# ID2image: Leakage of non-ID information into face descriptors and inversion from descriptors to images

Mingrui Li[●], William A. P. Smith[●], and Patrik Huber[●]

University of York, York YO10 5DD, UK
{ml1652,william.smith,patrik.huber}@york.ac.uk

**Abstract.** Embedding a face image to a descriptor vector using a deep CNN is a widely used technique in face recognition. Via several possible training strategies, such embeddings are supposed to capture only identity information. Information about the environment (such as background and lighting) or changeable aspects of the face (such as pose, expression, presence of glasses, hat etc.) should be discarded since they are not useful for recognition. In this paper, we present a surprising result that this is not the case. We show that non-ID attributes, as well as landmark positions and the image histogram can be recovered from the ID embedding of state-of-the-art face embedding networks (VGGFace2 and ArcFace). In fact, these non-ID attributes can be predicted from ID embeddings with similar accuracy to a prediction from the original image. Going further, we present an optimisation strategy that uses a generative model (specifically StyleGAN2 for faces) to recover images from an ID embedding. We show photorealistic inversion from ID embedding to face image in which not only is the ID realistically reconstructed but the pose, lighting and background/apparel to some extent as well.

## 1 Introduction

State-of-the-art face recognition relies on the use of deep neural networks (usually CNNs) to embed a face image to an identity vector [6,26,23]. A measure of distance in this embedding space is used to represent dissimilarity in identity. The goal of training such networks is to minimise the within-class scatter while maximising the between-class scatter for all identities. The former goal necessitates that the embedding should depend only on the identity of the person in the image. Environmental conditions such as the lighting, background and properties of the camera as well as changeable aspects of the face such as pose, expression and the presence of accessories should not affect this embedding (i.e. should not introduce within-class scatter). In other words, the embedding network should learn invariance to these factors. In this paper, we explore a number of unanswered questions about some of the properties of these embeddings.

**ID to image inversion**     First, we ask whether it is possible to create an image from an ID vector that correctly recreates the identity of the person.

Conventional text-based passwords are usually passed through a cryptographic hash function for storage [3]. Since these are pseudo one-way functions, it is extremely difficult to invert an encrypted password to cleartext, with brute force or dictionary attacks the only options. For password verification, only the encrypted version need be stored and a leak is not critical due to the difficulty of inversion. It is tempting to assume that ID embeddings from face images possess similar characteristics. For example, the Face ID facial recognition system developed by Apple Inc. makes this argument in its advertising to reassure users:

> *Face ID doesn't store an image of your face. Instead of storing an image, Face ID saves a mathematical value created from the characteristics of your features. It's impossible for anyone else to recreate your likeness from this.* [14]

However, since the embedding of a face image to an ID vector is a noisy process, it cannot be assumed that an identical embedding will arise from any image of the same person's face. Therefore, the ID vector cannot be passed through a hash function for storage since the hashed value will change dramatically with small changes in the ID vector. From a security perspective, this means that raw ID vectors must be stored for future identification. From an inversion perspective, it means that similar images map to similar ID vectors. Also, embeddings based on deep neural networks are computed via a differentiable function. This means that it is possible to optimise an image representation to minimise an ID vector loss, thereby reconstructing an image with the desired identity (subject to suitable regularisation, which we achieve via a generative model).

**Non-ID information leakage**     Second, we ask whether ID embeddings truly contain only ID-related information. The engineering of training datasets, network architectures and loss functions has been widely studied in the face recognition literature in order to satisfy the goal of invariance to non-ID factors in the input image. Datasets are created specifically to introduce lots of variation in these non-ID factors. Then, by designing a loss function that encourages the same ID to embed to the same point, invariance to these factors is hopefully learnt. In this paper, we investigate to what extent this has been achieved. In particular, we ask: how well do modern face embedding networks successfully remove non-ID factors when embedding a face image?

**Inversion including non-ID factors**     Finally, we connect these two lines of investigation by asking whether it is possible to recover an image from an ID vector that not only captures the correct face identity but also non-ID characteristics of the actual image that was used to compute the ID vector. The possibility of ID vector to image inversion raises privacy and security questions.

For the reasons mentioned above, ID descriptor vectors cannot be stored securely hashed. This means that any third party with whom identity must be verified is receiving an encoding of a face image from which an image can be recovered. Where non-ID information leaks into this representation, it means the image itself can potentially be recovered. This could, for example, leak un-

intended information in the background of the image or maybe an unflattering image that the user would not wish to be made public.

## 2   Related work

**GAN inversion**   General adversarial networks (GANs) have become a popular tool in computer vision and machine learning. In 2019 and 2020, Karras et al. [15,16] set the benchmark with StyleGAN and StyleGAN2, able to create photo-realistic face images at high resolution. A defining feature of StyleGAN is in its generator architecture. Instead of feeding an input latent code $z \in Z$ only to the beginning of the network, a mapping network first transforms it to an intermediate latent code $w \in W$. Affine transforms then produce styles that control the layers of a synthesis network. In addition to that, stochastic variation is facilitated by providing random noise maps to the synthesis network.

