

Unsupervised Face Recognition using Unlabeled Synthetic Data

Fadi Boutros^{1,2}, Marcel Klemt¹, Meiling Fang^{1,2}, Arjan Kuijper^{1,2} and Naser Damer^{1,2}

¹Fraunhofer Institute for Computer Graphics Research IGD, Darmstadt, Germany

²Department of Computer Science, TU Darmstadt, Darmstadt, Germany

Email: fadi.boutros@igd.fraunhofer.de

Abstract—Over the past years, the main research innovations in face recognition focused on training deep neural networks on large-scale identity-labeled datasets using variations of multi-class classification losses. However, many of these datasets are retreated by their creators due to increased privacy and ethical concerns. Very recently, privacy-friendly synthetic data has been proposed as an alternative to privacy-sensitive authentic data to comply with privacy regulations and to ensure the continuity of face recognition research. In this paper, we propose an unsupervised face recognition model based on unlabeled synthetic data (USynthFace). Our proposed USynthFace learns to maximize the similarity between two augmented images of the same synthetic instance. We enable this by a large set of geometric and color transformations in addition to GAN-based augmentation that contributes to the USynthFace model training. We also conduct numerous empirical studies on different components of our USynthFace. With the proposed set of augmentation operations, we proved the effectiveness of our USynthFace in achieving relatively high recognition accuracies using unlabeled synthetic data. The training code and pretrained model are publicly available under <https://github.com/fdbtrs/Unsupervised-Face-Recognition-using-Unlabeled-Synthetic-Data>.

I. INTRODUCTION

The evolution in deep learning network architectures, training losses, and availability of large-scale identity-labeled training datasets are behind the major advances in recognition accuracy by the recent state-of-the-art (SOTA) face recognition (FR) models. The main FR works focus on proposing novel FR training losses, especially margin-penalty based softmax loss e.g. ArcFace [14], CurricularFace [23] or ElasticFace [4], to train deep neural network [19] on large-scale identity-labeled dataset [37], [17], [8]. This research direction is driven by the availability of large-scale identity-labeled datasets and the high recognition performance achieved by margin-penalty softmax losses. Recently, there were an increase concerns about collecting, maintaining, redistributing and using biometric data due to legal and ethical privacy issues in some countries [26], [31]. Especially that many of FR datasets such as VGGFace2 [8] have been collected from the web without the proper consent of subjects. Privacy regulations such as the General Data

Protection Regulation (GDPR) [31] classify biometric data as personal data. They grant the right to individuals to withdraw their consent to use or store their personal data. Practically, maintaining such regulations is challenging, especially in the case that such data is collected from the web and is widely distributed.

To overcome this challenge, the use of privacy-friendly synthetic data as an alternative to authentic data in biometrics development has recently attracted attention [6], [5], [13], [12]. In the field of FR, two main previous works proposed the use of synthetically generated data by Generative Adversarial Network (GAN) [16] to develop FR models. SynFace [29] investigated the different behavior of FR models trained on authentic and synthetic images and proposed identity and domain mixup to reduce the performance gap of FR models trained on synthetic data in comparison to FR models trained on authentic data. SFace [7] proposed the use of synthetic data to develop FR models, including presenting a public synthetic database, FR training protocols, detailed analyses of the identity transfer from generator training to the generated data, the identifiability of the authentic data in the trained models. SynFace [29] and SFace [7] mainly focused on using synthetic data to train supervised FR to learn multi-class classification problems.

This work presents contributions toward developing unsupervised FR using privacy-friendly synthetic data (USynthFace). Unlike previous works that require synthetic labeled data [7] or mixing up authentic with synthetic data [29], our proposed framework does not require labeled data or involve authentic data in the FR model training. Thus, it takes full advantage of privacy-friendly synthetic data and does not require a special GAN model to generate labeled data. This work is the first to propose unsupervised FR using synthetic data. Our unsupervised FR training paradigm is based on the concept of the Momentum Contrast [18] and contrastive learning [33] for unsupervised representation learning. The main idea of our approach is to extract a pair of feature representations from two augmented views of the same instance. Then, learning to enhance the similarity between this pair to be higher than the similarity to any other instance. As such learning paradigm mainly depends on augmenting the training sample, we propose a large set of augmentation operations based on geometric and color transformations, as well as controlled GAN augmentation to simulate different realistic appearance variations, i.e. pose, illumination, and expression. We also provide sensitivity studies on all the

This research work has been funded by the German Federal Ministry of Education and Research and the Hessen State Ministry for Higher Education, Research and the Arts within their joint support of the National Research Center for Applied Cybersecurity ATHENE. This work has been partially funded by the German Federal Ministry of Education and Research (BMBF) through the Software Campus Project.

