# FRCSyn Challenge at WACV 2024:
# Face Recognition Challenge in the Era of Synthetic Data

Pietro Melzi[1]     Ruben Tolosana[1]     Ruben Vera-Rodriguez[1]     Minchul Kim[2]
Christian Rathgeb[3]     Xiaoming Liu[2]     Ivan DeAndres-Tame[1]     Aythami Morales[1]
Julian Fierrez[1]     Javier Ortega-Garcia[1]     Weisong Zhao[4,5]     Xiangyu Zhu[6,7]     Zheyu Yan[6]
Xiao-Yu Zhang[4,5]     Jinlin Wu[8]     Zhen Lei[6,7,8]     Suvidha Tripathi[9]     Mahak Kothari[9]
Md Haider Zama[9]     Debayan Deb[9]     Bernardo Biesseck[10,11]     Pedro Vidal[10]
Roger Granada[12]     Guilherme Fickel[12]     Gustavo Führ[12]     David Menotti[10]
Alexander Unnervik[13,14]     Anjith George[13]     Christophe Ecabert[13]
Hatef Otroshi Shahreza[13,14]     Parsa Rahimi[13,14]     Sébastien Marcel[13,15]     Ioannis Sarridis[16]
Christos Koutlis[16]     Georgia Baltsou[16]     Symeon Papadopoulos[16]     Christos Diou[17]
Nicolò Di Domenico[18]     Guido Borghi[18]     Lorenzo Pellegrini[18]     Enrique Mas-Candela[19]

Ángela Sánchez-Pérez[19]     Andrea Atzori[20]     Fadi Boutros[21,22]     Naser Damer[21,22]
Gianni Fenu[20]     Mirko Marras[20]

[1]Universidad Autonoma de Madrid, Spain [2]Michigan State University, US [3]Hochschule Darmstadt, Germany

[4]IIE, CAS, China [5]School of Cyber Security, UCAS, China [6]MAIS, CASIA, China

[7]School of Artificial Intelligence, UCAS, China [8]CAIR, HKISI, CAS, China [9]LENS, Inc., USA

[10]Federal University of Paraná, Curitiba, PR, Brazil [11]Federal Institute of Mato Grosso, Pontes e Lacerda, Brazil

[12]unico - idTech, Brazil [13]Idiap Research Institute, Switzerland [14]École Polytechnique Fédérale de Lausanne, Switzerland

[15]Université de Lausanne, Switzerland [16]Centre for Research and Technology Hellas, Greece

[17]Harokopio University of Athens, Greece [18]University of Bologna, Cesena Campus, Italy [19]Facephi, Spain

[20]University of Cagliari, Italy [21]Fraunhofer IGD, Germany [22]TU Darmstadt, Germany

## Abstract

*Despite the widespread adoption of face recognition technology around the world, and its remarkable performance on current benchmarks, there are still several challenges that must be covered in more detail. This paper offers an overview of the Face Recognition Challenge in the Era of Synthetic Data (FRCSyn) organized at WACV 2024. This is the first international challenge aiming to explore the use of synthetic data in face recognition to address existing limitations in the technology. Specifically, the FRCSyn Challenge targets concerns related to data privacy issues, demographic biases, generalization to unseen scenarios, and performance limitations in challenging scenarios, including significant age disparities between enrollment and testing, pose variations, and occlusions. The results achieved in the FRCSyn Challenge, together with the proposed benchmark, contribute significantly to the application of synthetic data to improve face recognition technology.*

## 1. Introduction

Facial images represent the most popular data for biometric recognition nowadays, finding extensive applications in surveillance, government offices, and smartphone authentication [29], among others. Numerous studies in the literature have contributed to the development of state-of-the-art (SOTA) Face Recognition (FR) technologies, demonstrating exceptional performance on standard benchmarks [14, 25]. The success of these technologies is attributed to the advent of Deep Learning (DL) and the formulation of highly effective loss functions based on margin loss, capable of generating highly discriminative features [43]. As a result, FR systems have significantly advanced, achieving astonishing results on well-recognized databases, such as LFW [21].

However, FR still encounters numerous challenges due to factors such as variations in facial images concerning pose, aging, expressions, and occlusions, giving rise to significant issues in the field [1, 29, 45]. The application of

(a) DCFace [26].



(b) GANDiffFace [27].

Figure 1. Examples of synthetic identities (one for each row) and intra-class variations for different demographic groups.

DL introduces additional concerns, including limited training data, noisy labeling, imbalanced data related to different identities and demographic groups, and low resolution, among other issues [16]. Deploying FR systems that remain resilient to these challenges and generalize well to unseen conditions is a difficult task. For instance, training data often exhibit significant imbalances across demographic groups [43] and may fail to represent the full spectrum of possible occlusions in real-world scenarios [47]. Various limitations associated with established databases and benchmarks are discussed in [2].

In recent years, several approaches have been presented in the literature for the generation of face synthetic content [3, 15, 49] for different applications such as FR [8, 26, 28] and digital face manipulations, a.k.a. DeepFakes [33, 35, 38]. These synthetic data offer several advantages over real-world databases. Firstly, synthetic databases provide a promising solution to address privacy concerns associated with real data, often collected from individuals without their knowledge or consent through various online sources [32]. Secondly, synthetic face generators have the potential to produce large amounts of data, especially valuable following the discontinuation of established databases due to privacy concerns [19] and the enforcement of regulations like the EU-GDPR, which requires informed consent for collecting and using personal data [39]. Finally, when the synthesis process is controllable, it becomes relatively straightforward to create databases with the desired characteristics (e.g., demographic groups, age, pose, etc.) and their corresponding labels, without additional human efforts. This contrasts with real-world databases, which may not adequately represent diverse demographic groups [30], among many other aspects.