Recent work has shown that GANs can encode a rich set of semantics in their latent space. In addition to generating images, recently, attempts have been made to invert the GAN generating process from the image back to latent space for the purpose of image manipulation or analysis, which is widely known as GAN inversion. To accomplish that, most existing works either learn an extra encoder or regressor separate from the GAN (e.g. [5,25]), or directly optimise the latent code to fit a target image (e.g. [32]), or a combination of the two (i.e. initialising the optimisation with the result from a regressor, e.g. [34]).

For StyleGAN specifically, recently, several works have shown that it is possible to retrieve the latent code $w$ of a target image [27,28]. The works show that inverting to the latent space $W$ is easier than to $Z$. However, accurately reconstructing a target image is still an ongoing challenge. In another recent work, Abdal et al. [1,2] propose a framework to project an image into the latent space $W+$, where $W+$ contains separates latent vectors for the specific scales of StyleGAN, hence effectively reconstructing different levels of features in the target image. Finally, Yin et al. [33] introduce a deep-inversion method that inverts a pretrained neural network to generate synthesised class-conditional input images for data-free knowledge transfer.

What all these methods have in common is that the inversion they do usually means going from an image to a latent code. Much different to that, we investigate *ID-descriptor to image* inversion, which uses StyleGAN as a generative model (i.e. ID→StyleGAN code→image).

**ID descriptor information & inversion**    In a seminal work, in 2015, Mahendran and Vedaldi [21] set out to analyse the visual information contained in both shallow (e.g. HOG) and deep feature representations, to investigate the question: *given an encoding of an image, to which extent is it possible to reconstruct the image itself*. They propose an optimisation method to invert representations using gradient descent. Among their findings are that networks retain rich information even at deep levels and that a progressively more invariant and abstract notion of the image content is formed in the network. As face identity

descriptors usually are made up of the final layer of a deep network, it would therefore be the most invariant and abstract representation.

Few works so far have investigated the inversion of a face descriptor back to a face image. Genova et al. [13] train an image to 3DMM parameter regressor in an unsupervised manner, where the key idea is an ID loss between ID descriptor of the original image and the ID descriptor of a 3DMM image rendered with a differentiable renderer. The only trainable part of their system is the *ID descriptor to 3DMM parameter* regressor. This constitutes an ID inversion network - but they do not reconstruct the image itself, so no non-ID information is recovered.

Following on from the general idea of an ID loss, there are a number of works that employ an identity loss in 3D face model fitting, for example GAN-FIT [12]. Cole et al. [7] also perform an ID-only inversion for the purposes of face frontalisation. They assume that the face encoder successfully removes all non-ID information and train to reconstruct only frontal image landmarks and textures. Since they do not use a GAN, their results are not photorealistic. Some recent works [24,10] have investigated a so-called *black-box attack*, which is whether given only an ID descriptor, and no access to an attacked model, one could reconstruct an image of the face, and they have presented encouraging (technically) as well as concerning (from a privacy perspective) results.

In terms of non-ID information contained in face descriptors, in an early work, Kumar et al. [18] came to the perhaps surprising result that using a so-called *inverse crop*, where the face is cut out of an image, leads to surprisingly high face recognition rates on LFW. It should be noted though that these inverse crops do contain hair and part of ears/chin, and that LFW is a fairly simple dataset (recognition rates achieved are over 99%, even a decade ago) - so it is largely unexplored, if any background or other non-ID information is present in face descriptors, especially in today's state-of-the-art networks. To the best of our knowledge, no work so far has investigated if any non-ID properties can be recovered from identity descriptors.

**Privacy leakage & Adversarial learning**   Past studies have shown that face recognition networks encode soft biometrics (e.g. age, race and gender) while training [11]. This indicates that face descriptors are at risk of privacy leaks. Privacy leakage in face representation is a vital issue, as it can raise several concerns associated with the unauthorised extraction of an individual's information. More recent studies have shown that the extraction of such soft biometric data can improve the performance of facial recognition systems [22,9].

Inspired by GANs, several adversarial learning methods have been proposed that remove or mitigate sensitive information stored in the training data or learn somehow obfuscated features. Alvi et al. [4] are inspired by domain and task adaptation methods, they propose a joint learning and unlearning method to remove bias from neural network embedding. To unlearn the bias, the authors adversarially minimise the classification loss and confusion loss. The confusion loss is computed by calculating the cross-entropy between classifier output and uniform distribution. Li et al. [19] deployed several neural networks to simulate an attacker that attempts to reconstruct entire raw data. An obfuscator is
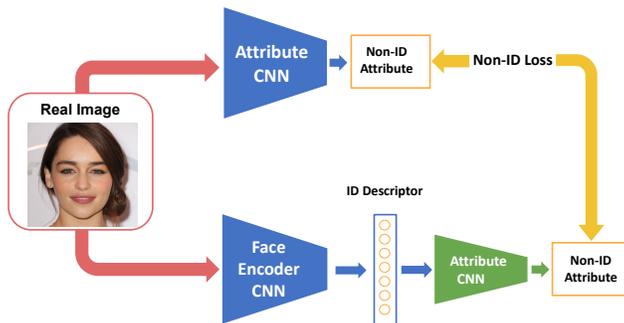
**Fig. 1.** Non-ID attribute regression via an ID bottleneck. Only the green component is trained: an MLP that maps an ID vector to the appropriate attribute (such as expression, landmarks, image histogram etc). The labels either come from a pretrained (and fixed) attribute estimation network that takes an image as input, they are provided as manually assigned labels or they are computed directly from the image.