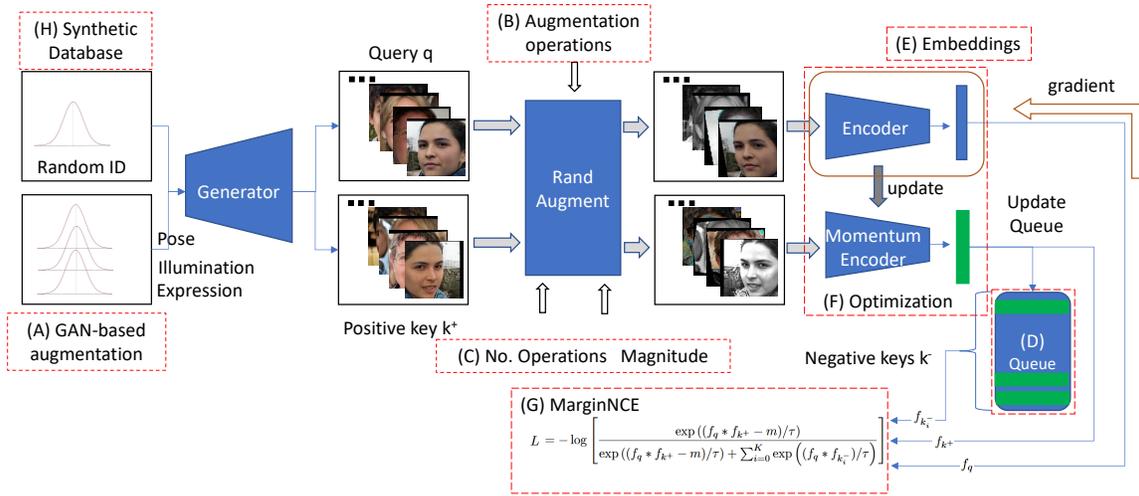


Fig. 1: An overview of our unsupervised FR training paradigm, USynthFace. Synthetic samples are generated and augmented by a generative model to output a query (q) and a positive key k^+ . Then, these samples are augmented with color and geometric transformation through RandAugment. The augmented q and k^+ are processed by the encoder and the momentum encoder, respectively. Then, the queue is updated where the current mini-batch is enqueued and the oldest mini-batch from the queue is dequeued. Finally, the contrastive loss is calculated based on q , k^+ and the negative keys (k^-). Gradients are only back-propagated through the encoder. To maintain key consistency in the queue and to avoid rapidly updating the key encoder, updating the key encoder (momentum encoder) is driven by momentum update with the query encoder.

components of our unsupervised FR framework. USynthFace achieved relatively high verification performances on several benchmarks. Our USynthFace also achieved very competitive results to supervised FR trained on synthetic data and outperformed them on several benchmarks. For example, our USynthFace achieved SOTA accuracy (92.23%) on Labeled Face on the Wild (LFW) for FR trained on synthetic data.

II. METHODOLOGY

This section presents our unsupervised FR training framework and its components based on synthetic data (USynthFace). Figure 1 illustrates the pipeline of our proposed framework. Synthetic face images are conditionally generated using conditional GAN with different random identities, pose, illumination and expression to output a query (q) and a positive key (k^+), which are two instances sampled using same random identity latent vector with different pose, illumination and expression. Then, the query and the positive key are further augmented with geometric and color transformations. The query images are then processed by the encoder and the positive keys are processed by the momentum encoder. The resulting feature representations of the momentum encoder are then pushed into the queue and the ones from the oldest batch are dequeued. Finally, the feature representations of the query, the positive key and the negative keys (retrieved from the queue) are utilized to calculate the contrastive loss.

A. Unsupervised face recognition

a) Unsupervised representation learning: We utilize in this work the concept of the Momentum Contrast (MoCo)[18] for unsupervised representation learning. MoCo

uses contrastive learning to maximize the similarity between feature representations of positive pairs and minimize the similarity between feature representations of negative pairs. Positive pairs are two versions of the same instance augmented using geometric or color transformations, while pairing with any other images (of different instances) is considered as negative pair, which might be perceived as self-supervised learning. MoCo introduced a dynamic queue to form a larger amount of negative pairs than forming negative pairs only from the current batch [36], [21]. Moreover, unlike memory bank where the feature encodings are produced from different training stages [34], MoCo introduced momentum encoder to maintain consistent feature representations.

Consider a query image q encoded into f_q , a positive key of the same instance of q , augmented as k^+ and encoded into f_{k^+} along with a set of negative keys $\{k_i^-\}_{i=1}^K$ (encoded into $\{f_{k_i^-}\}_{i=1}^K$) that are retrieved from the queue. A contrastive loss that guides the model to enhance the similarity between f_q and f_{k^+} to be larger than the similarity between f_q and $\{f_{k_i^-}\}_{i=1}^K$ can be measured using MarginNCE [35] as follows:

$$L = -\log \frac{\exp((f_q * f_{k^+} - m)/\tau)}{\exp((f_q * f_{k^+} - m)/\tau) + \sum_{i=1}^K \exp((f_q * f_{k_i^-})/\tau)} \quad (1)$$

where τ is a temperature hyper-parameter that controls the entropy of the distribution [20] and m is a margin penalty used to encourage the model to learn discriminative feature representations. A detailed sensitivity study on the optimal margin selection is provided in Section IV-I. Once the loss function is calculated, the loss gradients are only back-propagated through the encoder (query encoder θ_{enc}). To avoid rapidly updating the key encoder which might break the queue consistency [18], the weights of the momentum en-

coder (key encoder θ_{mom_enc}) are slowly updated by evolving the query encoder [18] with a momentum coefficient, as follows: $\theta_{mom_enc} \leftarrow mc * \theta_{mom_enc} + (1 - mc) * \theta_{enc}$, where $mc \in [0, 1]$ is a momentum coefficient.