These advantages have motivated an initial exploration of the application of face synthetic data to current FR systems. Innovative generative frameworks have been introduced to synthesize databases suitable for training FR systems, including Generative Adversarial Networks (GANs) [6, 34] and 3D models [3]. While these synthetic databases

advance in the field, some have limitations that impact FR systems performance compared to those trained with real data. Specifically, databases synthesized with GANs provide limited representations of intra-class variations [34], and those synthesized with 3D models lack realism. Recently, Diffusion models have been employed to generate synthetic databases with enhanced intra-class variations, effectively mitigating some limitations observed in prior synthetic databases [26,27]. This is supported by various recent works involving Diffusion models [5,22,49].

To evaluate the effectiveness of novel synthetic databases generated using Diffusion models for training FR systems, this paper analyzes the results achieved in the "Face Recognition Challenge in the Era of Synthetic Data (FRCSyn)" organized at WACV 2024[1]. This challenge is designed to comprehensively analyze the following research questions:

1. Can synthetic data effectively replace real data for training FR systems, and what are the limits of FR technology exclusively trained with synthetic data?

2. Can the utilization of synthetic data be beneficial in addressing and mitigating the existing limitations within FR technology?

In the proposed FRCSyn Challenge, we have designed specific tasks and sub-tasks to address these questions. In addition, we have released to the participants two novel synthetic databases created using two state-of-the-art Diffusion methods: DCFace [26] and GANDiffFace [27]. These databases have been generated with a particular focus on tackling common challenges in FR, including imbalanced demographic distributions, pose variation, expression diversity, and the presence of occlusion (see Figure 1).

The proposed FRCSyn Challenge provides valuable insights for the future of FR and the utilization of synthetic data, with a specific emphasis on quantifying the performance gap between training FR systems with real and synthetic data. In addition, the FRCSyn Challenge proposes

---

[1]https://frcsyn.github.io/

| Database | Framework | Use | # Id | # Img/Id |
|---|---|---|---|---|
| DCFace [26] | DCFace | Train | 10K | 50 |
| GANDiffFace [27] | GANDiffFace | Train | 10K | 50 |
| CASIA-WebFace [46] | Real-world | Train | 10.5K | 47 |
| FFHQ [24] | Real-world | Train | 70K | 1 |
| BUPT-BalancedFace [42] | Real-world | Eval | 24K | 45 |
| AgeDB [31] | Real-world | Eval | 570 | 29 |
| CFP-FP [36] | Real-world | Eval | 500 | 14 |
| ROF [17] | Real-world | Eval | 180 | 31 |

Table 1. Details of the databases considered in the FRCSyn Challenge. Id = Identities, Img = Images.

standard benchmarks that are easily reproducible for the research community. The reminder of the paper is organized as follows. Section 2 provides details about the databases considered in the FRCSyn Challenge. In Section 3, we outline the proposed tasks and sub-tasks, the experimental protocol, and metrics used in the challenge. In Section 4, we provide a description of the top-5 FR systems proposed in the FRCSyn Challenge for each sub-task. Section 5 presents the results achieved in the different tasks and sub-tasks of the challenge. Finally, in Section 6, we draw the conclusions from the FRCSyn Challenge and highlight potential future research directions in the field.

## 2. FRCSyn Challenge: Databases

Table 1 provides details of the public databases considered in the FRCSyn Challenge. Participants were instructed to download all necessary databases for the FRCSyn Challenge upon registration. Permission for redistributing these databases was obtained from the owners.

**Synthetic Databases:** For the training of the proposed FR systems, we provide access to two synthetic databases generated using recent frameworks based on Diffusion models:

- **DCFace** [26]. This framework comprises: *i)* a sampling stage for generating synthetic identities $X_{ID}$, and *ii)* a mixing stage for generating images $X_{ID,sty}$ with the same identities $X_{ID}$ from the sampling stage and styles selected from a "style bank" of images $X_{sty}$.

- **GANDiffFace** [27]. This framework combines GANs and Diffusion models to generate fully-synthetic FR databases with desired properties such as human face realism, controllable demographic distributions, and realistic intra-class variations.

Figure 1 provides examples of the synthetic face images created using DCFace and GANDiffFace approaches. These synthetic databases represent a diverse range of demographic groups, including variations in ethnicity, gender, and age. The synthesis process considers typical variations in FR, including pose, facial expression, illumination, and

occlusions. In the FRCSyn Challenge, synthetic data are exclusively utilized in the training stage, replicating realistic operational scenarios.

**Real Databases:** For the training of FR systems (depending on the sub-task, please see Section 3), participants are allowed to use two real databases: *i)* **CASIA-WebFace** [46], a database containing $494,414$ face images of $10,575$ real identities collected from the web, and *ii)* **FFHQ** [24], a database designed for face applications, containing $70,000$ high-quality face images with considerable variation in terms of age, ethnicity and image background. These real databases are chosen as they are used to train the generative frameworks of DCFace and GANDiffFace, respectively. This strategy enables a direct comparison between the traditional approach of training FR systems using only real data and the novel approach explored in this challenge, using synthetic data. Despite not being specifically designed for face recognition, the FFHQ database can be considered in the proposed challenge for various purposes, such as training a model for feature extraction and applying domain adaptation, among other possibilities.

For the final evaluation of the proposed FR systems, we consider four real databases: *i)* **BUPT-BalancedFace** [42], *ii)* **AgeDB** [31], *iii)* **CFP-FP** [36], and *iv)* **ROF** [17]. BUPT-BalancedFace [42] is designed to address performance disparities across different ethnic groups. We relabel it according to the FairFace classifier [23], which provides labels for ethnicity and gender. We then consider the eight demographic groups obtained from all possible combinations of four ethnic groups (Asian, Black, Indian, and White) and two genders (Female and Male). We recognize that these groups do not comprehensively represent the entire spectrum of real world ethnic diversity. The selection of these categories, while imperfect, is primarily driven by the need to align with the demographic categorizations used in BUPT-BalancedFace [42] for facilitating easier and more consistent evaluation. The other three databases, *i.e.,* AgeDB [31], CFP-FP [36], and ROF [17], are real-world databases widely employed to benchmark FR systems in terms of age variations, pose variations, and presence of occlusions. It is important to highlight that, as different real databases are considered for training and evaluation, we also intend to analyse the generalization ability of the proposed FR systems.