trained to hide privacy-related information, which ensures that the attacker cannot train using features provided by the obfuscator to accurately infer privacy attributes or reconstruct the raw data. Dhar et al. [9] propose to adversarially minimise gender predictability to reduce gender bias from pretrained face descriptors. They show that gender bias in face recognition is correlated with the power of face descriptors to predict gender. The face descriptors with low gender predictability generally demonstrated lower gender bias in face verification.

## 3  Non-ID attribute prediction from ID

We begin by exploring to what extent we are able to estimate non-ID "attributes" from an ID descriptor provided by a pretrained face encoder CNN. We use "attribute" here in very general terms, including image-based attributes such as landmark positions and colour histograms and non-ID face attributes such as the presence or absence of a smile, glasses or hat. For each attribute, we train an MLP that maps from an ID descriptor to the target attribute. All of our MLPs are trained on CelebA [20], which includes 202,599 celebrity face images with 40 binary attributes. We aligned and cropped the original CelebA to a VGGFace2 compatible version and scaled all images to resolution 224. We train in a supervised fashion using either labelled real data or synthetically generated data. In all cases, we embed images $\mathbf{i}$ to an ID descriptor, $\mathbf{d} = \text{VGG}(\mathbf{i})$ using the VGGface2 face encoder [6], where $\mathbf{d} \in \mathbb{R}^{2048}$. Fig. 1 illustrates our proposed Non-ID attribute regression framework.

### 3.1  Discrete Binary Attributes

We begin by estimating discrete binary attributes. These have been manually labelled as part of the CelebA [20] dataset. Then we train an MLP to classify

| | Attribute | | | Landmarks | Histogram |
|---|---|---|---|---|---|
| | Smiling | Wearing_Hat | Eyeglasses | mean | EMD |
| From ID (VGGFace2 [6]) | 91.0% | 99.0% | 99.7% | 9.3% | 2.68 |
| From ID (ArcFace [8]) | 81.7% | 96.7% | 96.7% | 7.3% | 2.45 |
| From image [20] | 92% | 99% | 99% | 8.2% | 0.0 |
| Baseline | 50.4% | 96.7% | 94.0% | 11% | 2.97 |

**Table 1.** Quantitative results for attribute prediction (discrete binary attributes, landmarks and image histogram) from ID vectors (row 1 and 2) and images (row 3). In row 4 we show baseline performance in which we simply always predict the most common class, the mean landmarks or the mean histogram respectively. For attribute prediction, higher is better, and for the landmark and histogram prediction, lower is better.

the binary class using Binary Cross Entropy loss, Adam with learning rate $10^{-3}$, for 20 epochs. We train separate networks for the smiling, glasses and wearing hat attributes. Our MLP consists of 3 fully connected layers with 256 hidden neurons and a sigmoid output layer.

### 3.2   Histogram regression

In this section, we examine whether background colour and lighting information can be restored from ID by predicting the histogram of the RGB channel. Histograms of RGB intensities provide a global summary of an image that encapsulates not only ID but also environment-related features such as background, camera settings and lighting. For each image, we compute a ground truth hard histogram and use this as the label for training. Here, the three colour values are binned into fixed width bins, with one histogram per colour channel. We found that a very small MLP provides best performance for the task of image histogram regression from ID vector. We use 2 fully connected layers with ReLU activation and 8 neurons per layer. We apply softmax to the output layer such that the output represents a normalised histogram. We use $N = 10$ bins. We train this network using mean squared error loss, Adam with learning rate of $10^{-6}$, and batch size 32, for 20 epochs.

### 3.3   Landmark regression

Finally, we attempt to directly regress the coordinates of 68 face landmarks from the ID vector. We apply the dlib landmark detector [17] to all images in the CelebA dataset [20]. We use these as pseudo ground truth labels for our regressor. We use a three-layer MLP to predict landmarks from ID vector with 256 neurons per hidden layer and 136 outputs for the 2D coordinates of each landmark. For comparison, we also train a more conventional image to landmark regressor using a CNN. We use a simple architecture comprising 5 convolutional layers followed by 2 fully connected layers. The activation function is ReLUs. Both image to landmark and ID vector to landmark networks are trained using

**Fig. 2.** Examples of correctly (left) and incorrectly (right) classified samples. We show the original images but note that the classification is done *only on the ID vectors derived from these images.*

mean squared error loss and the SGD optimiser with learning rate $10^{-3}$, batch size 16, and for 150 epochs.

