optimal weights are estimated such that the deidentified and the original image have as many common attributes as possible, but at the same time are classified as two different subjects.

A $q$-far deidentification approach, which is also AAM based, but does not follow the $k$-same principle (since the surrogate faces are mutually different), was proposed by Samaržija and Ribarić in [17] and combines face deidentification with pose estimation. The authors cover different facial orientations by fitting multiple AAMs and achieve anonymity by replacing the original faces with surrogates that are sufficiently far (i.e., $q$-far) from the initial identities.

Sim and Zhang present a method for controllable face deidentification in [5]. They demonstrate a high degree of control over different attributes (such as identity, gender, age, femininity or race) of the deidentified faces and similar to our approach are able to alter or retain specific aspects of the target appearance.

Another example related to the AAM-based techniques was recently proposed by Sun et al. [18]. Here, the authors propose the $k$-diff-furthest algorithm, which is related to the $k$-anonymity model and $k$-same family of techniques, but differs in its ability to track individuals in the deidentified video, since the deidentified faces have distinguishable properties. The experimental evaluation performed by the authors shows that the algorithm is capable of maintaining the diversity of the deidentified faces and keeps them as distinguishable as their original faces. However, the approach does not deal with the data utility aspect, e.g., expression preservation.
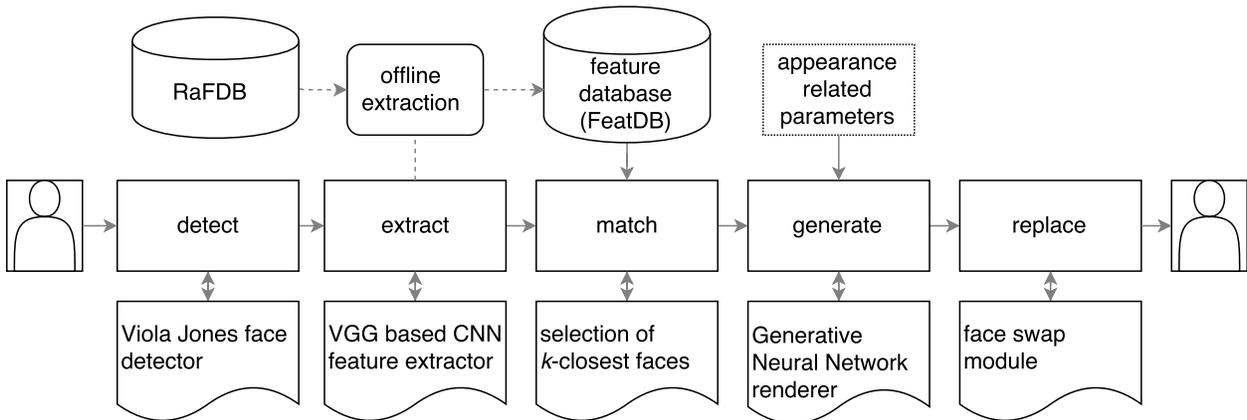
Different from AAM-based deidentification approaches, Brkić et al. [19] propose a deidentification method based on style-transfer. The authors describe a pipeline that enables altering the appearance of faces in videos in such a way that an artistic style replacement is performed on the input data, thus making automatic recognition more difficult. Another interesting deidentification approach was presented by Chriskos et al. in [20], which, in contrary to most deidentification methods, hinders the recognition only for automatic recognition algorithms, but not human observers. The authors utilize projections on hyperspheres in order to defeat classifiers, while preserving enough visual information to enable human viewers to correctly identify individuals.

## 3. Deidentification with generative deep neural networks

Here we present a detailed description of our deidentification pipeline, which exploits generative neural networks (GNNs) to conceal the identity of people in the image data.

## 3.1. Overview

A block-diagram of our deidentification pipeline is presented in Fig. 1. The procedure starts with a face detection step that takes an image or video frame as input and then locates candidate facial regions in the input data for deidentification. For each detected region, a feature vector is computed using a state-of-the-art deep face recognition network, i.e., the VGG network from [21], and matched against a fixed gallery of $M$ subjects. Based on this matching procedure, the $k$-closest identities ($k \ll M$) from the gallery data are selected and fed to our generative network to synthesize an artificial face with visual characteristics of the selected $k$ identities. Finally, the artificially generated (deidentified) face is blended into the input image (or frame) to conceal the original identity.



**Fig. 1.** *Block diagram of our deidentification pipeline. The procedure uses a generative neural network to generate synthetic faces that can be used for deidentification. Each generated face is a combination of $k$ identities from the gallery data that are closest (i.e., most similar in the feature space) to the input face.*

One appealing characteristic of our deidentification pipeline is the flexibility of the GNN, which is able to synthesize high-resolution, realistic-looking faces under various appearances. Here, the generation process is governed by a small number of appearance-related parameters that control the visual characteristics of the synthesized faces, such as pose, skin color, gender, identity, facial expression, and alike. Thus, with this setup, we are able to generate artificial faces with predefined identities, facial expressions, gender, and so forth, or alter any of these at the time. For instance, if our goal is to preserve facial expressions of faces, we could automatically recognize facial expressions from the input image and use the recognition result as input to the GNN. The network would then generate a synthetic image with the predefined expression. A similar procedure could be used to retain or alter any visual characteristic of the input faces and contribute towards the preservation of data utility, which is one of the main goals of contemporary deidentification technology.

Even though our approach is similar in nature to the $k$-same family of algorithms [3], [4], [10] that implement the $k$-anonymity protection model [11], there are important differences that invalidate some of the $k$-anonymity model assumptions. For example, our technique does not operate on a subject-specific set of images (with one image per subject only) nor is it limited to closed set scenarios. Thus, the anonymity guarantees associated with the $k$-same family of algorithms do not apply to our approach, so we use extensive experimental validation to demonstrate the feasibility of the developed deidentification pipeline.

### 3.2. Face detection and target identity estimation

Our deidentification procedure starts with a standard face detection step using the off-the-shelf Viola-Jones face detector from OpenCV [22]. The detector process the input image or video frame and returns bounding boxes of all detected faces. Each detected region is then processed separately in a sequential manner and a 4096-dimensional feature vector is extracted from each region using the pre-trained 16-layer VGG face network from [21]. For this step, the output of the last fully-connected layer of the VGG face network is considered as a feature vector. Each computed feature vector is matched against a gallery of feature vectors using the cosine similarity.