Algorithm 1 USynthFace training pipeline

```

 $Z_{id} \leftarrow$  sample  $I$  vectors from  $N(0, 1)$ 
 $RA \leftarrow$  RandAugment( $N, M$ )
while  $e < num\_epochs$  do
  shuffle  $Z_{id}$ 
  for all  $z_{id}$  in  $Z_{id}$  do
    for  $i$  in  $[0, 1]$  do
       $z_{pose}, z_{expr}, z_{illu} \sim N(0, 1)$ 
       $z_{app}^{(i)} \leftarrow z_{pose} \parallel z_{expr} \parallel z_{illu}$ 
    end for
     $q \leftarrow G(z_{id} \parallel z_{app}^{(0)})$ 
     $k^+ \leftarrow G(z_{id} \parallel z_{app}^{(1)})$ 
     $q \leftarrow RA(q)$ 
     $k^+ \leftarrow RA(k^+)$ 
     $f_q \leftarrow enc(q)$ 
     $\theta_{mom\_enc} \leftarrow mc * \theta_{mom\_enc} + (1 - mc) * \theta_{enc}$ 
     $f_{k^+} \leftarrow mom\_enc(k^+)$ 
     $queue \leftarrow update(f_{k^+}, queue)$ 
     $f_{k_i^-} \leftarrow queue$ 
     $l \leftarrow -\log \frac{\exp((f_q * f_{k^+} - m) / \tau)}{\exp((f_q * f_{k^+} - m) / \tau) + \sum_{i=1}^K \exp((f_q * f_{k_i^-} - m) / \tau)}$ 
    backward(enc, l)
  end for
end while

```



Fig. 2: Samples of augmented images. The first row shows synthetically generated images with a same identity latent vector and the augmented version of these images with RandAugment are presented in the second row.

b) Synthetic data generation: The training dataset of our unsupervised model is synthetically generated using Generative Adversarial Network (GAN) [16]. Specifically, we used DiscoFaceGAN (DFG) [15] to conditionally generate I images with different identity, pose, illumination, and expression. DFG presented a 3D morphable face model (3DMM) [3] to the StyleGAN model [24], enabling disentanglement of identity, pose, expression and illumination in the latent space to conditionally generate realistic images with varying attributes. The conditional image generation serves as data augmentation for our unsupervised learning model as presented in the next section.

c) Data augmentation: Unsupervised representation learning approaches such as MoCo [18], AMDIM [2], and SimCLR [9] are heavily dependant on data augmentation to construct positive pairs of the same instance. Previous approaches [9], [2], [18] utilized geometric and color transformation for augmenting the training data. In this work,

we propose to enrich the conventional data augmentation operations, i.e. geometric and color transformations with GAN-based augmentations generated by a conditional generative model. The conventional data augmentation method is based on RandAugment [11]. The search space of RandAugment has two hyperparameters N and M , where N is the number of transformations applied sequentially to each input image and M is the magnitude of each transformation. Transformation operations include: Horizontal-flipping, Rotate, Translate-x, Translate-y, Shear-x, Shear-y, Sharpness, AutoContrast, Contrast, Solarize, Posterize, Equalize, Color, Brightness, ResizedCrop and Grayscale. Sensitive studies on the effect of each operation on FR verification performance and the selection of RandAugment optimal hyper-parameters are provided in Sections IV-D and IV-E, respectively. To simulate more variations that occur in real images, we also utilize DFG [15] to augment training images with different pose, facial expression, and illumination. To generate such images, we first randomly sample 128-D vector from a normal Gaussian distribution $N(0, 1)$, which represents the identity information [15]. The expression, illumination and pose attributes are controlled by three latent vectors (32-D, 16-D and 3-D, respectively) and are disentangled from the identity latent vector. Two augmented views of the same image (and thus identity) can be generated by fixing the identity latent and randomly modifying the attribute latent vectors. Formally, two augmented images, query q and positive key k^+ , can be generated by sampling two appearance vectors as follows:

$$z_{app}^{(0)} = z_{pose} \parallel z_{expr} \parallel z_{ill}, \{z_{pose}, z_{expr}, z_{ill}\} \sim N(0, 1) \quad (2)$$

and

$$z_{app}^{(1)} = z_{pose} \parallel z_{expr} \parallel z_{ill}, \{z_{pose}, z_{expr}, z_{ill}\} \sim N(0, 1). \quad (3)$$

Each of $z_{app}^{(0)}$ and $z_{app}^{(1)}$ are then concatenated with identity latent vector $z_{id} \sim N(0, 1)$ (randomly sampled) to generate augmented q and k^+ , as follows:

$$q = G(z_{id} \parallel z_{app}^{(0)}) \quad (4)$$

and

$$k^+ = G(z_{id} \parallel z_{app}^{(1)}). \quad (5)$$

Augmentation	EER	FMR10	FMR100	FMR1000
GAN-based	0.0110	0.0019	0.0118	0.0558
RandAugment	0.0967	0.0955	0.1520	0.1915
GAN-based + RandAugment	0.1650	0.2038	0.3547	0.4681

TABLE I: The effect of augmentation on identity in the images indicated by the verification performances as EER, FMR10, FMR100 and FMR1000 on three constructed datasets using GAN-based, RandAugment, and GAN-based with RandAugment augmentations. GAN-based with RandAugment augmentation results in bigger effects on identity, providing more challenging samples for the FR training.

III. EXPERIMENTAL SETUPS

This section presents the experimental settings followed in this paper.

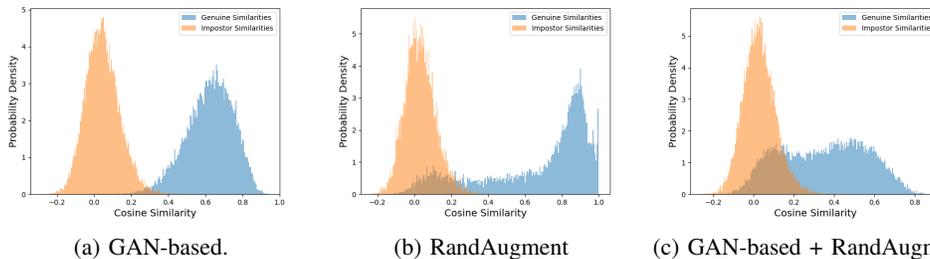


Fig. 3: The genuine (blue) and imposter (orange) score distributions of three different data augmentation settings. The genuine pairs are constructed using GAN-based augmentation (Figure 3a), RandAugment (Figure 3b), and GAN-based with RandAugment (Figure 3c). The biggest effect is noticed when combining both augmentations.

