# GBDF: Gender Balanced DeepFake Dataset Towards Fair DeepFake Detection

Aakash Varma Nadimpalli and Ajita Rattani*

School of Computing
Wichita State University, USA
axnadimpalli@shockers.wichita.edu, ajita.rattani@wichita.edu

**Abstract.** Facial forgery by deepfakes has raised severe societal concerns. Several solutions have been proposed by the vision community to effectively combat the misinformation on the internet via automated deepfake detection systems. Recent studies have demonstrated that facial analysis-based deep learning models can discriminate based on protected attributes. For the commercial adoption and massive roll-out of the deepfake detection technology, it is vital to evaluate and understand the fairness (the absence of any prejudice or favoritism) of deepfake detectors across demographic variations such as gender and race. As the performance differential of deepfake detectors between demographic subgroups would impact millions of people of the deprived sub-group. This paper aims to evaluate the fairness of the deepfake detectors across males and females. However, existing deepfake datasets are not annotated with demographic labels to facilitate fairness analysis. To this aim, we manually annotated existing popular deepfake datasets with gender labels and evaluated the performance differential of current deepfake detectors across gender. Our analysis on the gender-labeled version of the datasets suggests (a) current deepfake datasets have skewed distribution across gender, and (b) commonly adopted deepfake detectors obtain unequal performance across gender with mostly males outperforming females. Finally, we contributed a gender-balanced and annotated deepfake dataset, GBDF, to mitigate the performance differential and to promote research and development towards fairness-aware deep fake detectors. The GBDF dataset is publicly available at: https://github.com/aakash4305/GBDF

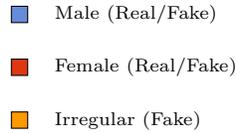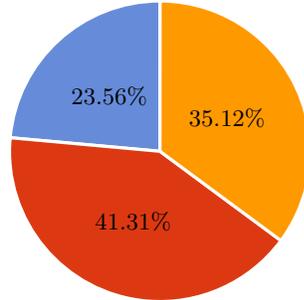**Keywords:** DeepFakes · Fairness and Bias in AI · Facial Analysis.

## 1 Introduction

With the advances in deep generative models, synthetic media have become so realistic that they are often indiscernible from authentic content for human eyes. However, synthetic media generation techniques used by malicious users to deceive pose a severe societal and political threat. In this context, Deepfakes - facial forgery technique that depicts human subjects with altered identities

---
* Corresponding author

or malicious actions using various deep fake generation techniques- has been flagged as a top AI threat [26,11,31,19,33]. Deep fakes have been used to commit fraud, falsify evidence, manipulate public debates, and destabilize political processes [9,31].

**FaceForensics++ Distribution**

**CelebDF Distribution**

23.56%

35.12%

41.31%

30%

70%

■ Male (Real/Fake)

■ Female (Real/Fake)

■ Irregular (Fake)

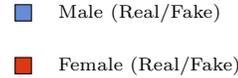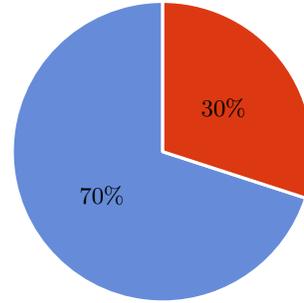■ Male (Real/Fake)

■ Female (Real/Fake)

**Fig. 1.** Illustration of the distribution of videos in Face Forensics++ and Celeb-DF Dataset across gender. The percentage distribution of videos belonging to males (real/fake), females (real/fake) and those classified as irregular swaps is shown.

To mitigate the risk posed by deep fakes, the vision community has developed a series of effective deep fake detection methods [26,11,31] trained on large-scale deepfake datasets. The popular deep fake detection methods include convolutional neural networks (CNN) for detecting visual artifacts [22] and blending boundaries [19], mouth movement analysis [14] and behavioral biometrics [2]. The popular publicly available deep fake datasets include Celeb-DF [21], Face-Forensics++ [28], DeeperForensics-1.0 [16] and DFDC [12] for research and development in this field.