## 3. FRCSyn Challenge: Setup

### 3.1. Tasks

The FRCSyn Challenge has been hosted on Codalab[2], an open-source framework for running scientific competi-

---

[2] https://codalab.lisn.upsaclay.fr/competitions/15485

| |
|---|
| **Task 1:** synthetic data for **demographic bias mitigation** |
|    Baseline: training only with CASIA-WebFace [46] and FFHQ [24]; |
|    Metrics: accuracy (for each demographic group); |
|    Ranking: average vs SD of accuracy, see Section 3.3 for more details. |
| **Sub-Task 1.1:** training exclusively with **synthetic** databases |
|    Train: DCFace [26] and GANDiffFace [27]; |
|    Eval: BUPT-BalancedFace [42]. |
| **Sub-Task 1.2:** training with **real and synthetic** databases |
|    Train: CASIA-WebFace, FFHQ, DCFace, and GANDiffFace; |
|    Eval: BUPT-BalancedFace. |
| **Task 2:** synthetic data for **overall performance improvement** |
|    Baseline: training only with CASIA-WebFace and FFHQ; |
|    Metrics: accuracy (for each evaluation database); |
|    Ranking: average accuracy. |
| **Sub-Task 2.1:** training exclusively with **synthetic** databases |
|    Train: DCFace and GANDiffFace; |
|    Eval: BUPT-BalancedFace, AgeDB [31], CFP-FP [36], and ROF [17]. |
| **Sub-Task 2.2:** training with **real and synthetic** databases |
|    Train: CASIA-WebFace, FFHQ, DCFace, and GANDiffFace; |
|    Eval: BUPT-BalancedFace, AgeDB, CFP-FP, and ROF. |

Table 2. Tasks and sub-tasks proposed in FRCSyn Challenge with their respective metrics and databases. SD = Standard Deviation.

tions and benchmarks. It aims to explore the application of synthetic data into the training of FR systems, with a specific focus on addressing two critical aspects in current FR technology: *i)* mitigating demographic bias, and *ii)* enhancing overall performance under challenging conditions that include variations in age and pose, the presence of occlusions, and diverse demographic groups. To investigate these two areas, the FRCSyn Challenge considers two distinct tasks, each comprising two sub-tasks. Sub-tasks have been designed to consider different approaches for training FR systems: *i)* utilizing solely synthetic data, and *ii)* involving a combination of real and synthetic data. Consequently, the FRCSyn Challenge comprises a total of four sub-tasks. A summary is provided in Table 2. For each sub-task, we specify the databases allowed for training FR systems. Nevertheless, participants have the flexibility to decide whether and how to utilize each database in the training process.

**Task 1:** The first proposed task explores the use of synthetic data to address demographic biases in FR systems. To evaluate the proposed systems, we create lists of mated and non-mated comparisons derived from individuals in the BUPT-BalancedFace database [42]. We consider the eight demographic groups described in Section 2, obtained from the combination of four ethnic groups with two genders. For non-mated comparisons, we exclusively focus on pairs of individuals belonging to the same demographic group, as these are more relevant than non-mated comparisons between individuals of different demographic groups.

**Task 2:** The second proposed task explores the application of synthetic data to enhance overall performance in FR under challenging conditions. To assess the proposed sys-

tems, we use lists of mated and non-mated comparisons derived from individuals included in the four databases indicated in Section 2, namely BUPT-BalancedFace [42], AgeDB [31], CFP-FP [36], and ROF [17]. Each database allows the evaluation of specific challenging conditions for FR, including diverse demographic groups, aging, pose variations, and presence of occlusions.

## 3.2. Experimental protocol

**Training:** The four sub-tasks proposed in the FRCSyn Challenge are mutually independent. This means that participants have the freedom to participate in any number of sub-tasks of their choice. For each selected sub-task, participants are expected to propose a FR system and train it twice: *i)* using authorized real databases only, *i.e.,* CASIA-WebFace [46] and FFHQ [24], and *ii)* in accordance with the specific requirements of the chosen sub-task, as summarized in Table 2. According to this protocol, participants provide both the *baseline system* and the *proposed system* for the specific sub-task. The baseline system plays a critical role in evaluating the impact of synthetic data on training and serves as a reference point for comparing against the conventional practice of training solely with real databases. To maintain consistency, the baseline FR system, trained exclusively with real data, and the proposed FR system, trained according to the specifications of the selected sub-task, must have the same architecture.

**Evaluation:** In each sub-task, participants are provided with comparison files containing both mated and non-mated comparisons, which are used to evaluate the performance of their proposed FR system. In Task 1 there is a single comparison file containing balanced comparisons of different demographic groups, while in Task 2 there are four comparison files, one for each real database considered. The evaluation process occurs twice for each sub-task to assess: *i)* the baseline system trained exclusively with real databases, and *ii)* the proposed system trained in accordance with the sub-task specifications. For the evaluation of each sub-task, participants must submit through Codalab platform two files per database (one for the baseline and one for the proposed system), including the score and the binary decision (mated/non-mated) for each comparison listed in the comparison files. The organizers retain the right to disqualify participants to uphold the integrity of the evaluation process if anomalous results are detected or if participants fail to adhere to the challenge's rules.

**Restrictions:** Participants have the freedom to choose the FR system for each task, provided that the system's number of Floating Point Operations Per Second (FLOPs) does not exceed 25 GFLOPs. This threshold has been established to facilitate the exploration of innovative architectures and

encourage the use of diverse models while preventing the dominance of excessively large models. Participants are also free to utilize their preferred training modality, with the requirement that only the specified databases are used for training. This means that no additional databases can be employed during the training phase, such as to establish verification thresholds. Generative models cannot be utilized to generate supplementary data. Participants are allowed to use non-face databases for pre-training purposes and employ traditional data augmentation techniques using the authorized training databases.