A. Dataset

a) Training dataset: We employ a pretrained DFG model to synthetically generate facial images as discussed in Section II-A.0.b. The model is trained on Flickr-Faces-HQ dataset (FFHQ) [24] that contains 70k images of the size 1024×1024 pixels collected from Flickr and encompass variation in ethnicity, age, image background, and accessories [24]. We opt to generate 100K images from the DFG model, each from different identity latent representations. During the training phase, we augmented these images with conventional augmentation transformations as well as with GAN-based augmentation i.e. pose, illumination and expression (as detailed in Section II-A.0.c). All training data are aligned and cropped to 112×112 using similarity transformation [14], [4], [23] based on detected facial landmarks by Multi-task Cascaded Convolutional Networks (MTCNN) [38]. All images are then normalized to have pixel values between -1 and 1.

b) Evaluation datasets: We used in this paper the following datasets as evaluation benchmarks for our ablation studies: Labeled Faces in the Wild (LFW) [22], AgeDB-30 [27], Celebrities in Frontal to Profile in the Wild (CFP-FP) [30], Cross-Age LFW (CA-LFW) [40], and Cross-Pose LFW (CP-LFW) [39]. The verification accuracy is reported for each of the considered benchmarks following their defined protocols. In all ablation studies in this paper, the overall verification performance is based on the sum of the performance ranking Borda count on LFW, AgeDB-30, CFP-FP, CA-LFW and CP-LFW.

Model	R-R (%)		R-S (%)	
	>FMR100_Th	>FMR1000_Th	>FMR100_Th	>FMR1000_Th
ArcFace	2.6857	0.5664	3.1015	0.5827
CurricularFace	1.9137	0.4024	2.0284	0.3741
ElasticFace	2.0538	0.2951	2.3518	0.3130

TABLE II: Percentage of comparison scores that are larger than different operation thresholds at FMR100 and FMR1000 for R-R and R-S imposter comparison. The percentage number indicates how many comparisons are falsely matched as genuine. The low percentage for the R-S setting and its similarity to R-R indicates that the identities of the authentic data is not linked to these of the synthetic data.

B. Model training setup

The network architecture of the encoder model is ResNet-50 [19], which is one of the widely used architectures in

recent SOTA FR [14], [1], [23], [4]. Following [18], the momentum encoder is updated with a momentum coefficient of 0.999 [18] and the temperature value τ of contrastive loss is set to 0.07 [18]. The feature representation dimensions is initially set to 512-D in the results presented in Sections IV-D, IV-E, IV-F. Later, in Section IV-G we present an ablation study on the optimal feature representation dimensions of 128, 256, 512, and 1024-D. The queue size is set to 32768 based on sensitivity study presented in Section IV-F. An optimizer Stochastic Gradient Descent (SGD) is used with initial learning rate of 0.1. The momentum is set to 0.9 and the weight decay to $5e-4$. The learning rate is divided by 10 after 8, 16, 24, and 32 epochs. The models presented in Sections IV-D, IV-E, IV-F, IV-G are trained for 40 epochs in total with a batch size of 512 on 100K synthetic images. All models are implemented using PyTorch [28] and trained on two CPU 16 core Intel Xeon Gold 5218 and four NVIDIA Quadro RTX6000 GPUs.

IV. RESULTS

This section presents the achieved results by our proposed USynthFace and its components, including: 1) Studying the effect of different data augmentation operations on identity preservation. 2) Analysing identity-shared information between synthetic data and original generative model authentic training data. 3) Ablation studies of different components of our framework (Figure 1), where we provide extensive experiment evaluations of the components of our framework (marked with red rectangles in Figure 1). 4) Comparison with SOTA synthetic-based FR.

A. To which degree does data augmentation effect identity information?

We evaluated the effect of augmenting face images on the identity in the face image. The two augmented versions of the same image (instance) were considered as a genuine pair and pairing with any other image of different instances is considered as an imposter pair. We used SOTA FR model ElasticFace (ElasticFace-Arc) [4] to extract representation features of our synthetic data. The achieved verification performances are reported as Equal Error Rate (EER), FMR10, FMR100, and FMR1000, which are the lowest false non-match rate (FNMR) for a false match rate (FMR) $\leq 10.0\%$, $\leq 1.0\%$ and $\leq 0.1\%$ respectively, along with plotting the genuine-imposter score distributions. We made the

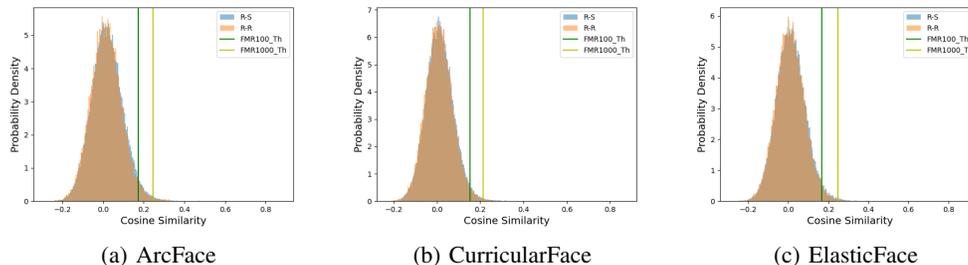


Fig. 4: The score distributions of the two settings, the authentic data R-R and the cross-dataset R-S, achieved by ArcFace [14], CurricularFace [23], and ElasticFace [4]. The highly overlapping score distributions indicate an extremely weak identity relation between the authentic training data and the generated synthetic data.

following observation: 1) GAN-based augmentations (pose, illumination and expression) preserve to large degree the identity information of the augmented sample (0.0110 EER) as shown in Table I and Figure 3a. 2) Color and geometric transformations through RandAugment lead to degradation in verification performance (0.0967 EER) in comparison to GAN-based augmentation. 3) As expected, combining GAN-based with RandAugment achieve the lowest verification performance (0.1650 EER) in comparison to the GAN-based (0.0110 EER) and RandAugment (0.0967 EER). However, combining GAN-based with RandAugment (thus creating challenging genuine pairs) significantly improved our unsupervised model on the considered benchmarks as will be shown in Section IV-E.