### 3.4    Results

We now evaluate our prediction of non-ID attributes from ID vectors, as provided by the VGGFace2 network. We show quantitative results for all attributes in Table 1. The evaluation images we used are 15k test images from CelebA with the remaining 188k images for training. To validate our results, we repeat the experiments with ID vectors generated by the ArcFace network.

For discrete binary attributes, we show the percentage classified correctly. Our result regressed from the ID vector is shown in the first and second rows. This shows that non-ID leakage exists in different networks. For comparison in the third row we show the result from [20] computed *from the original image.* In the fourth row, we show the baseline performance obtained by always guessing the more common class for the binary attribute prediction, the mean landmarks for the landmark prediction, and the mean histogram for the histogram prediction. With the exception of the wearing hat attribute for the ArcFace embedding, we can see that we significantly outperform that baseline and, remarkably, match or even exceed the performance of an image-based method despite only having access to an ID vector that should be independent of these non-ID attributes. In Fig. 2 we show some examples of correctly and incorrectly classified samples. It is interesting to note that quite subtle smiles are encoded in the ID vectors such that we correctly classify them and, even in the case of the false positives shown, there are still smile-like features in the wrinkles around the mouth. Similarly, the false positives for wearing hat are in fact wearing headgear.

We now evaluate image-based attributes. In Table 1 we show quantitative results for landmark prediction in the fifth column and histogram prediction in the
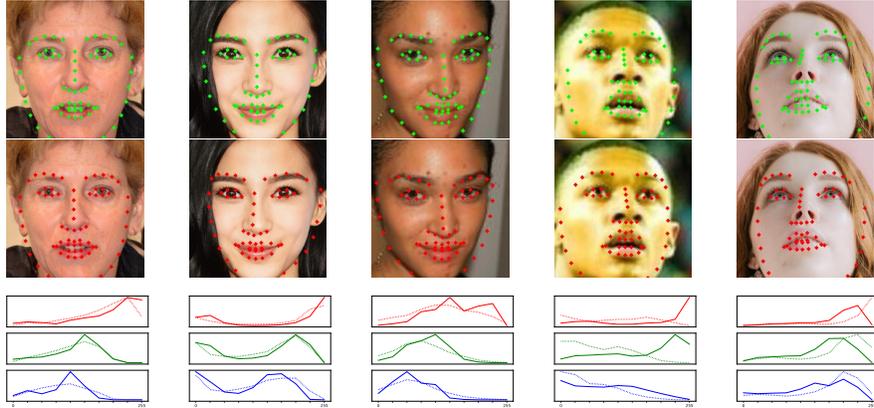
**Fig. 3.** Qualitative results for histogram and landmark regression. row 1: input image with ground truth landmarks, row 2: landmarks regressed from ID vector, row 3: ground truth image histograms (dotted) and histograms regressed from ID vector (solid).

sixth column. For the landmark error, we show Euclidean distance averaged over landmarks expressed as a percentage of the interocular distance. For histogram error, we show the Earth Mover's distance to ground truth. Our prediction from ID vector outperforms the baseline and is only marginally worse than prediction from images (in the case of landmarks). In the case of histogram, the prediction from images is exact. In Fig. 3 we show qualitative results for landmark and histogram estimation from ID vectors. In the first row we show the original image with ground truth (dlib) landmarks overlaid. In the second row we show the original images with the landmarks regressed from the ID vector overlaid. The landmarks are qualitatively convincing and clearly reconstruct pose - an entirely non-ID related property. In the third row we show the ground truth (dotted lines) and estimated (solid lines) RGB image histograms.

## 4    Image from ID with a generative model

We now investigate if, and how well, an image of a person can be recovered from an ID descriptor. Having shown that non-ID attributes can be estimated from a face descriptor, we also explore to what extent the original image itself, including non-ID information, can be reconstructed from an ID descriptor.

### 4.1    ID-only inversion

We begin by attempting to create an image that is recognisably the same person as the original but not necessarily similar to the original image. We pose this as an optimisation problem and use a generative face model to constrain the problem. Specifically, we restrict the solution to the space of images represented by the StyleGAN2 face model [16]. We denote by $\mathbf{i} = g(\mathbf{z})$ the face image that
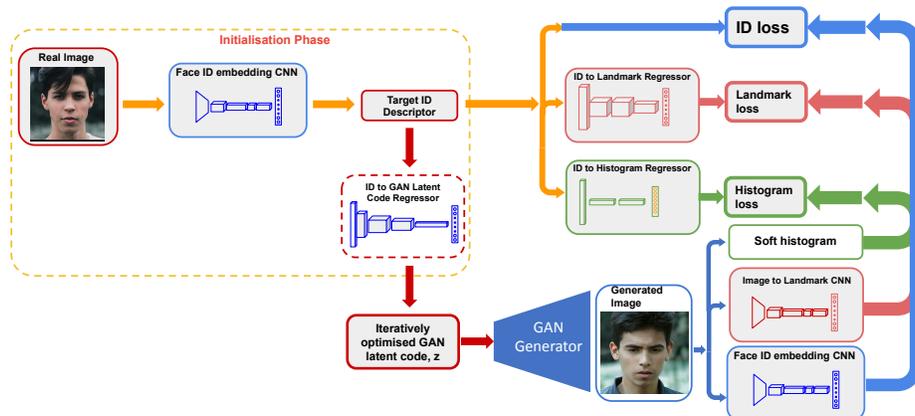
**Fig. 4.** Reconstructing an image from an ID descriptor, including preservation of non-ID properties (landmarks and image histogram). We assume we only have access to the ID descriptor of a real image. We initialise the optimisation using a regression network to predict GAN latent code from ID descriptor. We then iteratively optimise the GAN latent code in order to produce an image that matches the ID, landmarks and histogram predicted from the target ID descriptor using pretrained networks.