The matching procedure between the feature vector extracted from a region of the input image and the gallery of feature vectors results in an ordered list of similarity scores. Based on this list, we identify the $k$ most similar identities (where $k \ll M$) in our gallery and feed these to the generative network for surrogate-face generation. The idea of generating a synthetic surrogate face based on $k$ closest identities is similar in essence to the established family of $k$-same family of algorithms, except for the fact that the final face is in our case entirely generated by a GNN.

To generate the gallery for our deidentification approach, we process the images from the Radboud Faces Database (RaFDB) [23] with the VGG network during an offline extraction step and store templates of all $M$ identities of the RaFDB dataset in the so-called feature database, FeatDB in Fig. 1, of our pipeline. These feature vectors correspond to a finite set of facial identities that can be used for generating new (surrogate) faces for deidentification.

### 3.3. Face generation

The generative part used in our deidentification pipeline comprises a powerful GNN recently introduced by Dosovitskiy et al. in [7] for generating 2D images from 3D objects under different viewpoint angles and various basic transformations. The same architecture was later extended to another application involving face generation[1] by Michael D. Flynn. In this work, we use the same approach (and architecture) to GNNs and train our own network for deidentification.

As already suggested in the previous section, we want the network to be able to generate surrogate faces of identities or mixtures of identities contained in the feature database of our deidentification pipeline. Thus, we use RaFDB for GNN training. Once fully trained, the network is able to generate new artificial faces in accordance with the supplied identities, but also in line with other appearance-related parameters that are exposed during the training stage. The generation process can be described as follows:

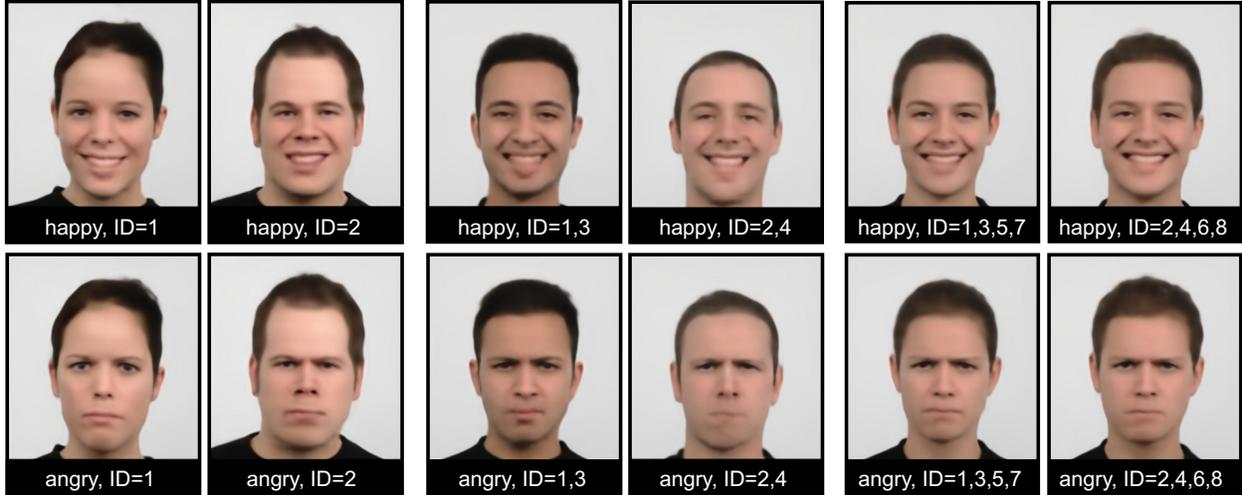$$\mathbf{x} = \text{GNN}(\mathbf{y}, \mathbf{z}), \tag{1}$$

where $\mathbf{x}$ is the output of the generative network, GNN, $\mathbf{y}$ stands for an identity-related parameter vector that encodes information about the $k$-closest identities returned by our matching procedure, and $\mathbf{z}$ denotes a parameter vector that guides the generation process and affects specific characteristics of the visual appearance of the generated output.

In the above equation, the vector $\mathbf{z}$ can in general relate to any appearance characteristic that is appropriately annotated in the training data. Since RaFDB is annotated with respect to different facial expressions, we train our network to generate faces with different identities as well as facial expressions. However, the number of appearance-related parameters exposed by the network is not limited and is in general defined by the labels available with the training data.

The generative network consists of fully-connected and deconvolutional layers as described in detail in [7]. Each deconvolution layer includes one upscaling layer followed by a convolution

---

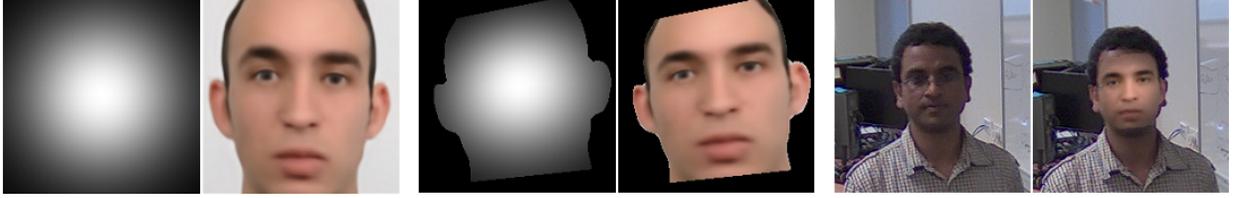[1]https://github.com/zo7/deconvfaces

**Fig. 2.** *Sample outputs from the generative network – identity mixing with $k$ identities using $k = 1$, $k = 2$ and $k = 4$ is displayed from left to the right, respectively. With an increasing number of identities the generated faces converge towards an average appearance, yet they are still realistic and without ghosting effects. The first label below each images refers to the generated facial expression and the second to the identity / identities used during the generation process.*

layer. Stabilization of the error during training is achieved by adding batch normalization layers and leaky-rectifier-unit-activation functions between the deconvolution layers. The loss of the training procedure is calculated as the pixel-wise mean square difference between the ground truth image and the artificially generated image. The training is performed with 456 images from the RaFDB dataset on a desktop PC with Intel(R) Core(TM) i7-5820K CPU (3.30GHz), 32GB of RAM, utilizing a TitanX GPU and takes around 24 hours.