Augmentation	LFW	AgeDB-30	CFP-FP	CA-LFW	CP-LFW
HF	73.12	50.95	60.99	60.05	56.13
HF+GAN-based	81.53	53.65	67.21	65.03	64.22

TABLE III: Verification accuracies (%) of two data augmentation settings on five different FR benchmarks. HF refers to horizontal-flipping. Adding GAN-based augmentations enhances the accuracy of the resulting FR model, which will be referred to as "baseline". Higher accuracy in bold.

B. Does the synthetic data share identity information with the GAN authentic training data?

Driven by privacy concerns, we answered this question by conducting an N:N evaluation where references were compared to probes from the GAN authentic training dataset (noted as R-R) and N:M evaluation where authentic references from the GAN training dataset were compared to synthetic probes generated by GAN generator model (noted as R-S). In this experiment, feature representations were obtained from ArcFace [14], CurricularFace [23] and ElasticFace [4], respectively¹. As authentic and synthetic datasets do not have identity label, we calculated the operation thresholds at FMR100 (FMR100_Th) and FMR1000 (FMR1000_Th) for each of the evaluated models on LFW dataset. The comparison scores below the operational threshold were considered as non-match, i.e. of a different identity and the ones that were higher than the operational threshold

¹The network architecture of ElasticFace-Arc [4], ArcFace [14] and CurricularFace [23] is ResNet100 trained on MS1MV2 [17] by the corresponding authors (model publicly available).

were considered as match, i.e. of the same identity. Figure 4 shows score distributions of the cross-dataset (R-S) and authentic data (R-R) with two operational thresholds FMR100_Th and FMR1000_Th. It can be clearly noticed that R-R and R-S score distributions are highly overlapped and only a few samples are considered matched, i.e. achieved comparison scores higher than the operational threshold. This observation is complementary to the previous findings in [32] and [7], as these works also reported that the identity relation between the GAN training dataset and synthetic data is weak.

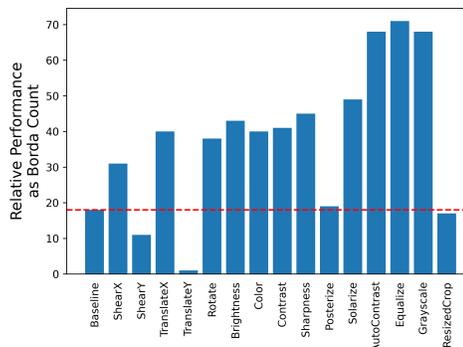


Fig. 5: Augmentation operation selection for USynthFace indicated by the Borda Count (Table IV) as verification performance. 12 out of 15 candidate operations outperformed baseline operations, composing the search space for RandAugment.

C. Impact of GAN-based augmentation

We evaluated the impact of GAN-based augmentation on our USynthFace by training and evaluating USynthFace model with widely used augmentation operation in FR [4], [14], [23], [25], horizontal-flipping. This model is considered as a baseline in this study. Then, we trained a second instance of the baseline model with GAN-based augmentation, i.e. pose, illumination and expression (in addition to horizontal-flipping). The achieved verification performances on the considered evaluation benchmarks are presented in Table III. One can clearly notice that including GAN-based augmentation in the model training significantly improved the verification accuracies in comparison to the baseline model.

D. Impact of conventional data augmentation

In this study, we evaluated the achieved verification performance by introducing different geometric/color transfor-

Augmentation Operation	LFW	AgeDB30	CFP-FP	CALFW	CPLFW	Borda Count
Baseline	81.53	53.65	67.21	65.03	64.22	18
ShearX	81.45	53.63	68.33	65.13	64.43	31
ShearY	80.63	53.12	67.57	63.45	64.28	11
TranslateX	81.82	53.77	69.4	64.78	64.98	40
TranslateY	76.55	52.77	67.39	61.03	62.73	1
Rotate	82.02	54.03	67.91	64.82	65.17	38
Brightness	82.17	54.28	67.64	65.07	65.23	43
Color	81.73	56.12	67.71	66.05	64.93	40
Contrast	82.98	53.30	68.07	64.85	65.4	41
Sharpness	81.92	55.37	68.06	65.55	64.95	45
Posterize	80.47	54.12	67.61	64.48	64.33	19
Solarize	81.87	54.85	69.00	66.23	64.93	49
AutoContrast	84.05	57.88	69.03	67.37	66.73	68
Equalize	85.23	58.57	69.63	67.45	65.42	71
Grayscale	83.05	63.87	68.51	67.78	65.68	68
ResizedCrop	78.92	53.02	68.33	63.35	64.37	17

TABLE IV: Impact of different conventional augmentation operations given as verification accuracies (%) of the trained models. The borda count shows that many augmentations improve beyond the baseline model.