arises via the generator from the latent code $\mathbf{z}$. Suppose we are given a target VGG descriptor, $\mathbf{d}$, then we wish to solve the following optimisation problem:

$$\min_{\mathbf{z}} L_{\text{ID}}(\mathbf{z}), \quad \text{where } L_{\text{ID}}(\mathbf{z}) = \|\text{VGG}(g(\mathbf{z})) - \mathbf{d}\|_2^2. \tag{1}$$

In practice, this optimisation is prone to convergence on local minima and sensitive to initialisation. For this reason, we train a network that we use for initialisation that regresses a StyleGAN2 latent code directly from an ID descriptor. Our network comprises an MLP with 3 hidden layers with ReLU activation, 2,048 units per hidden layer. We train this network using synthetic data obtained by randomly sampling images from StyleGAN2, passing these through the face encoder network and then using the resulting ID descriptor and random GAN latent code as an input/output training pair. We subsequently optimise (1) using the Adam optimiser with a learning rate 0.001.

### 4.2 Image reconstruction

The results of the above process successfully produce an image with the correct identity. However, they often fail to reconstruct certain features of the original image, for example the pose and expression of the face, the lighting in the image, the background and the presence of apparel. We have shown that, with suitable supervision and training, it is possible to extract some of these properties from weak signals that find their way into the ID descriptor. Once reconstructed, we now show that these can be used to provide additional, direct supervision to

**Fig. 5.** Direct regression versus ID loss optimisation. Top row: input images. Middle row: output of ID to StyleGAN2 latent code regression network. Bottom row: after subsequent optimisation of StyleGAN2 latent code to minimise $L_{\mathrm{ID}}$.

the inversion problem. Essentially, we ask that not only the ID be reconstructed but that additional, non-ID, features estimated from the ID descriptor also be reconstructed (specifically landmarks and image histogram).

**Landmarks**     From the target ID descriptor, $\mathbf{d}$, we use the pretrained regression network described in Section 3.3 to compute approximate target landmarks, $f_{\mathrm{ID}\to\mathrm{landmarks}}(\mathbf{d})$. During reconstruction, we compare the target landmarks with those extracted from the current image reconstruction using the pretrained image to landmark regression CNN, $f_{\mathrm{image}\to\mathrm{landmarks}}(\mathbf{i})$:

$$L_{\mathrm{landmarks}}(\mathbf{z}) = \|f_{\mathrm{image}\to\mathrm{landmarks}}(g(\mathbf{z})) - f_{\mathrm{ID}\to\mathrm{landmarks}}(\mathbf{d})\|_2^2. \tag{2}$$

**Soft histogram**     For the histogram reconstruction loss, we follow a similar strategy. We use the pretrained regression network described in Section 3.2 to compute an approximate target histogram, $f_{\mathrm{ID}\to\mathrm{histogram}}(\mathbf{d})$. The exact histogram of the reconstructed image is discrete and therefore not differentiable. For this reason, we use a differentiable soft approximation of the image histogram.

The idea is to use sigmoid to softly assign values to bins. Consider a vector $\mathbf{x} \in \mathbb{R}^M$ of $M$ values. We wish to compute a soft histogram $H(\mathbf{x}) \in \mathbb{R}^N$ which softly assigns all values in $\mathbf{x}$ to $N \in \mathbb{Z}^+$ histogram bins. We specify minimum and maximum values (we use $\min = 0$ and $\max = 255$ for image histograms) and the bin width by $\delta = \frac{\max - \min}{N}$. The $i$th bin centre is given by $c_i = \min + \delta(i - 0.5)$. Then, the value of the $k$th bin in $H$ is:

$$H(\mathbf{x})_k = \sum_{j=1}^{M} f(\mathbf{x}_j - c_k + \delta/2) - f(\mathbf{x}_j - c_k - \delta/2), \tag{3}$$
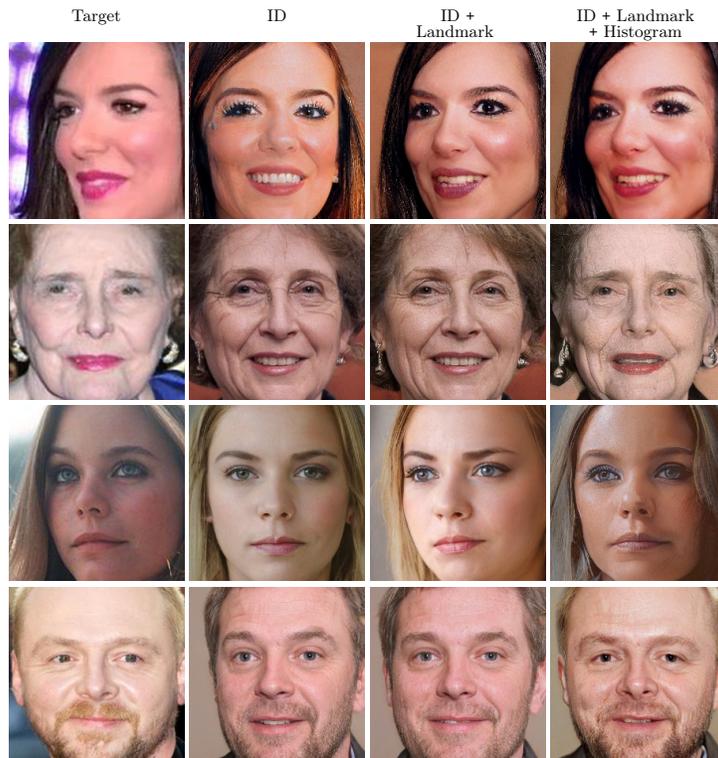
**Fig. 6.** Ablation study. We show inversion results with only ID loss, ID and landmark losses and all three proposed losses.

where $f$ is an assignment function. In a hard (non-differentiable) histogram, $f$ is the Heaviside step function. In our soft histogram, we use sigmoid, $f(x) = \textbf{Sigmoid}(\sigma z)$, with parameter $\sigma$ which controls the softness of the bins. When $\sigma$ is very large, the soft histogram approaches the hard histogram but the gradient vanishes, while small $\sigma$ yields a very soft histogram that badly approximates the true histogram. We use $\sigma = 1.85$ in our experiments. To compute a soft image histogram, we apply (3) to all values in one colour channel of an image, yielding three histograms. Now we can write the histogram loss as: $L_{\text{histogram}}(\mathbf{z}) = \|H(g(\mathbf{z})) - f_{\text{ID}\rightarrow\text{histogram}}(\mathbf{d})\|_2^2$.

**Image reconstruction**     We now pose the image reconstruction problem as optimising the weighted sum of the ID, landmark and histogram losses:

$$\min_{\mathbf{z}} w_1 L_{\text{ID}}(\mathbf{z}) + w_2 L_{\text{landmarks}}(\mathbf{z}) + w_3 L_{\text{histogram}}(\mathbf{z}), \qquad (4)$$

where we use $w_1 = 1$, $w_2 = 0.0006$ and $w_3 = 0.01$ in our experiments. This process is visualised in Fig. 4.

**Fig. 7.** Reconstruction of the same person twice under very different lighting/pose/expression. For each case, we show the original images for the same person (i.e. same ID) in the first row and the corresponding reconstructions in the second row.

### 4.3    Qualitative results

We now present results of inversion from ID to image. We begin with an ablation study of ID-only inversion in Fig. 5. The results show that ID loss optimisation significantly improves over direct regression. The identities in the bottom row are clearly a better visual match to those in the top row. We then follow with an ablation study of our full inversion pipeline in Fig. 6. We show input images in the first column and results with various combinations of losses in columns 2-4. We initialise with our ID to GAN latent code regressor. Then we iteratively optimise only ID loss (column 2 - as in Section 4.1), ID loss and landmark loss (column 3) and all of ID, histogram and landmark losses (column 4). The result in column 2 convincingly reconstructs the ID of the original person but the pose and lighting are wrong. Introducing landmark loss largely corrects the pose (though note StyleGAN2 is biased towards frontal poses which means large pose angles are often underestimated). Introducing histogram loss yields similar lighting and skin tone producing an image similar to the original.

Next, we illustrate that our approach is capable of reconstructing different images of the same person under different conditions. In Fig. 7 we show pairs of real images of the same person in row one. From left to right, these exhibit different lighting, expression and pose. We show our full inversion result in row two. Even though both original images should yield the same ID descriptor, there is enough leaked information that we are able to convincingly reconstruct lighting, expression and pose.

Finally, we show additional inversion results in Fig. 8. The last column shows a failure case in which the pose is incorrectly reconstructed. This occurs when estimated landmark accuracy is low and is further compounded by the StyleGAN2 bias towards frontal faces.

### 4.4    Quantitative results

We quantitatively evaluate the reconstructed images, comparing them to the original images. To facilitate that, we calculate the mean squared error (MSE),

**Fig. 8.** Additional inversion results. We show original target image (top row), reconstructions using only ID loss (middle) and full reconstruction result (bottom). The last column shows a failure case.

peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) [31] between each reconstructed image and the original image, on the MoFA-Test dataset [29], containing 84 images and 78 identities. Table 2 shows an ablation study of reconstructing only with ID loss, with ID and landmark loss, and with all losses. It can be seen that overall, the extra losses help to recreate the actual image and not just the identity of a person.