With the current training dataset, our network is able to interpolate between different identities as well as different facial expressions. By combining identities, it can generate new averaged faces, almost without ghosting effect as shown in Fig. 2. By using the appearance-related parameter vector, **z**, it can also preserve the utility of the input data, e.g., in the form of facial expressions, which can be retained or altered in regards to the original input. A few sample faces, generated with our GNN are shown in Fig. 2. Here, the first block of images shows artificial images generated based on a single identity (i.e., $k = 1$), the second block shows images computed based on two identities (i.e., $k = 2$) and the last block shows sample faces generated from four identities (i.e., $k = 4$). Note that the generated faces get closer to an average appearance as the number of identities increases, yet they still appear realistic and feature no ghosting effects.

## 3.4. Face replacement

The last step of our deidentification pipeline is face replacement, during which the generated surrogate face is blended into the input image. The face replacement step starts with facial-landmark estimation using the approach from [24]. The landmarks are detected in both, the generated face and the detected original, input face. Using both sets of landmarks, we then estimate a perspective transformation that aligns the landmarks of the artificially generated face and the landmarks of the original face using RANSAC. The generated face image is then warped using this transformation in order to adjust the synthetic content over the landmarks of the original image. This correction is

**Fig. 3.** *Illustration of the replacement procedure (from left to right): the Gaussian mask, the artificially generated face image, the modified Gaussian mask, geometrically corrected synthetic face without background, sample frame, and deidentified frame.*

needed in all cases, where faces in the input images are not entirely frontal.

Following the geometric corrections, we apply a second post-processing procedure that discards the background of the generated faces. During this step, simple skin-color segmentation is performed using the *upper* and *lower* boundaries in the HSV color-space that define the skin intensities, i.e., *lower*=$[0, 10, 20]$, *upper*=$[200, 255, 255]$. Pixels with values within the defined range are retained and the rest is discarded. Erosion and dilation are then used to remove possible isolated regions that do not belong to the facial area. With this step we make sure that most of the background around the generated facial area is removed and only the facial region without the gray-colored background, as seen in Fig. 3, is swapped during deidentification.

In the last step, the warped and segmented synthetic face image is blended with the original image. Blending is performed with a Gaussian kernel mask:

$$g(x, y) = e^{-\left((x-\mu_x)^2 + (y-\mu_y)^2\right)/2\sigma^2}, \tag{2}$$

where $\mu_x = s/2$, $\mu_y = s/2$, $\sigma = s/6$, $s = \min(w, h)$, $w$ and $h$ stand for the dimensions of the generated image, and $x$ and $y$ denote image coordinates. This online generated kernel then serves as a weight mask when blending the original and generated image pixels. The kernel is warped using the same homography transformation in order to ensure the best possible face alignment and a suitable level of naturalness of the final output. The replacement procedure is illustrated in Fig. 3. Here, the first image shows the Gaussian weight mask, the second image shows the initial output of the generative network, the third image shows the Gaussian mask modified with the result of the geometric correction and segmentation step, the fourth image presents the adjusted (synthetic face) and the last two images depict an original and deidentified frame from our test dataset.

## 4. Experiments and results

In this section we present experimental results aimed at demonstrating the merits of our deidentification pipeline. We first discuss the experimental dataset and performance metrics and then present qualitative as well as quantitative results.

### 4.1. Dataset, experimental setup and performance measures

To evaluate the performance of our deidentification approach, we use the ChokePoint dataset [9], which contains video footage captured in a typical surveillance scenario. People in the videos were recorded while walking through a portal above which an array of 3 cameras was placed. The

ChokePoint videos exhibit variations across illumination conditions, pose, image sharpness, and alike and are well suited for studying the performance of deidentification technology.

The ChokePoint dataset contains 48 video sequences with a total of 64,204 frames. The videos feature 25 subjects walking through the first portal and 29 subjects walking through the second portal. For our experiments, we partitioned the video sequences into two distinct subsets. The first subset contained 24 sequences with people in mostly frontal poses, while the second subset contained the remaining 24 sequences with people in less frontal poses, i.e., profile frames. We refer to the former subset as *original* and to the latter as *profile* from hereon. The video sequences from the original subset were subjected to our deidentification approach and stored for the experimental evaluation.

To measure the efficacy of the developed deidentification pipeline, we conduct four types of verification experiments with a 10-fold cross-validation protocol. During each fold, we perform 300 legitimate (matching, client) and 300 illegitimate (non-matching, impostor) verification attempts. The different types of verification experiments are briefly outlined below:

- **Original vs. original:** In this experiment we sample 300 image pairs for the legitimate verification attempts and 300 image pairs for the illegitimate verification attempts for each experimental fold from video sequences of the original subset. The goal of this experiment is to establish the baseline performance of the recognition techniques considered in our experiments. Since video frames are sampled from the same set of videos, this experiment may be biased towards higher performances, since the appearance variability between frames is limited.

- **Original vs. profile:** In this experiment we construct the image pairs for the legitimate and illegitimate verification attempts of each fold from images taken from the original and profile subsets. Here, the first image in the pair is always sampled from the original subset and the second is always sampled from the profile subset. Because the two subsets contain distinct video sequences, this experiment better reflects the baseline performance of the recognition techniques considered in our experiments.

- **Deidentified vs. original:** This experiment is equivalent to the *original vs. original* experiment with the difference that the first image of each image pair is replaced with its deidentified version. Thus, the experiment is meant to measure the efficacy and performance of the proposed deidentification procedure. All verification attempts of all 10 cross-validation folds in this experiment have a direct correspondence in the original vs. original experiment and, therefore, clearly demonstrate the effect of deidentification on the verification performance.

- **Deidentified vs. profile:** The last experiment follows the same approach as the *original vs. profile* experiment, but replaces the video frames from the original subset with its deidentified version. The goal of this experiment is again to demonstrate the feasibility of our deidentification approach.