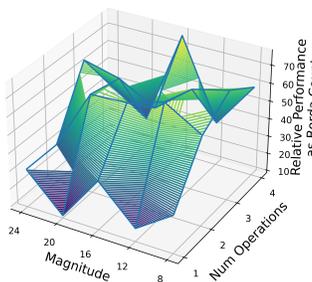


Fig. 6: Hyperparameter of RandAugment selection for USynthFace in terms of Borda Count (Table V) as a verification performance.

mations to model training. The baseline model is noted as "Baseline" and trained with horizontal-flipping and GAN-based augmentations. Table IV and Figure 5 present the achieved results by different models, each was trained with a single candidate augmentation operation (in addition to horizontal-flipping and GAN-based operations). The candidate operation is included in the final augmentation space if it has led to improvement in overall verification performances (in terms of Borda count) in comparison to the baseline model. Out of 15 candidate operations, 12 operations outperformed the baseline operation. These operations are included in the search space of RandAugment.

E. Conventional data augmentation through RandAugment

The augmentation operations from the previous study (Section IV-D) are used to build the search space for RandAugment. We evaluate in this section by randomly augmenting the training samples with multiple operations, i.e. 1, 2, 3 or 4 and with different magnitudes, i.e. 8, 12, 16, 20 or 24. In total, we trained and evaluated 20 models (4 different numbers of operations and five possible magnitudes). The best verification performance is achieved by randomly applying 4 operations (sequentially) with the magnitude of 16 as shown in Table V and Figure 6.

Number of Operations	Magnitude	LFW	AgeDB30	CFP-FP	CALFW	CPLFW	Borda Count
1	8	86.43	63.48	72.97	68.83	68.02	30
	12	86.35	62.23	73.04	68.73	68.13	16
	16	86.50	65.13	73.39	68.37	68.37	36
	20	86.27	62.72	72.20	68.70	67.93	10
2	8	86.48	62.80	73.64	68.72	68.63	30
	12	86.95	63.50	73.73	69.13	68.38	59
	16	87.00	63.33	74.31	69.43	68.75	71
	20	86.93	63.20	74.36	68.88	68.92	65
3	8	86.98	62.98	74.31	69.17	68.77	64
	12	86.55	62.97	74.17	68.52	68.15	30
	16	86.92	63.95	74.01	69.67	68.88	70
	20	86.60	63.20	75.11	68.80	68.88	65
4	8	86.68	63.22	74.21	68.75	68.42	49
	12	86.57	63.82	75.59	68.50	68.77	60
	16	86.57	63.13	74.54	68.88	69.15	65
	20	86.70	62.82	74.19	69.07	69.25	59
5	8	86.53	62.67	74.37	68.47	68.63	33
	12	86.93	64.15	74.51	69.08	68.80	77
	16	86.25	62.72	74.57	68.08	68.02	24
	24	86.00	64.33	74.41	67.52	68.27	37

TABLE V: Sensitivity study on RandAugmen hyperparameters. Verification accuracies (%) of different settings on five FR benchmarks. Randomly applying 4 operations with the magnitude of 16 obtained the best overall performance. The highest Borda Count (accuracy) is in bold.

Queue Size	LFW	AgeDB30	CFP-FP	CALFW	CPLFW	Borda Count
512	86.22	63.43	74.44	67.73	68.53	8
1024	86.55	62.58	73.97	67.87	67.85	7
2048	86.25	63.97	73.67	68.63	68.57	14
4096	86.97	63.25	74.84	68.82	68.72	24
8192	86.50	63.60	75.30	68.30	68.67	22
16384	86.47	62.72	73.67	68.85	69.07	17
32768	86.93	64.15	74.51	69.08	68.80	30
65536	87.02	63.57	74.47	67.87	68.18	18

TABLE VI: Verification accuracies (%) of using different queue sizes on five FR benchmarks. The highest accuracy indicated by the highest Borda Count is in bold.

Feature Dimensionality	LFW	AgeDB-30	CFP-FP	CA-LFW	CP-LFW	Borda Count
128	86.35	63.55	73.93	68.52	68.12	9
256	86.52	63.02	74.23	68.48	68.33	10
512	86.93	64.15	74.51	69.08	68.80	20
1024	86.65	63.52	73.71	68.75	68.15	11

TABLE VII: Verification accuracies (%) achieved by models with different feature representation dimensionality on five FR benchmarks. The best overall verification accuracy is achieved by feature representation of 512-D.

LR-Scheduler	Maximal Epochs	LFW	AgeDB-30	CFP-FP	CA-LFW	CP-LFW
step-based	40	86.93	64.15	74.51	69.08	68.80
plateau-based	200	91.52	69.30	78.46	75.35	71.93

TABLE VIII: Verification accuracies (%) of different LR schedulers on five FR benchmarks. Models trained with plateau-based LR scheduler and more epochs yield better performances than step-based scheduler.

Margin	LFW	AgeDB-30	CFP-FP	CA-LFW	CP-LFW	Borda Count
0.00	91.52	69.30	78.46	75.35	71.93	12
0.05	91.30	70.37	78.73	75.52	71.58	13
0.10	92.12	71.08	78.19	76.15	71.95	22
0.15	91.83	70.78	78.11	76.18	71.50	17
0.20	91.65	70.75	77.80	75.93	71.37	11

TABLE IX: Verification accuracies (%) of different margin values for the MarginNCE loss on five FR benchmarks. A margin value of 0.10 leads to the best overall performance and thus is used in the final experimental setting.