Second, we test how well the reconstructed images (using all losses) preserve identity on the MoFA-Test dataset. We use cosine similarity on VGGFace [23] (as opposed to VGGFace2, which is used for our inversion) embeddings as a measure of how well a method was able to reconstruct the identity. Fig. 9 shows the distribution of similarity scores of our method, compared with Genova et al. [13], Tran et al. [30], and MoFA [29]. Note that these three methods solve a different problem: reconstruction with a 3D morphable model *given the original image*. However, Genova et al. [13] do this via an ID bottleneck meaning the comparison is meaningful. With an average similarity score of 0.77, we significantly outperform all other methods (0.40 for Genova et al., 0.22 for Tran et al., and 0.18 for MoFA). This is particularly notable given that we reconstruct the image *only from an ID vector*. The difference is likely partly down to using a generative model (StyleGAN2) that is much more powerful than a 3DMM.

## 5   Conclusion

Our results show that, indeed, non-identity information finds its way into state-of-the-art face descriptor embedding networks like VGGFace2 and ArcFace. We have shown that, given the possibility of being able to query the network, it is possible to not only reconstruct an image of a person's face encoded in a
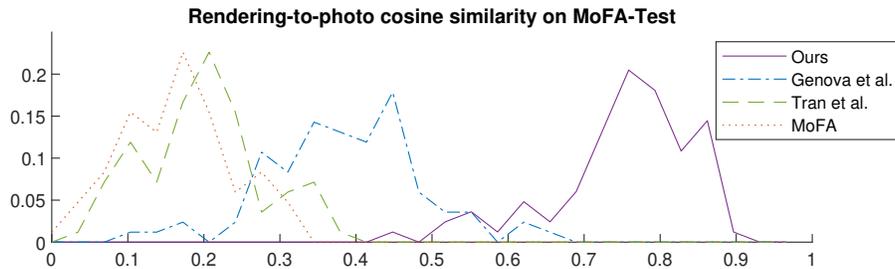
**Rendering-to-photo cosine similarity on MoFA-Test**



**Fig. 9.** Distribution of VGGFace cosine similarity for MoFA-Test. We show the distribution of similarity scores of our method, Genova et al. [13], Tran et al. [30], and MoFA [29] for the original images and their corresponding reconstruction.

|                        | ID only | ID+LMs | ID+LMs+Hist |
|------------------------|---------|--------|-------------|
| MSE (lower better)     | 2.05    | 2.14   | 2.03        |
| PSNR (higher better)   | 12.98   | 13.21  | 13.35       |
| SSIM (higher better)   | 0.14    | 0.14   | 0.15        |

**Table 2.** Quantitative evaluation on MOFA-test, comparing reconstruction with only ID loss, ID and landmark loss, and ID, landmark and histogram loss.

descriptor but also non-ID attributes as well as landmark positions and the image histogram. Being able to reconstruct not just an image with the right ID but the actual original input image has privacy and security implications.

There are many important avenues for future work. First, it is important to replicate these results on other face embedding networks (our initial experiments suggest that our findings indeed transfer between networks). Second, the inversion performance could likely be improved by including other explicit non-ID features. Thirdly, our current work only tried to inversion the background colour by predicting the histogram of RGB intensity from ID, and the subsequent work could be extended by restoring the solid objects on the background to provide a more intuitive view of the background information leaks. Finally, and most interestingly, we believe that our work provides a route to improving face recognition performance while also alleviating privacy concerns. Non-ID information is a nuisance factor for face recognition. It means that some of the capacity of the embedding space is wasted on useless information and that distance measures incorrectly observe identity dissimilarity when in fact the difference is due to non-ID factors. This is in addition to privacy and security concerns related to the leakage of image information into ID vectors. In future we will introduce an adversarial loss into face recognition training that penalises the inclusion of non-ID information in the embedding. We will also attempt to reconstruct non-id information from black-box features, without access to the model's architecture.