We report performance with standard performance metrics and graphs. Specifically, we present Receiver Operating Characteristic (ROC) curves [25], [26], which plot the value of the verification rate (VER) against the false acceptance rate (FAR) for different values of the decision threshold, and report a number of scalar performance metrics for all experiments, i.e., the equal error rate (EER), which is the operating point on a ROC curve, where FAR and 1-VER are equal, the verification rate at 1% FAR (VER-1) and the area under the ROC curve (AUC) [27]. Because we

**Fig. 4.** *Qualitative examples of deidentified frames. Each row shows a few example frames from a video sequence of the ChokePoint dataset (left image of each pair) and the corresponding deidentification results (right image of each pair). Note how the generative network is able to generate realistic renderings of faces for deidentification.*

use a 10-fold cross validation protocol, we report all metrics in the form of the mean and standard deviation computed over all experimental folds [28].

## 4.2. Qualitative evaluation

We first demonstrate the efficacy of our deidentification approach with a few qualitative examples in Fig. 4. Here, each row shows a few frames from a video sequence of the ChokePoint dataset and the corresponding deidentification result. In each image pair, the left image represents the original frame and the right image its deidentified counterpart. We can see that for the most part the deidenfied faces generated by the generative network appear natural and realistic. In this case, substituted identities are generated from 2 most similar identities (i.e. $k = 2$).

One key characteristic of our deidentification approach is the flexibility that the generative network offers when producing synthetic face images for deidentification. The generation process can be parameterized with respect to the desired target appearance of the synthetic face, which makes it possible to generate faces with different characteristics (in terms of facial expression, skin color, gender, etc.) and is important when trying to retain non-identity-related information in the deidentified data. In video conferencing applications, for example, one may want to protect the privacy of the conference participants by hiding their identity, but still preserve the information that facial expressions convey during the conversation. In customer-profiling applications, the focus is typically on the demographics of the customers (such as gender or age distributions) and not on the identity. With our deidentification approach we are able to conceal the identity of people in the image data and retain (or alter) certain aspects of the facial appearance. This characteristic is demonstrated in Fig. 5. Here, the first column shows a few sample frames from a video sequence of the ChokePont dataset and the second, third, fourth and fifth column show three different deidentification results that were rendered with a "happy", "angry", "surprised" and "neutral" facial expression, respectively. While we only show results for different facial expressions, our deidentification pipeline is in general able to generate variations of synthetic faces in accordance with any appearance-related label of the training data. Thus, if the data used for training contains images annotated with respect to facial expressions, we are able to generate faces with different facial expressions, if the data contains labels for gender, we can synthesize faces belonging to males or females and so forth. The number of different appearance variations our approach can cover is only limited by the number of available labels.
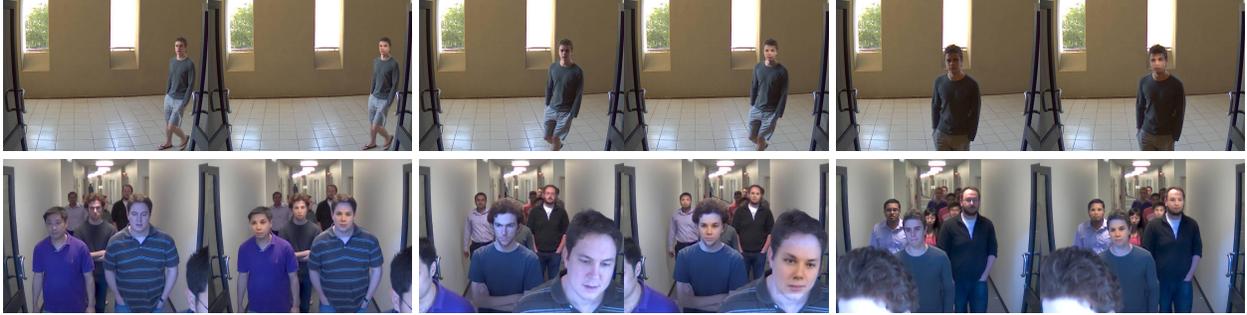
**Fig. 5.** *Deidentified frames rendered with different facial expressions. Images in the first column represent original frames from a video sequences of the ChokePoint dataset. The second, third, fourth and fifth columns show deidentified frames rendered with a "happy", "angry", "surprised" and "neutral" expression, respectively. As can be seen, our deidentification approach is highly flexible and is able to retain or alter specific aspects of the deidentified faces.*

In Fig. 6 we show some examples of visually less pleasing (or problematic) deidentification results. The image artifacts visible here are a consequence of different scene conditions (see the fourth image in the second row of Fig. 6 for an extreme example) and can be ascribed to our replacement procedure. These artifacts could be alleviated by a more elaborate face-replacement approach exploiting, for example, Poisson blending or color-profile matching. However, this would affect the speed of our pipeline, which currently runs at around 12 frames per second if processing the sequence with only one subject present at the time and around 5 frames per second if executing it on a sequence involving multiple subjects simultaneously present in a scene. These framerates were achieved on a desktop PC with Intel(R) Core(TM) i7-6700K CPU (4.00GHz) and 32GB of RAM. Another cause of image artifacts are extreme facial poses when people exit the scene (this is common to all sequences), which result in visible misalignment between the superimposed (deidentified) faces and the original facial areas.

Among the main limitations of our deidentification approach is the persistence of identity in the deidentified data. As we are dealing with video footage, the deidentification procedure should ideally produce the same (consistent) result for all frames of the given video sequence. In other words, the facial area of a given subject should be replaced with an artificially generated face of the same target identity over the entire duration of each video. However, due to changes in facial appearance, our matching module occasionally returns inconsistent results and causes changes in the target identity of the deidentified frames. This effect is demonstrated in the first row of Fig. 6, where an identity change can be observed in the last image pair due to variations in the scene's illumination despite the fact that the same subject is being deidentified. Nevertheless, because the target identity for deidentification is determined with a state-of-the-art face recognition model (i.e., VGG [21]), our procedure is able to assign a consistent target identity most of the time for all test videos considered in our experiments.

In Section 5, where we discuss possible directions for future work, we propose some possible

**Fig. 6.** *Qualitative examples of problematic deidentification results. The upper row shows an identity switch in the last frame due to a change in the scenes illumination. The lower row shows difficulties due to the presence of multiple people, some of which occlude faces in the background. Misalignment between the original and surrogate faces is also visible in the second image pair of the lower row, which happens due to extreme viewing angles when people exit the scene.*

improvements of our deidentification pipeline, which address most of the existing issues of the current implementation.