Method	Unsupervised	Identities	Samples per Identity	Total	LFW	AgeDB-30	CFP-FP	CA-LFW	CP-LFW
SynFace [29]	✗	10K	50	500K	88.98	-	-	-	-
SynFace (w/IM) [29]	✗	10K	50	500K	91.97	-	-	-	-
SFace-10 [7]	✗	10,575	10	105K	87.13	63.30	68.84	73.47	66.82
SFace-20 [7]	✗	10,575	20	211K	90.50	69.17	73.33	76.35	71.17
SFace-40 [7]	✗	10,575	40	423K	91.43	69.87	73.10	76.92	73.42
SFace-60 [7]	✗	10,575	60	634K	91.87	71.68	73.86	77.93	73.20
USynthFace (ours)	✓	100K	1	100K	91.52	69.30	78.46	75.35	71.93
USynthFace (ours)	✓	200K	1	200K	91.93	71.23	78.03	76.73	72.27
USynthFace (ours)	✓	400K	1	400K	92.23	71.62	78.56	77.05	72.03

TABLE X: Verification accuracies (%) on five different FR benchmarks achieved by the supervised and SOTA SynFace [29] and SFace [7] models, and our USynthFace model trained on the synthetic training databases of different sizes. The bold number refers to the highest performance on each benchmark. Nothing that the authors of SynFace [29] only provided evaluation results on LFW. Our unsupervised USynthFace model obtained very competitive and even better results than supervised synthetic-based FR models.

F. Analyses of the queue size

In the previous section, we evaluated several augmentation methods for training our USynthFace model. In this section, we study varying the queue size of momentum contrast. Our achieved results in Table VI pointed out that maintaining a queue of 32768 negative keys leads to the highest overall verification performance on the considered evaluation benchmarks. Noting that increasing the queue size to 65536 did not improve the overall verification performances.

G. Study of feature representation dimensionality

Based on the optimal queue size and augmentation methods, we evaluate in this section different learned feature representation dimensionalities i.e. 128, 256, 512 and 1024. All models in this section are trained with GAN-based augmentation and RandAugment with 4 sequential operations and a magnitude of 16 as well as a queue size of 32668. The achieved results of utilizing different feature representation dimensionality are presented in Table VII where the best overall verification performance is achieved using feature representation of 512-D.

H. Training optimization

The presented results are achieved so far by training our USynthFace models using a step-based learning rate schedule. Previous works [10] on unsupervised representation learning pointed out that increasing the number of training epochs is beneficial for improving unsupervised model accuracy. To provide complete evaluation results, we study increasing the number of epochs to a maximum of 200 [10] and using a plateau-based learning scheduler. The initial learning rate is set to 0.1 and it divided by 10 when the average validation accuracy does not improve for 10 consecutive epochs. The training is stopped when the average validation accuracy does not improve for 20 consecutive epochs with maximum of 200 epochs. Using the listed training settings, the model training stopped after 91 training epochs. The achieved results, in this case, are presented in Table VIII, pointing out that increasing the training epochs significantly improved the verification performance on all considered benchmarks.

I. Impact of different margins in MarginNCE

We study in this section different margin values (0, 0.05, 0.1, 0.15 and 0.2) for MarginNCE loss (Eq. 1). The presented results in this section were obtained by training four different models with the optimal observed training settings from the previous experiments. It can be noticed from the achieved results in Table IX that the overall verification performance is improved by increasing margin values from 0 to 0.1. However, when we increase the margin value to 0.15 or 0.20, the overall verification performances are slightly degraded.

J. Study of training database size

Given that there is no restriction on the number of synthetic samples that can be generated using the generative model, we increased the training dataset size from 100K to 200K and to 400K images. Then, we trained two instances of our unsupervised model with the new constructed datasets. The achieved results are presented in Table X together with other SOTA synthetic-based FR models. One can notice that increasing the training dataset sizes to 200K and 400K images slightly improved the overall verification performance in comparison to the model trained with 100K images.

K. Comparison with the SOTA synthetic-based FR

We compared the achieved verification performance by our USynthFace with the recent SOTA FR that proposed the use of synthetic data in FR training (Table X). Noting that this is the first work that proposed to train FR with privacy-friendly synthetic data in an unsupervised fashion. SynFace [29] and SFace [7] are trained with supervised learning to learn multi-class classification using margin-penalty softmax loss. SynFace only reported the verification performance on LFW dataset. On LFW and CFP-FP datasets, our unsupervised model outperformed SFace and SynFace. On AgeDB-30, CA-LFW and CP-LFW, our unsupervised model achieved very competitive results to SFace, even though USynthFace training is unsupervised.

V. CONCLUSION

We presented in this work a novel unsupervised face recognition solution trained on unlabeled synthetic data. The

unsupervised training is based on creating positive pairs to unlabeled synthetic face images of random identities through well-studied augmentations. We proposed not only to use conventional data augmentations in our USynthFace model training, but also introduced GAN-based augmentation to the training pipeline, enhancing the variability in the synthetic face image appearances. This has been complemented with a set of empirical studies on the validity of the different components of our USynthFace and their design choices. With a simple yet effective training paradigm, our USynthFace advanced the SOTA performance on a number of the evaluation benchmarks, in comparison to the recent face recognition models trained on synthetic data, while being the only one trained in an unsupervised manner.