# References

1. Abdal, R., Qin, Y., Wonka, P.: Image2stylegan: How to embed images into the stylegan latent space? In: Proceedings of the IEEE international conference on computer vision. pp. 4432–4441 (2019)
2. Abdal, R., Qin, Y., Wonka, P.: Image2stylegan++: How to edit the embedded images? In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8296–8305 (2020)
3. Al-Kuwari, S., Davenport, J.H., Bradford, R.J.: Cryptographic hash functions: Recent design trends and security notions. IACR Cryptol. ePrint Arch. **2011**, 565 (2011)
4. Alvi, M., Zisserman, A., Nellåker, C.: Turning a blind eye: Explicit removal of biases and variation from deep neural network embeddings. In: Proc. European Conference on Computer Vision (ECCV) (2018)
5. Bayat, N., Khazaie, V.R., Mohsenzadeh, Y.: Inverse mapping of face gans. arXiv preprint arXiv:2009.05671 (2020)
6. Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A.: VGGface2: A dataset for recognising faces across pose and age. In: 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018). pp. 67–74. IEEE (2018)
7. Cole, F., Belanger, D., Krishnan, D., Sarna, A., Mosseri, I., Freeman, W.T.: Synthesizing normalized faces from facial identity features. In: CVPR. pp. 3386–3395. IEEE Computer Society (2017)
8. Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4690–4699 (2019)
9. Dhar, P., Gleason, J., Souri, H., Castillo, C.D., Chellappa, R.: Towards gender-neutral face descriptors for mitigating bias in face recognition (2020)
10. Duong, C.N., Truong, T.D., Luu, K., Quach, K.G., Bui, H., Roy, K.: Vec2face: Unveil human faces from their blackbox features in face recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6132–6141 (2020)
11. Fu, Y., Guo, G., Huang, T.S.: Age synthesis and estimation via faces: A survey. IEEE transactions on pattern analysis and machine intelligence **32**(11), 1955–1976 (2010)
12. Gecer, B., Ploumpis, S., Kotsia, I., Zafeiriou, S.: GANFIT: generative adversarial network fitting for high fidelity 3d face reconstruction. In: CVPR. pp. 1155–1164. Computer Vision Foundation / IEEE (2019)
13. Genova, K., Cole, F., Maschinot, A., Sarna, A., Vlasic, D., Freeman, W.T.: Unsupervised training for 3D morphable model regression. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 8377–8386 (2018)
14. Inc., A.: There's more to iPhone. `https://www.apple.com/iphone/more` (2021), [Online; accessed 13-August-2021]
15. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR. pp. 4401–4410. Computer Vision Foundation / IEEE (2019)
16. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: CVPR. pp. 8107–8116. IEEE (2020)
17. Kazemi, V., Sullivan, J.: One millisecond face alignment with an ensemble of regression trees. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1867–1874 (2014)

18. Kumar, N., Berg, A.C., Belhumeur, P.N., Nayar, S.K.: Attribute and simile classifiers for face verification. In: ICCV. pp. 365–372. IEEE Computer Society (2009)
19. Li, A., Guo, J., Yang, H., Chen, Y.: Deepobfuscator: Adversarial training framework for privacy-preserving image classification. arXiv preprint arXiv:1909.04126 (2019)
20. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: Proceedings of the IEEE international conference on computer vision. pp. 3730–3738 (2015)
21. Mahendran, A., Vedaldi, A.: Understanding deep image representations by inverting them. In: CVPR. pp. 5188–5196. IEEE Computer Society (2015)
22. Mirjalili, V., Raschka, S., Namboodiri, A., Ross, A.: Semi-adversarial networks: Convolutional autoencoders for imparting privacy to face images. In: 2018 International Conference on Biometrics (ICB). pp. 82–89. IEEE (2018)
23. Parkhi, O.M., Vedaldi, A., Zisserman, A.: Deep face recognition. In: Proceedings of the British Machine Vision Conference (BMVC). BMVA Press (2015)
24. Razzhigaev, A., Kireev, K., Kaziakhmedov, E., Tursynbek, N., Petiushko, A.: Black-box face recovery from identity features. In: ECCV Workshops (5). Lecture Notes in Computer Science, vol. 12539, pp. 462–475. Springer (2020)
25. Richardson, E., Alaluf, Y., Patashnik, O., Nitzan, Y., Azar, Y., Shapiro, S., Cohen-Or, D.: Encoding in style: a stylegan encoder for image-to-image translation. arXiv preprint arXiv:2008.00951 (2020)
26. Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 815–823 (2015)
27. Shen, Y., Gu, J., Tang, X., Zhou, B.: Interpreting the latent space of GANs for semantic face editing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9243–9252 (2020)
28. Tewari, A., Elgharib, M., Bharaj, G., Bernard, F., Seidel, H.P., Pérez, P., Zöllhofer, M., Theobalt, C.: StyleRig: Rigging StyleGAN for 3d control over portrait images. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
29. Tewari, A., Zollhöfer, M., Kim, H., Garrido, P., Bernard, F., Pérez, P., Theobalt, C.: Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In: IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017. pp. 3735–3744. IEEE Computer Society (2017)
30. Tran, A.T., Hassner, T., Masi, I., Medioni, G.G.: Regressing robust and discriminative 3d morphable models with a very deep neural network. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017. pp. 1493–1502. IEEE Computer Society (2017)
31. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing **13**(4), 600–612 (2004)
32. Wei, T., Chen, D., Zhou, W., Liao, J., Zhang, W., Yuan, L., Hua, G., Yu, N.: A simple baseline for stylegan inversion. CoRR **abs/2104.07661** (2021)
33. Yin, H., Molchanov, P., Alvarez, J.M., Li, Z., Mallya, A., Hoiem, D., Jha, N.K., Kautz, J.: Dreaming to distill: Data-free knowledge transfer via deepinversion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8715–8724 (2020)
34. Zhu, J., Shen, Y., Zhao, D., Zhou, B.: In-domain gan inversion for real image editing. arXiv preprint arXiv:2004.00049 (2020)