### 4.3. Automatic and manual reidentification

The quality and efficacy of deidentification techniques is typically measured through reidentification experiments [29], where the goal is to evaluate the risk of successfully identifying a person from deidentified data. This risk is commonly assessed with automatic and manual recognition experiments.

As outlined in Section 4.1, we perform a number of verification experiments in a 10-fold cross validation protocol towards the risk assessment and consider three state-of-the-art automatic recognition approaches from the literature. Specifically, we use *i)* the open-source implementation of the 16-layer VGG face network from [21] – VGG from hereon, *ii)* our own 24-layer implementation of the SqueezeNet network from [30] trained on around 2.5 million images (i.e., the VGG network training data) – SqueezeNet from hereon, and *iii)* the 4SF algorithm from OpenBR (version 1.1) [31] – OpenBR from hereon. For the two networks, we use the output from the last fully connected layer of each network as a feature vector and compute a similarity score for an image pair as the cosine angle between the two corresponding feature vectors. For OpenBR we use the default matching option.
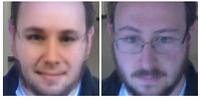
We also conduct manual recognition experiments using a similar 10-fold cross-validation protocol as with automatic techniques, but limit the extend of comparisons to 5% of the automatic experiments. Thus, 30 comparisons are performed during each fold resulting in a total of 300 verification attempts (150 legitimate and 150 illegitimate experiments) in each experimental run. Each verification experiment was evaluated by one human evaluator, i.e., four evaluators covered four verification experiments. To produce similarity scores needed for generating performance metrics and ROC curves, we manually assign a similarity score from a five-point scale to each comparison in accordance with the methodology proposed in [32].

Similar to other existing works on face deidentification, our approach tries to conceal the identity of people by replacing the detected facial areas with a synthetically generated surrogate faces. However, identity cues can also be extracted from contextual information that is not directly related to facial appearance. For example, the facial outline, hair-style, or even clothing can represent a

give-away that recent recognition techniques based on deep models as well as humans may be able to pick up. To explore this issue, we conduct two sets of experiments:

- With context: here we feed the facial area to the recognition technique directly as it is detected by the face detector. Thus, the facial area also contains contextual information about the shape of the head, hair style and alike. A comparison of two images with context is illustrated in the last column of Table 1 (first row).

- Without context: here we trim the bounding box returned by the face detector on each side by 10%, the facial areas used for the recognition experiments are therefore cropped tighter and contain only little contextual information. A sample comparison of two images as used in this set of experiments is shown in the last column of Table 1 (second row).
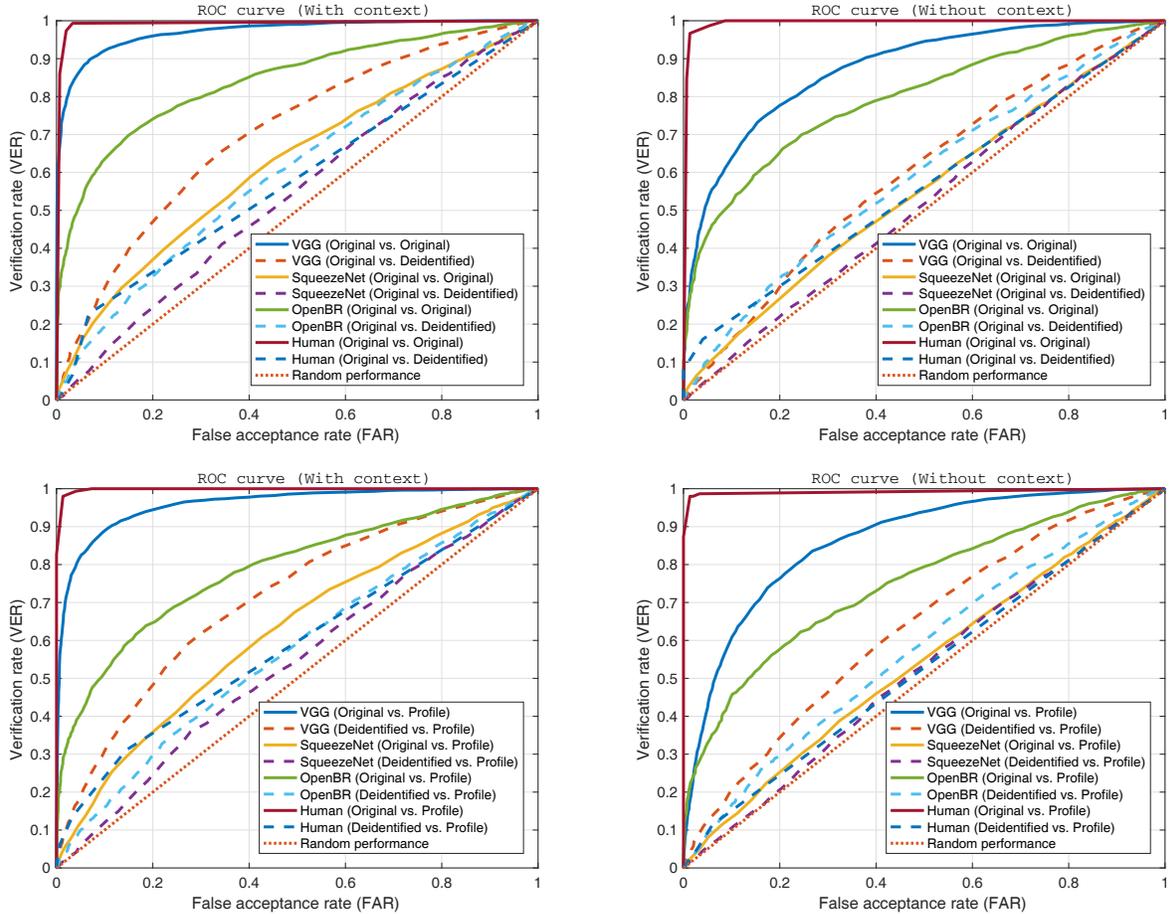
**Table 1** Quantitative results of the experiments. Average values and standard deviations over 10-fold are presented for all performance metrics.