## 3.3. Metrics

We evaluate FR systems using a protocol based on lists of mated and non-mated comparisons for each sub-task and database. From the binary decisions provided by participants, we calculate verification accuracy. This approach is straightforward and allows participants to choose the preferred threshold for their systems. Additionally, we calculate the gap to real (GAP) [26] as follows: $\text{GAP} = (\text{REAL} - \text{SYN})/\text{SYN}$, with REAL representing the verification accuracy of the baseline system and SYN the verification accuracy of the proposed system, trained with synthetic (or real + synthetic) data. Other metrics such as False Non-Match Rate (FNMR) at different operational points, which are very popular for the analysis of FR systems in real-world applications, can be computed from the scores provided by participants. Comprehensive evaluations of the proposed systems will be conducted in subsequent studies, including FNMRs and metrics for each demographic group and database used for evaluation. Next, we explain how participants are ranked in the different tasks.

**Task 1:** To rank participants and determine the winners of Sub-Tasks 1.1 and 1.2, we closely examine the trade-off between the average (AVG) and standard deviation (SD) of the verification accuracy across the eight demographic groups defined in Section 2. We define the trade-off metric (TO) as follows: $\text{TO} = \text{AVG} - \text{SD}$. This metric corresponds to plotting the average accuracy on the x-axis and the standard deviation on the y-axis in 2D space. We draw multiple 45-degree parallel lines to find the winning team whose performance falls to the far right side of these lines. With this proposed metric, we reward FR systems that achieve good levels of performance and fairness simultaneously, unlike common benchmarks based only on recognition performance. The standard deviation of verification accuracy across demographic groups is a common metric for assessing bias and should be reported by any work addressing demographic bias mitigation.

**Task 2:** To rank participants and determine the winners of Sub-Tasks 2.1 and 2.2, we consider the average verifica-

| Team | Affiliations | Country | Sub-Tasks |
|---|---|---|---|
| CBSR | 4-8 | China | 1.2 - 2.2 |
| LENS | 9 | USA | all |
| BOVIFOCR-UFPR | 10-12 | Brazil | all |
| Idiap | 13-15 | Switzerland | all |
| MeVer | 16,17 | Greece | all |
| BioLab | 18 | Italy | 2.1 |
| Aphi | 19 | Spain | 1.1 - 2.1 |
| UNICA-FRAUN-HOFER IGD | 20-22 | Italy, Germany | 1.2 - 2.2 |

Table 3. Description of the top-5 best teams ordered by the affiliation number. The numbers reported in the column 'affiliations' refer to the ones provided in the title page.

tion accuracy across the four databases used for evaluation, described in Section 2. This approach allows us to evaluate four challenging aspects of FR simultaneously: *i)* pose variations, *ii)* aging, *iii)* presence of occlusions, and *iv)* diverse demographic groups, providing a comprehensive evaluation of FR systems in real operational scenarios.

## 4. FRCSyn Challenge: Description of Systems

The FRCSyn Challenge received significant interest, with 67 international teams correctly registered, comprising research groups from both industry and academia. These teams work in various domains, including FR, generative AI, and other aspects of computer vision, such as demographic fairness and domain adaptation. Finally, we received submissions from 15 teams, receiving all sub-tasks high attention. The submitting teams are geographically distributed, with six teams from Europe, five teams from Asia, and four teams from America. Table 3 provides a general overview of the top-5 best teams, including the sub-tasks they participated. Next, we describe briefly the FR systems proposed for each team.

**CBSR (Sub-Tasks 1.2 and 2.2):** They first trained a recognition model using CASIA-WebFace [46]. They extracted features for images in FFHQ [24] and clustered them using the DBSCAN [18] for pseudo labels. Then, they removed the samples in FFHQ that are similar to CASIA-WebFace with a cosine similarity threshold of 0.6 and merged the two to train a new model $F$. They utilized $F$ to de-overlap DCFace [26] and GANDiffFace [27] from CASIA-WebFace and FFHQ. Subsequently, they conducted the intra-class clustering for all databases using DBSCAN (similarity threshold of 0.3) and removed the samples that were separate from the class center. They merged the cleansed databases and trained IResNet-100 with mask and sunglasses augmentation and AdaFace loss [25]. They trained two recognition models using occlusion augmentation with 10% and 30% probability, respectively. They finally submitted the average similarity prediction of the two

models. The threshold was determined by the 10-fold optimal threshold in the validation set.

They constructed different validation sets for different evaluation tasks. For AgeDB [31], they randomly sampled pairs from the training databases. For CFP-FP [36], they added randomly positioned vertical bar masks to the images to simulate the self-occlusion due to pose. For ROF [17], they detected face landmarks [41] and added mask and sunglasses to images. For BUPT-BalancedFace [42], they randomly sampled pairs from DCFace with GANDiffFace because they have balanced demographic groups. All validation sets consisted of $12,000$ image pairs containing $6,000$ positive pairs and $6,000$ negative pairs. Code available[3].

**LENS (All sub-tasks):** For sub-tasks using only synthetic data (*i.e.,* 1.1 and 2.1), they observed that since the evaluation data are real databases, they needed an approach that makes the architecture robust to domain shifts between synthetic training data and real test data. For the same, they utilized the augmentations and AdaFace loss introduced in [25]. The augmentations like Crop, Photometric jittering, and Low-res scaling from [25] helped to create more robust images similar to the real domain, effectively improving performance. They further enhanced the features by using an ensemble of two models, with different styles of augmenting databases like randomly selecting four from set of Identity, Spatial transformations, Brightness, Color, Contrast, Sharpness, Posterize, Solarize, AutoContrast, Equalize, Grayscale, ResizedCrop augmentations in each iteration, inspired from [5]. The features of the two models were then combined to create a feature set of length $1024$. The same method was repeated for Sub-Tasks 1.2 and 2.2.