REFERENCES

- [1] X. An, X. Zhu, Y. Gao, Y. Xiao, Y. Zhao, Z. Feng, L. Wu, B. Qin, M. Zhang, D. Zhang, and Y. Fu. Partial FC: training 10 million identities on a single machine. In *ICCVW*, pages 1445–1449. IEEE, 2021.
- [2] P. Bachman, R. D. Hjelm, and W. Buchwalter. Learning representations by maximizing mutual information across views. In *NeurIPS*, pages 15509–15519, 2019.
- [3] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *SIGGRAPH*, pages 187–194. ACM, 1999.
- [4] F. Boutros, N. Damer, F. Kirchbuchner, and A. Kuijper. Elasticface: Elastic margin loss for deep face recognition. In *CVPR Workshops*, pages 1577–1586. IEEE, 2022.
- [5] F. Boutros, N. Damer, and A. Kuijper. Quantface: Towards lightweight face recognition by synthetic data low-bit quantization. In *26th ICPVR 2022, August 21-25, 2022*. IEEE, 2022.
- [6] F. Boutros, N. Damer, K. B. Raja, R. Ramachandra, F. Kirchbuchner, and A. Kuijper. Iris and periocular biometrics for head mounted displays: Segmentation, recognition, and synthetic data generation. *Image Vis. Comput.*, 104:104007, 2020.
- [7] F. Boutros, M. Huber, P. Siebke, T. Rieber, and N. Damer. Sface: Privacy-friendly and accurate face recognition using synthetic data. *CoRR*, abs/2206.10520, 2022.
- [8] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *FG*, pages 67–74. IEEE Computer Society, 2018.
- [9] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 2020.
- [10] X. Chen, H. Fan, R. B. Girshick, and K. He. Improved baselines with momentum contrastive learning. *CoRR*, abs/2003.04297, 2020.
- [11] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPR Workshops*, pages 3008–3017. Computer Vision Foundation / IEEE, 2020.
- [12] N. Damer, F. Boutros, F. Kirchbuchner, and A. Kuijper. D-id-net: Two-stage domain and identity learning for identity-preserving image generation from semantic segmentation. In *ICCV Workshops*, pages 3677–3682. IEEE, 2019.
- [13] N. Damer, C. A. F. López, M. Fang, N. Spiller, M. V. Pham, and F. Boutros. Privacy-friendly synthetic data for the development of face morphing attack detectors. In *CVPR Workshops*, pages 1605–1616. IEEE, 2022.
- [14] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, pages 4690–4699. Computer Vision Foundation / IEEE, 2019.
- [15] Y. Deng, J. Yang, D. Chen, F. Wen, and X. Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *CVPR*, pages 5153–5162. Computer Vision Foundation / IEEE, 2020.
- [16] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [17] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *ECCV (3)*, volume 9907 of *Lecture Notes in Computer Science*, pages 87–102. Springer, 2016.
- [18] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9726–9735. Computer Vision Foundation / IEEE, 2020.
- [19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778. IEEE Computer Society, 2016.
- [20] G. E. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015.
- [21] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio. Learning deep representations by mutual information estimation and maximization. In *ICLR*. OpenReview.net, 2019.
- [22] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst, 10 2007.
- [23] Y. Huang, Y. Wang, Y. Tai, X. Liu, P. Shen, S. Li, J. Li, and F. Huang. Curricularface: Adaptive curriculum learning loss for deep face recognition. In *CVPR*, pages 5900–5909. Computer Vision Foundation / IEEE, 2020.
- [24] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410. Computer Vision Foundation / IEEE, 2019.
- [25] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song. Sphereface: Deep hypersphere embedding for face recognition. In *CVPR*, pages 6738–6746. IEEE Computer Society, 2017.
- [26] B. Meden, P. Rot, P. Terhörst, N. Damer, A. Kuijper, W. J. Scheirer, A. Ross, P. Peer, and V. Struc. Privacy-enhancing face biometrics: A comprehensive survey. *IEEE Trans. Inf. Forensics Secur.*, 16:4147–4183, 2021.
- [27] S. Moschoglou, A. Papaioannou, C. Sagonas, J. Deng, I. Kotsia, and S. Zafeiriou. Agedb: The first manually collected, in-the-wild age database. In *CVPR Workshops*, pages 1997–2005. IEEE Computer Society, 2017.
- [28] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8024–8035, 2019.
- [29] H. Qiu, B. Yu, D. Gong, Z. Li, W. Liu, and D. Tao. Synface: Face recognition with synthetic data. In *ICCV*, pages 10860–10870. IEEE, 2021.
- [30] S. Sengupta, J. Chen, C. D. Castillo, V. M. Patel, R. Chellappa, and D. W. Jacobs. Frontal to profile face verification in the wild. In *WACV*, pages 1–9. IEEE Computer Society, 2016.
- [31] The European Parliament and the Council of the European Union. Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (General Data Protection Regulation), 2016.
- [32] P. J. Tinsley, A. Czajka, and P. J. Flynn. This face does not exist... but it might be yours! identity leakage in generative models. In *WACV*, pages 1319–1327. IEEE, 2021.
- [33] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.
- [34] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, pages 3733–3742. Computer Vision Foundation / IEEE Computer Society, 2018.
- [35] J. Xie, X. Zhan, Z. Liu, Y. Ong, and C. C. Loy. Delving into inter-image invariance for unsupervised visual representations. *CoRR*, abs/2008.11702, 2020.
- [36] M. Ye, X. Zhang, P. C. Yuen, and S. Chang. Unsupervised embedding learning via invariant and spreading instance feature. In *CVPR*, pages 6210–6219. Computer Vision Foundation / IEEE, 2019.
- [37] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *CoRR*, abs/1411.7923, 2014.
- [38] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multi-task cascaded convolutional networks. *CoRR*, abs/1604.02878, 2016.
- [39] T. Zheng and W. Deng. Cross-pose lfw: A database for studying cross-pose face recognition in unconstrained environments. Technical Report 18-01, Beijing University of Posts and Telecommunications, 02 2018.
- [40] T. Zheng, W. Deng, and J. Hu. Cross-age LFW: A database for studying cross-age face recognition in unconstrained environments. *CoRR*, abs/1708.08197, 2017.