| Test description | | Original-to-original | | Original-to-profile | | Deidentified-to-original | | Deidentified-to-profile | | Context illustration |
|---|---|---|---|---|---|---|---|---|---|---|
| Metric (in %) | | EER | VER-1 | EER | VER-1 | EER | VER-1 | EER | VER-1 | |
| *Context* | VGG | $8.7 \pm 1.0$ | $70.7 \pm 5.7$ | $10.5 \pm 1.3$ | $56.0 \pm 10.3$ | $34.4 \pm 2.2$ | $4.2 \pm 3.0$ | $34.5 \pm 1.2$ | $5.5 \pm 3.2$ |  |
| | SqueezeNet | $40.8 \pm 2.8$ | $4.4 \pm 2.1$ | $40.9 \pm 1.9$ | $3.3 \pm 1.6$ | $47.3 \pm 2.3$ | $0.8 \pm 0.7$ | $47.4 \pm 2.9$ | $1.2 \pm 1.0$ | |
| | OpenBR | $23.6 \pm 2.2$ | $34.5 \pm 6.8$ | $28.3 \pm 1.7$ | $24.2 \pm 6.3$ | $42.8 \pm 1.8$ | $2.8 \pm 1.8$ | $45.3 \pm 2.8$ | $2.1 \pm 1.4$ | |
| | Human | $2.0 \pm 2.8$ | $n/a$ | $1.0 \pm 2.2$ | $n/a$ | $42.0 \pm 7.6$ | $n/a$ | $41.8 \pm 7.2$ | $n/a$ | |
| *No Context* | VGG | $21.5 \pm 2.9$ | $26.1 \pm 6.3$ | $21.8 \pm 1.7$ | $13.1 \pm 5.7$ | $43.3 \pm 2.2$ | $1.6 \pm 1.3$ | $40.6 \pm 2.0$ | $3.4 \pm 1.5$ |  |
| | SqueezeNet | $47.0 \pm 2.1$ | $3.8 \pm 1.3$ | $47.1 \pm 2.1$ | $1.8 \pm 0.9$ | $49.4 \pm 1.8$ | $0.9 \pm 0.7$ | $48.4 \pm 2.1$ | $1.4 \pm 1.3$ | |
| | OpenBR | $27.8 \pm 2.0$ | $19.1 \pm 6.6$ | $32.2 \pm 3.1$ | $15.9 \pm 5.7$ | $43.9 \pm 2.9$ | $1.9 \pm 1.5$ | $45.2 \pm 2.9$ | $1.1 \pm 0.9$ | |
| | Human | $2.3 \pm 3.2$ | $n/a$ | $1.7 \pm 2.8$ | $n/a$ | $44.0 \pm 8.4$ | $n/a$ | $47.3 \pm 5.8$ | $n/a$ | |

Numerical results of the experiments are presented in Table 1. Note that VER-1 values are not reported for the manual experiments (denoted as Human) because of an insufficient number of manually graded image comparisons. As expected, the results with contextual information are significantly better than those without contextual information for all experiments when non-deidentified images are used. When the verification attempts are conducted with deidentified images, contextual information still contributes to a higher performance in all experiments, but the differences between images with and without context are smaller. As also evidenced by the ROC curves of the experiments in Fig. 7, the best performing automatic technique, the VGG network, is able to ensure a recognition performance well above random with an EER of 34.4% for the *deidentified-vs-original* experiment and an EER of 34.5% for the *deidentified-vs-profile* experiment when context is available. If no contextual information is present, the VGG performance drops to an EER of 43.3% and 40.6% for the same experiments, respectively. These observations suggest that contextual information is important and may be exploited by contemporary recognition techniques to boost performance. Thus, care needs to be taken to appropriately conceal, modify or remove contextual information from the data as well.

Another interesting observation that can be made from the ROC plots in Fig. 7 is the drop in performance for the manual experiments. On the unaltered images, human performance is close to perfect for all experiments. However, after deidentification human performance drops to (more or less) random if no contextual information is present and is only slightly better than chance if contextual information is available.
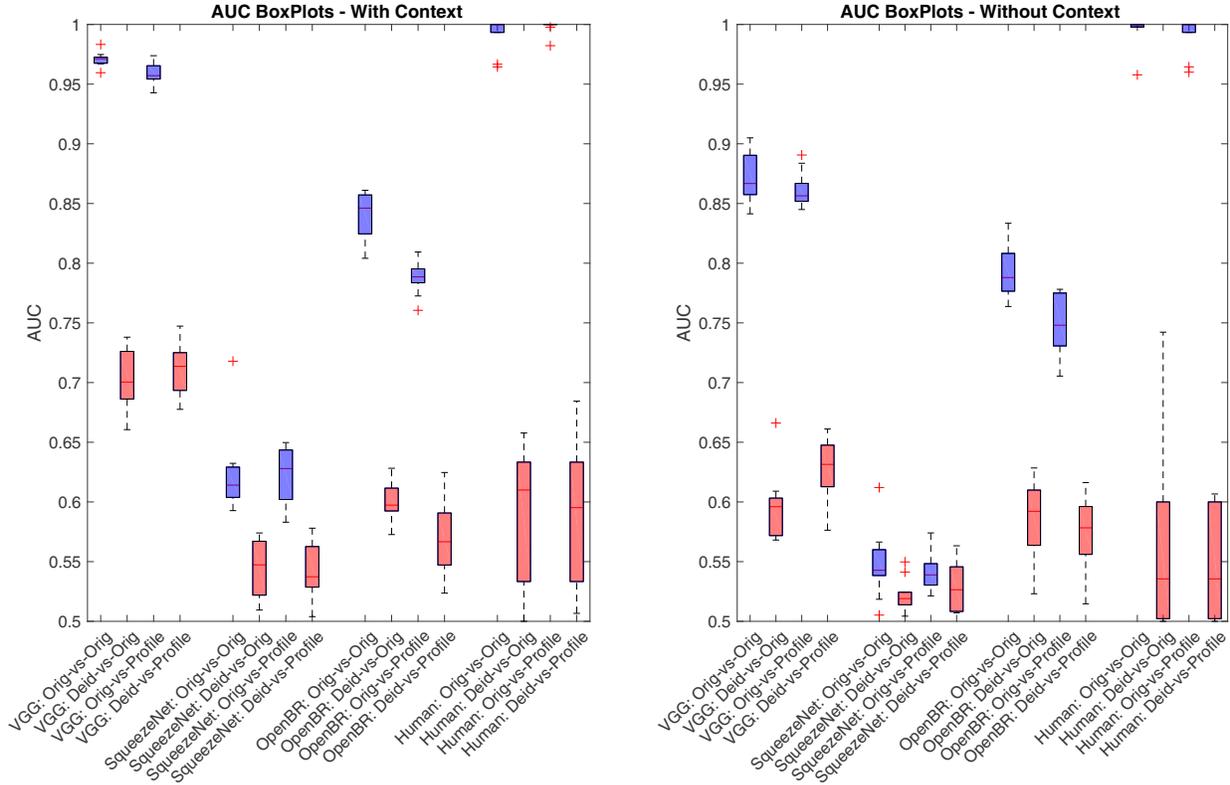
The observations made so far are also supported by the box-and-whiskers plots of the AUC values computed from the 10 experimental folds in Fig. 8. Here, a value of 0.5 indicates random