After cropping and alignment, they divided their total data in the ratio $80:20$ for training and validation, respectively. For training the baseline model and Sub-Tasks 1.2 and 2.2, they utilized CASIA-WebFace [46] for the real database and skipped FFHQ [24]. They adopted the architecture of ResNet-50 [20] (R50) backbone for all the sub-tasks for its lesser number of parameters and suitability when the size of the databases is not huge. They used AdaFace loss from [25].

**BOVIFOCR-UFPR (All sub-tasks):** Inspired by Zhang *et al.* [48], they reduced bias in Sub-Task 1.1 by creating a multi-task collaborative model composed of two backbones $B(x)$ and $R(e)$, which produced the embeddings $e \in R^{512}$ and $g \in R^{256}$, respectively. This schema forced $B(x)$ to learn less biased features across different ethnic groups. ResNet100 and ResNet18 [20] architectures were used as $B(x)$ and $R(e)$. Each training sample $x_i$ contained two labels $y_i$ (to compute the subject loss $L_S$ [14]), and $w_i$, (to

compute the ethnic group loss $L_E$ [14]). Their total loss was $L_T = \lambda_S L_S + \lambda_E L_E$. In Sub-Task 2.1 they employed ArcFace [14] as their loss function and Resnet100 [20] as the backbone, which is one of the top-performing models for deep FR [11]. They trained the network using the Insight-Face library for 26 epochs. The images used for training were augmented using Random Flip with a probability of $0.5$. They used DCFace [26] as the training database in this sub-task, which provided the most accurate feature vectors on the validation set.

**Idiap (All sub-tasks):** The primary strategy for all tasks and sub-tasks was the fusion of features from two models, chosen for its potential to enhance accuracy and reduce bias. These models compute a mean feature vector via a feature fusion approach and undergo independent training to maximize the differences between them, to improve fusion results. For preprocessing, RetinaFace [13] was used to detect facial landmarks across all evaluation sets, and a similarity transform aligned five key facial points to a standard template before cropping and resizing images to $112 \times 112$ pixels, with pixel values normalized between $[-1, 1]$.

The models were based on iResNet-50 and iResNet-101 architectures. Training utilized specific databases for each track, with the iResNet-101 leveraging CosFace loss [40] and the iResNet-50 using AdaFace loss [25]. Training ran for approximately $60,000$ batches of size 256, with learning rate adjustments at set intervals. Training data underwent further preprocessing, including random cropping and augmentations in resolution, brightness, contrast, and saturation. The final model checkpoint was taken after the last training step. A subset of the training data was used to determine the optimal threshold for maximizing verification accuracy, using a 10-fold cross-validation approach based on a random selection of identities and comparison pairs.

**MeVer (All sub-tasks):** Their proposed system utilized the sub-center ArcFace loss [12] to mitigate noise, which occurs in synthetic training data [9]. Comprising three CNNs, the proposed system adapted various margins within the ArcFace loss [14], aligning with relevant literature, indicating different demographic groups require different margin considerations [44]. Final embeddings were obtained by combining the outputs of three ResNet-50 [20] models each trained with $4$, $5$, and $5$ subcenters and margins of $0.45$, $0.47$, and $0.50$. Prediction involved computing the Euclidean distance between feature vectors, utilizing thresholds of $1.5$ and $1.35$ for tasks involving synthetic-only and mixed synthetic-real training data, respectively. The training procedure involved a batch size of 256, an initial learning rate of $0.1$ that decayed by a factor of 10 at steps 75k, 127.5k, and 165k over 180k total training steps. Optimizing with stochastic gradient descent (SGD), momentum was set

at 0.9, and weight decay at 0.0005. Data preprocessing involved an MTCNN [50], resizing all data to $112 \times 112$, and employing color jittering and random horizontal flip augmentations. Task-wise, both synthetic databases were utilized, while the CASIA-WebFace database was specific to Sub-Tasks 1.2 and 2.2. Validation included 800 synthetic identities and $1,000$ identities from CASIA-WebFace for the tasks involving synthetic-only and mixed synthetic-real databases, respectively. Code available[4].

**BioLab (Sub-Task 2.1):** The model selected for the Sub-Task 2.1 is a customized ResNet-101 [14, 20], which had been trained using the margin-based AdaFace loss [25], whose advantage is its resilience when training data contain low-quality images with unrecognizable faces. According to their assumption, this ensured that the model's performance remained unaffected when exposed to GAN-related visual glitches and artifacts. Their baseline model was trained employing the CASIA-WebFace database [46]. Differently, the proposed model employed both DCFace [26] and GANDiffFace [27]. In both cases they built the validation set by generating couples from the first classes of the training sets, which were excluded from training. They applied data augmentation on the training set. Following [25], the pipeline consisted of random horizontal flips, random crop-and-resize, and random color jittering on saturation and value channels. Each transformation had a probability of 20% of being applied. Finally, the model was optimized with cross entropy loss and SGD with an initial learning rate of 0.05. Learning rate scheduling was employed to improve training stability. For face verification, the dissimilarity between embeddings was measured employing the cosine distance. Its threshold was computed to maximize the accuracy on the validation set (*i.e.,* using a non-overlapping partition of the training databases), following the same idea described in the LFW protocol [21]. Code available[5].

**Aphi (Sub-Tasks 1.1 and 2.1):** In their approach, they used an EfficientNetV2-S [37] architecture to produce a 512-D deep embedding trained with ArcFace [14] loss function. They modified the backbone network by reducing the first layer's stride from 2 to 1 to enhance the preservation of spatial features. The output of the backbone network was projected with a $1 \times 1$ convolutional layer and normalized with batch normalization. These features were flattened and fed into a fully connected layer which produces the deep embedding. The weights of the model were optimized through the SGD algorithm with a momentum of 0.9 and a weight decay of $1e^{-4}$ during 20 epochs and a learning rate starting at 0.1 and decayed through a polynomial

scheduler. The model was trained with the images aligned using a proprietary algorithm, resized to $112 \times 112$, and normalized in the range of $-1$ to 1. To prevent overfitting, they applied data augmentation techniques during training, including Gaussian Blur, Random Scale, Hue-Saturation adjustments, and Horizontal Flip transformations as well as dropout with a rate of 0.2 before the deep embedding projection. To train the baseline model, they made use of CASIA-WebFace [46] and for their proposed model, they employed the synthetic database DCFace [26].

**UNICA-FRAUNHOFER IGD (Sub-Tasks 1.2 and 2.2):** The presented solution utilized ResNet100 [20] as network architecture as it is one of the most widely used architectures in state-of-the-art FR approaches [4]. Training and validation images were aligned and cropped to $112 \times 112$ using five-points landmarks extracted with MTCNN. The network's outputs were 512-D feature representations. The presented solution, submitted to Sub-Tasks 1.2 and 2.2, relies on training the ResNet100 network with CosFace as a loss function with a margin penalty value of 0.35 and a scale parameter of 64 [40]. The model was trained for 40 epochs with a batch size of 512 and an initial learning rate of 0.1. The learning rate was divided by 10 after 10, 22, 30, and 40 epochs. During the training phase the training databases, CASIA-WebFace [46] and DCFace [26], provided by the competition organizers, were merged into one database with a total number of 20.572 identities. During the training phase, an extensive set of data augmentation operations based on RandAugment [7, 10] was applied only to the synthetic samples. The real samples were only augmented with horizontal flipping. Code available[6].

## 5. FRCSyn Challenge: Results

Table 4 presents the rankings for the different sub-tasks considered in the FRCSyn Challenge. In general, the rankings for Sub-Tasks 1.1 and 1.2 (bias mitigation), corresponding to the descending order of TO, closely align with the ascending order of SD (*i.e.,* from less to more biased FR systems). Notably, in Sub-Task 1.1, the top two classified teams, LENS (92.25% TO) and Idiap (91.88% TO), exhibit negative GAP values (-0.74% and -3.80%, respectively), indicating higher accuracy when training the FR system with synthetic data compared to real data. These results highlight the potential of DCFace [26] and GANDiffFace [27] synthetic data to reduce bias in current FR technology. The inclusion of real data in the training process (*i.e.,* Sub-Task 1.2) results in general in a simultaneous increase in AVG and reduction in SD, being the CBSR team the winner with a 95.25% TO (*i.e.,* 3% TO general improvement between Sub-Tasks 1.1 and 1.2). In addition, and as it happens

---

**Sub-Task 1.1 (Bias Mitigation): Synthetic Data**

| Pos. | Team | TO [%] | AVG [%] | SD [%] | GAP [%] |
|---|---|---|---|---|---|
| **1** | **LENS** | **92.25** | **93.54** | **1.28** | **-0.74** |
| 2 | Idiap | 91.88 | 93.41 | 1.53 | -3.80 |
| 3 | BOVIFOCR | 90.51 | 92.35 | 1.84 | 4.23 |
| 4 | MeVer | 87.51 | 89.62 | 2.11 | 5.68 |
| 5 | Aphi | 82.24 | 86.01 | 3.77 | 0.84 |

**Sub-Task 1.2 (Bias Mitigation): Synthetic + Real Data**

| Pos. | Team | TO [%] | AVG [%] | SD [%] | GAP [%] |
|---|---|---|---|---|---|
| **1** | **CBSR** | **95.25** | **96.45** | **1.20** | **-2.10** |
| 2 | LENS | 95.24 | 96.35 | 1.11 | -5.67 |
| 3 | MeVer | 93.87 | 95.44 | 1.56 | -0.78 |
| 4 | BOVIFOCR | 93.15 | 95.04 | 1.89 | 1.28 |
| 5 | UNICA | 91.03 | 94.06 | 3.03 | -10.62 |

**Sub-Task 2.1 (Overall Improvement): Synthetic Data**

| Pos. | Team | AVG [%] | GAP [%] |
|---|---|---|---|
| **1** | **BOVIFOCR** | **90.50** | **2.66** |
| 2 | LENS | 88.18 | 3.75 |
| 3 | Idiap | 86.39 | 6.39 |
| 4 | BioLab | 83.93 | 6.88 |
| 5 | MeVer | 83.45 | 3.20 |

**Sub-Task 2.2 (Overall Improvement): Synthetic + Real Data**

| Pos. | Team | AVG [%] | GAP [%] |
|---|---|---|---|
| **1** | **CBSR** | **94.95** | **-3.69** |
| 2 | LENS | 92.40 | -1.63 |
| 3 | Idiap | 91.74 | 0.00 |
| 4 | BOVIFOCR | 91.34 | 1.77 |
| 5 | MeVer | 87.60 | -1.57 |

Table 4. Ranking for the four sub-tasks, according to the metrics described in Section 3.3. TO = Trade-Off, AVG = Average accuracy, SD = Standard Deviation of accuracy, GAP = Gap to Real.

in Sub-Task 1.1, we can observe in Sub-Task 1.2 negative GAP values for the top teams (*e.g.,* -2.10% and -5.67% for the CBSR and LENS teams, respectively), evidencing that the combination of synthetic and real data (proposed system) outperforms FR systems trained only with real data (baseline system).