**Fig. 7.** *ROC curves of the verification experiments. The curves on the left show the results of the experiments with images with contextual information and the curves on the right show the results obtained without contextual information. The upper row shows experiments with unaltered and deidentified images from the original subset and the lower row shows experiments with unaltered and deidentified images from the profile subset. All results show that our approach is effective.*

performance. The blue box plots show results for unaltered images and the red box plots show the results for the deidentified images. It needs to be noted that even in cases, when the performance after deidentification is not exactly random, it is still significantly lower than that obtained with unaltered images for all tested techniques. The highest median AUC value of any experiment after deidentification (AUC = 0.719) is achieved by the VGG network when contextual information is available. However, while this value is significantly above random it is of limited use to applications requiring reliable face recognition.

In our last experiment, we compare our deidentification approach to existing deidentification techniques from the literature. Specifically, we report results for two naive methods, i.e., blurring and pixelation, which unlike techniques from the $k$-Same family can be applied to video data using the same experimental protocol as used in the previous experiments. The results of the comparison are generated with the best performing recognition approach from Table 1, i.e., the VGG network, and are presented in Table 2. As can be seen, the naive methods result in worse recognition performance than our approach and therefore appear to ensure better anonymity. However, these

**Fig. 8.** *AUC values from the 10 experimental folds presented in the form of box plots for all assessed techniques as well as human experiments. The blue plots show results with unaltered images, the red plots show experiments with deidentified images. The left box plots present the results with contextual information and the right box plots present the results without context. Note that after deidentification (red plots) the result are very close to 0.5 which indicated random performance.*

methods destroy most of the information content of the images and can to a certain extent also be bypassed as shown by the results of the parrot (or imitation) attack experiments.

On the right side of Table 2 we show some qualitative deidentification examples on a closed set of images from the XM2VTS dataset [33] (top row). Note that a closed set is required for the $k$-Same family of techniques to be applicable. In accordance with the $k$-anonymity scheme [11], we replace clusters of (in this case $k = 2$) images with the same surrogate face generated by our GNN (the clusters are color-coded in the image). The results of our deidentification approach (last row) are visually convincing and feature no ghosting effects, such as the images generated by the original $k$-Same approach from [3] (fourth row). With our approach it is also possible to retain certain aspects of the original data, which is not necessarily true for the blurred and pixelated images, shown in the second and third row of the image, respectively.

**Table 2** Deidentification performance with the VGG network. The left part of the table shows a comparison of a few existing (naive) deidentification techniques and the proposed approach in experiments on the ChokePoint dataset. The right part of the table presents a qualitative comparison of our approach with competing techniques from the literature on a closed set of images.

| Deidentification technique | Context | | No Context | | Qualitative comparison (closed set) | |
|---|---|---|---|---|---|---|
| | EER (in %) | VER-1 (in %) | EER (in %) | VER-1 (in %) | | |
| Pixelated | $45.1 \pm 1.8$ | $1.7 \pm 0.9$ | $47.3 \pm 2.7$ | $1.3 \pm 0.8$ |  | Original |
| Blurred | $37.0 \pm 1.4$ | $1.7 \pm 1.2$ | $43.0 \pm 2.1$ | $1.5 \pm 1.5$ | | Blurred |
| Pixelated (parrot attack) | $38.0 \pm 1.4$ | $3.0 \pm 1.6$ | $39.4 \pm 2.6$ | $2.7 \pm 1.3$ | | Pixelated |
| Blurred (parrot attack) | $32.4 \pm 1.7$ | $13.9 \pm 4.2$ | $35.8 \pm 2.1$ | $8.6 \pm 3.4$ | | $k$-Same |
| Ours | $34.3 \pm 2.2$ | $4.2 \pm 3.0$ | $43.2 \pm 2.2$ | $1.6 \pm 1.3$ | | This paper |
| No deidentification | $8.7 \pm 1.0$ | $70.7 \pm 5.7$ | $21.5 \pm 2.5$ | $26.1 \pm 6.3$ | | |

## 5. Conclusion

In this paper we have presented a novel approach to face deidentification using generative neural networks. The proposed approach was evaluated on the ChokePoint dataset with highly encouraging results. Our evaluation suggests that generative networks are a viable tool for face deidentification and that a high degree of anonymity can be ensured by swapping the original faces by artificially generated surrogate faces. Furthermore, our experiments show that due to the flexibility of the generative network it is possible to control the appearance of the generated surrogate faces and thus retain (or alter) only specific aspects of the input images – contributing significantly to the utility of the deidentified faces.

While our deidentification results are visibly convincing, additional improvements are possible. As part of our future work, we plan on including additional generator parameters to further capitalize on the utility of the deidentified faces. Other possible improvements include a better blending procedure that would improve the overall naturalness of the deidentified faces and remove artifacts. We will also consider incorporating a tracking scheme, which would improve the applicability of our approach on video data.