For Task 2, it is evident that the average accuracy across databases in Sub-Tasks 2.1 and 2.2 is lower than the accuracy achieved for BUPT-BalancedFace [42] in Sub-Tasks 1.1 and 1.2, emphasizing the additional challenges introduced by the other real databases considered for evaluation. Also, although good results are achieved in Sub-Task 2.1 when training only with synthetic data (90.50% AVG for BOVIFOCR-UFPR), the positive GAP values provided by the top-5 teams indicate that synthetic data alone currently struggles to completely replace real data for training FR systems in challenging conditions. Nevertheless, the negative GAP values provided by the top-2 teams in Sub-Task 2.2 also suggest that synthetic data combining with real data can mitigate existing limitations within FR technology.

Finally, analyzing the contributions of all the eight top teams, a notable trend emerges, showing the prevalence of well-established methodologies. ResNet backbones [20]

were chosen by seven teams, except for Aphi, which opted for EfficientNet [37]. The AdaFace [25] and ArcFace [14] loss functions were widely used, featuring in the approaches of CBSR, LENS, Idiap, and BioLab for the former, and BOVIFOCR-UFPR, MeVer, and Aphi for the latter. Idiap and UNICA-FRAUNHOFER IGD also considered the Cos-Face loss function [40]. Most of the teams integrated multiple networks into their proposed architectures for different objectives, *e.g.*, CBSR and LENS trained different networks with distinct augmentation techniques, while BOVIFOCR-UFPR and Idiap combined different loss functions. Some teams also addressed the challenges of domain shift between synthetic and real data, *e.g.*, LENS proposed solutions robust to domain shifts with consistent data augmentation, while CBSR implemented a range of strategies, including advanced data augmentation, identity clustering, and distinct thresholds for different databases. Notably, CBSR utilized all available databases for training, including FFHQ [24], unlike other teams. Excluding BOVIFOCR-UFPR, Aphi, and UNICA-FRAUNHOFER IGD, which exclusively used DCFace [26], the majority of teams employed both DCFace [26] and GANDiffFace [27], demonstrating the suitability of both generative frameworks.

## 6. Conclusion

The Face Recognition Challenge in the Era of Synthetic Data (FRCSyn) has provided a comprehensive analysis for the application of synthetic data to FR, addressing current limitations in the field. Within this challenge numerous approaches from different research groups have been proposed. These approaches can be compared across a variety of sub-tasks, with many being reproducible thanks to the materials made available by the participating teams. Future works will be oriented to a more detailed analysed of the results, including additional metrics and graphical representations. Furthermore, we are considering transforming the CodaLab platform into an ongoing competition, where new tasks and sub-tasks might be introduced.

## Acknowledgements

# References

[1] Insaf Adjabi, Abdeldjalil Ouahabi, Amir Benzaoui, and Abdelmalik Taleb-Ahmed. Past, present, and future of face recognition: A review. *Electronics*, 9(8):1188, 2020. 1

[2] Waqar Ali, Wenhong Tian, Salah Ud Din, Desire Iradukunda, and Abdullah Aman Khan. Classical and modern face recognition approaches: a complete review. *Multimedia tools and applications*, 80:4825–4880, 2021. 2

[3] Gwangbin Bae, Martin de La Gorce, Tadas Baltrušaitis, Charlie Hewitt, Dong Chen, Julien Valentin, Roberto Cipolla, and Jingjing Shen. DigiFace-1M: 1 Million Digital Face Images for Face Recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023. 2

[4] Fadi Boutros, Naser Damer, Florian Kirchbuchner, and Arjan Kuijper. ElasticFace: Elastic Margin Loss for Deep Face Recognition. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2022. 7

[5] Fadi Boutros, Jonas Henry Grebe, Arjan Kuijper, and Naser Damer. IDiff-Face: Synthetic-based Face Recognition through Fizzy Identity-Conditioned Diffusion Model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 2, 6

[6] Fadi Boutros, Marco Huber, Patrick Siebke, Tim Rieber, and Naser Damer. Sface: Privacy-friendly and accurate face recognition using synthetic data. In *Proceedings of the IEEE International Joint Conference on Biometrics*, 2022. 2

[7] Fadi Boutros, Marcel Klemt, Meiling Fang, Arjan Kuijper, and Naser Damer. Unsupervised face recognition using unlabeled synthetic data. In *Proceedings of the International Conference on Automatic Face and Gesture Recognition*, 2023. 7

[8] Fadi Boutros, Vitomir Struc, Julian Fierrez, and Naser Damer. Synthetic data for face recognition: Current state and future prospects. *Image and Vision Computing*, page 104688, 2023. 2

[9] Jiacheng Cheng, Tongliang Liu, Kotagiri Ramamohanarao, and Dacheng Tao. Learning with bounded instance and label-dependent label noise. In *Proceedings of the International Conference on Machine Learning*, 2020. 6

[10] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. RandAugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition Workshops*, 2020. 7

[11] Jiankang Deng, Jia Guo, Xiang An, Zheng Zhu, and Stefanos Zafeiriou. Masked face recognition challenge: The insightface track report. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 6

[12] Jiankang Deng, Jia Guo, Tongliang Liu, Mingming Gong, and Stefanos Zafeiriou. Sub-center ArcFace: Boosting Face Recognition by Large-scale Noisy Web Faces. In *Proceedings of the European Conference on Computer Vision*, 2020. 6

[13] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. RetinaFace: Single-shot Multi-level Face Localisation in the Wild. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2020. 6

[14] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 1, 6, 7, 8

[15] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2020. 2

[16] Hang Du, Hailin Shi, Dan Zeng, Xiao-Ping Zhang, and Tao Mei. The elements of end-to-end deep face recognition: A survey of recent advances. *ACM Computing Surveys (CSUR)*, 54(10s):1–42, 2022. 2

[17] Mustafa Ekrem Erak$ı$n, Uğur Demir, and Haz$ı$m Kemal Ekenel. On Recognizing Occluded Faces in the Wild. In *Proceedings of the International Conference of the Biometrics Special Interest Group*, 2021. 3, 4, 6

[18] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, volume 96, pages 226–231, 1996. 5

[19] Jules. Harvey, Adam. LaPlace. Exposing.ai. https://exposing.ai, 2021. 2

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016. 6, 7, 8

[21] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled Faces in the Wild: A Database for Studying Face Recognition in Unconstrained Environments. In *Proceedings of the Workshop on Faces in 'Real-Life' Images: Detection, Alignment, and Recognition*, 2008. 1, 7

[22] Manuel Kansy, Anton Raël, Graziana Mignone, Jacek Naruniec, Christopher Schroers, Markus Gross, and Romann M Weber. Controllable Inversion of Black-Box Face Recognition Models via Diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 2

[23] Kimmo Karkkainen and Jungseock Joo. FairFace: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021. 3

[24] Tero Karras, Samuli Laine, and Timo Aila. A Style-Based Generator Architecture for Generative Adversarial Networks. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2019. 3, 4, 5, 6, 8

[25] Minchul Kim, Anil K. Jain, and Xiaoming Liu. AdaFace: Quality Adaptive Margin for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 1, 5, 6, 7, 8

[26] Minchul Kim, Feng Liu, Anil Jain, and Xiaoming Liu. DC-Face: Synthetic Face Generation with Dual Condition Diffusion Model. In *Proceedings of the IEEE/CVF Conference on*

*Computer Vision and Pattern Recognition*, 2023. 2, 3, 4, 5, 6, 7, 8

[27] Pietro Melzi, Christian Rathgeb, Ruben Tolosana, Ruben Vera-Rodriguez, Dominik Lawatsch, Florian Domin, and Maxim Schaubert. GANDiffFace: Controllable Generation of Synthetic Datasets for Face Recognition with Realistic Variations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2023. 2, 3, 4, 5, 7, 8

[28] Pietro Melzi, Christian Rathgeb, Ruben Tolosana, Ruben Vera-Rodriguez, Aythami Morales, Dominik Lawatsch, Florian Domin, and Maxim Schaubert. Synthetic Data for the Mitigation of Demographic Biases in Face Recognition. In *Proceedings of the IEEE International Joint Conference on Biometrics*, 2023. 2

[29] Shervin Minaee, Amirali Abdolrashidi, Hang Su, Mohammed Bennamoun, and David Zhang. Biometrics recognition using deep learning: A survey. *Artificial Intelligence Review*, pages 1–49, 2023. 1

[30] Aythami Morales, Julian Fierrez, Ruben Vera-Rodriguez, and Ruben Tolosana. SensitiveNets: Learning agnostic representations with application to face images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(6):2158–2164, 2020. 2

[31] Stylianos Moschoglou, Athanasios Papaioannou, Christos Sagonas, Jiankang Deng, Irene Kotsia, and Stefanos Zafeiriou. AgeDB: The First Manually Collected, In-The-Wild Age Database. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition Workshops*, 2017. 3, 4, 6

[32] Madhumita Murgia and Max Harlow. Who's using your face? The ugly truth about facial recognition. *Financial Times*, 19, 2019. 2

[33] Joao C Neves, Ruben Tolosana, Ruben Vera-Rodriguez, Vasco Lopes, Hugo Proença, and Julian Fierrez. GANprintR: Improved Fakes and Evaluation of the State of the Art in Face Manipulation Detection. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):1038–1048, 2020. 2

[34] Haibo Qiu, Baosheng Yu, Dihong Gong, Zhifeng Li, Wei Liu, and Dacheng Tao. SynFace: Face Recognition With Synthetic Data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 2

[35] Christian Rathgeb, Ruben Tolosana, Ruben Vera-Rodriguez, and Christoph Busch. *Handbook of digital face manipulation and detection: from DeepFakes to morphing attacks*. Springer Nature, 2022. 2

[36] Soumyadip Sengupta, Jun-Cheng Chen, Carlos Castillo, Vishal M Patel, Rama Chellappa, and David W Jacobs. Frontal to profile face verification in the wild. In *Proceedings of the IEEE Winter conference on Applications of Computer Vision*, 2016. 3, 4, 6

[37] Mingxing Tan and Quoc Le. EfficientNetV2: Smaller Models and Faster Training. In *Proceedings of the International Conference on Machine Learning*, 2021. 7, 8

[38] Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information Fusion*, 64:131–148, 2020. 2

[39] Paul Voigt and Axel Von dem Bussche. The EU General Data Protection Regulation (GDPR). *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, 10(3152676):10–5555, 2017. 2

[40] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. CosFace: Large margin cosine loss for deep face recognition. In *Proceedings IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. 6, 7, 8

[41] Jun Wang, Yinglu Liu, Yibo Hu, Hailin Shi, and Tao Mei. FaceX-Zoo: A PyTorch Toolbox for Face Recognition. In *Proceedings of the 29th ACM International Conference on Multimedia*, 2021. 6

[42] Mei Wang and Weihong Deng. Mitigating bias in face recognition using skewness-aware reinforcement learning. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2020. 3, 4, 6, 8

[43] Mei Wang and Weihong Deng. Deep face recognition: A survey. *Neurocomputing*, 429:215–244, 2021. 1, 2

[44] Mei Wang, Yaobin Zhang, and Weihong Deng. Meta balanced network for fair face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8433–8448, 2021. 6

[45] David Wanyonyi and Turgay Celik. Open-source face recognition frameworks: A review of the landscape. *IEEE Access*, 10:50601–50623, 2022. 1

[46] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 3, 4, 5, 6, 7

[47] Dan Zeng, Raymond Veldhuis, and Luuk Spreeuwers. A survey of face recognition techniques under occlusion. *IET Biometrics*, 10(6):581–606, 2021. 2

[48] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2018. 6

[49] Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry Wu, Dmitry Lagun, Thabo Beeler, and Fernando De la Torre. ITI-GEN: Inclusive Text-to-Image Generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 2

[50] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503, 2016. 7