## Acknowledgement

## 6. References

[1] Ribarić, S., Ariyaeeinia, A., Pavešić, N.: 'De-identification for privacy protection in multimedia content: A survey', Signal Processing: Image Communication, **47**, 2016, pp. 131–151.

[2] Neustaedter, C., Greenberg, S., Boyle, M.: 'Blur filtration fails to preserve privacy for home-based video conferencing'. TOCHI, 2005, pp. 1–36.

[3] Newton, E.M., Sweeney, L., Malin, B.: 'Preserving privacy by de-identifying face images'. TKDE, 2005, **17**, 2, pp. 232–243.

[4] Gross, R., Airoldi, E., Malin, B., Sweeney, L.: 'Integrating utility into face de-identification'. PET, 2005, pp. 227–242.

[5] Sim, T., Zhang, L.: 'Controllable Face Privacy', AFGR, 2015, pp. 1–8.

[6] Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: 'Generative Adversarial Networks', arXiv:1406.2661, 2014.

[7] Dosovitskiy, A. and Sprigenberg, T. J., Brox T.: 'Learning to generate chairs with convolutional neural networks', CVPR, 2015, pp. 1538–1546.

[8] Larsen, A. B. L., Sonderby, S. K., Larochelle, H., Winther, O. : 'Autoencoding beyond pixels using a learned similarity metric', arXiv:1512.09300, 2015.

[9] Wong, Y., Chen, S., Mau, S., Sanderson, C., Lovell, B.C.: 'Patch-based Probabilistic Image Quality Assessment for Face Selection and Improved Video-based Face Recognition', CVPRW, 2011, pp. 81–88.

[10] Ribarić, S., Pavešić, N.: 'An Overview of Face De-identification in Still Images and Videos'. AFGR, 2015, pp. 1–6.

[11] Sweeney, L.: '$k$-anonymity: a model for protecting privacy'. Uncertainty, Fuzziness, and Knowledge-Based Systems, 2002, **10**, 5, pp. 557–570.

[12] Machanavajjhala, A., Kifer, D., Gehrke, J., Venkitasubramaniam, M.: '$l$-diversity: Privacy beyond $k$-anonymity'. TKDD, 2007, **1**, 1, article 3.

[13] Li, N., Li, T., Venkatasubramanian, S.: '$t$-closeness: Privacy beyond $k$-anonymity and $l$-diversity'. ICDE, 2007, pp. 106–115.

[14] Gross, R., Sweeney, L., de la Torre, F., Baker, S.: 'Semi-Supervised Learning of Multi-Factor Models for Face De-Identification', CVPR, 2008, pp. 1–8.

[15] Gross, R., Sweeney, L., de la Torre, F., Baker, S.: 'Model-Based Face De-Identification'. PRV, 2006, pp. 1–8.

[16] Jourabloo, A., Yin, X., Liu, X.: 'Attribute Preserved Face De-identification', ICB, 2015, pp. 278–285.

[17] Samaržija B., Ribarić, S.: 'An approach to the de-identification of faces in different poses'. MIPRO, 2014, pp. 1246–1251.

[18] Sun, Z., Meng, L., Ariyaeeinia A.: 'Distinguishable de-identified faces', AFGR, 2015, pp. 1–6.

[19] Brkić, K., Hrkać, T., Sikirić, I., Kalafatić, Z.: 'Towards neural art-based face de-identification in video data', SPLINE, 2016, pp. 1–5.

[20] Chriskos, P., Zoidi, O., Tefas, A., Pitas, I.: 'De-identifying facial images using projections on hyperspheres', AFGR, 2015, pp. 1–6.

[21] Parkhi, M. O., Vedaldi A., Zisserman A.: 'Deep Face Recognition', BMVC, 2015, article 41.

[22] Viola, P., Jones M.: 'Robust real-time face detection', IJCV, 2004, **57**, 2, pp. 137–154.

[23] Langner, O., Dotsch, R., Bijlstra, G., Wigboldus, D.H.J., Hawk, S.T., van Knippenberg, A.: 'Presentation and validation of the Radboud Faces Database', Cognition&Emotion, 2010, **24**, 8, pp. 1377–1388.

[24] Kazemi, V., Sullivan J.: 'One Millisecond Face Alignment with an Ensemble of Regression Trees', CVPR, 2014, pp. 1867–1874.

[25] Fawcett, T.: 'An introduction to ROC analysis'. PRL, 2006, **27**, 8, pp. 861–874.

[26] Peer, P., Emeršič, Ž., Bule, J., Žganec-Gros, J., Štruc, V.: 'Strategies for exploiting independent cloud implementations of biometric experts in multibiometric scenarios'. MPE, **2014**, pp. 1–15.

[27] Gajšek, R., Štruc, V., Dobrišek, S., Mihelič, F.: 'Emotion recognition using linear transformations in combination with video'. Interspeech, 2009, pp. 1967–1970.

[28] Emeršič, Ž., Štruc, V., Peer, P.: 'Ear Recognition: More than a Survey'. Neurocomputing, 2017, doi:10.1016/j.neucom.2016.08.139

[29] Garfinkel, S.: 'De-Identification of Personal Information', NISTIR 8053, 2015, pp. 1–46.

[30] Iandola, F.N., Han, S., Moskewicz, M.W., Ashraf, K., Dally, W.J., Keutzer, K.: 'SqueezeNet: AlexNet-level accuracy with $50\times$ fewer parameters and $< 0.5$ MB model size', arXiv:1602.07360, 2016.

[31] Klontz, J., Klare, B., Klum, S., Jain, A., Burge., M.: 'Open Source Biometric Recognition'. BTAS, 2013, pp. 1–8.

[32] Phillips, J., OToole, A.: 'Comparison of human and computer performance across face recognition experiments', IVC, 2014, **32**, 1, pp. 74-85.

[33] Messer, K., Kittler, J., Sadeghi, M., Marcel, S., Marcel, C., Bengio, S., Cardinaux, F., Sanderson, C., Czyz, J., Vandendorpe, L., Srisuk, S., Petrou, M., Kurutach, W., Kadyrov, A., Paredes, R., Kepenekci, B., Tek, F.B., Akar, G.B., Deravi, F., Mavity, N.: 'Face Verification Competition on the XM2VTS Database'. AVBPA, 2003, pp. 964–974.