Such efforts have been translated into creating **real-world impact** with Microsoft's release of Video Authenticator[1], an automated tool trained on the publicly available FaceForensics++ deepfake dataset, for detection of artificial manipulation in images and videos. Further, Facebook[2] has been advancing its methods to detect and ban AI-generated profiles, along with strengthening its policy on deepfakes and synthetic media. Recently, the Coalition for Content

---

[1] `https://blogs.microsoft.com/on-the-issues/2020/09/01/`
`disinformation-deepfakes-newsguard-video-authenticator/`

[2] `https://www.wired.com/story/facebook-removes-accounts-ai-generated-photos/`

Provenance and Authenticity (C2PA) has teamed up with Intel, and Adobe to develop new standards targeted at combating the proliferation of deepfakes[3].

While significant advances have been made towards accurate deepfake detection, very little is discussed on the fairness of these deepfake detectors across protected attributes (demographic variations) such as gender and race. Fairness is defined as the absence of any prejudice or favoritism towards an individual or a group based on their inherent or acquired characteristics [5,18]. For the *massive commercial roll-out* of deep fake detection technology, it is vital to examine the bias and fairness of this technology across demographics. This is to avoid any real-world consequences from a biased and flawed system toward a particular sub-group. *As in the common operating scenario, the social media data across gender and race would be audited at the mass level for authenticity via an automated deepfake detection system. Even the small performance differential of deepfake detectors across demographic sub-groups would impact millions of people belonging to the deprived sub-group.*

This draws attention to fairness and bias in AI-based facial analytics where unintended consequences from biased systems call for a thorough examination of the datasets and models [18,5,17,8,4]. Most of the published research in this domain suggests low performance for women, and dark-skinned people for facial attribute-based classification systems such as gender and age [8,17,30,24], and face recognition [5,4]. As biased datasets produce biased models, many of the efforts have been focused on developing gender and race-balanced datasets for various facial-analysis based applications. FairFace [17], a gender and race balanced facial attribute dataset, RFW [34], a racially balanced face recognition dataset, and a gender-balanced dataset developed from existing facial recognition datasets [4] are some of the examples.

This paper aims to examine the **fairness** of deepfake detectors across gender. However, current deepfake detection datasets are not annotated with demographic labels to facilitate the examination of bias. To this aim, deepfake datasets namely FaceForensics++,and Celeb-DF are *manually annotated* with gender labels. The fairness of popular deepfake detectors is evaluated on these datasets across gender. On manual annotation, we found that the gender distribution of the popular deepfake datasets is skewed. The large number of deepfakes in Faceforensics++ are irregular (in conformance with [32])- when a person's face is swapped with the face of another gender or race. This result in the loss of gender-specific information in the fake content. The popular Celeb-DF dataset distribution is heavily skewed towards males (70%).

Figure 1 shows the distribution of videos across gender for the popular Face-Forensics++ [28] and Celeb-DF [21] deepfake datasets. The deepfake detectors evaluated on these skewed datasets along with irregular swaps mostly obtain lower performance for females over males. Finally, we introduced a gender-balanced and annotated deepfake dataset, GBDF, developed from FaceForensics++, Celeb-DF, and DeeperForensics-1.0 and consisting of $10,000$ videos. This balanced dataset aims to mitigate the performance differential of deepfake

---

[3] https://c2pa.org/post/release_1_pr/

detectors due to existing gender unbalanced training sets along with irregular swaps. The dataset information is available to the vision community to promote further research and development in this field. Note that according to ISO/IEC 22116 [7], the term "sex", understood as "the state of being male or female" would be more appropriate instead of "gender" in the context of this study. However, in consistency with the existing studies [8,4], the term gender is used in this paper. To the best of our knowledge, the only study in [32] evaluates the bias of three popular CNN-based deepfake detectors trained on Faceforensics++ across gender and race. The test bed was created using UTKFace and RFW datasets and the deepfakes were generated using the Face X-ray model. The authors reported performance differences for dark-skinned people and emphasized the importance of benchmark representation and auditing for increased demographic transparency.

The main **contributions** of the paper are as follows:

1. Gender label annotation of the popular deepfake datasets namely, FaceForensics++ and Celeb-DF to facilitate analysis of the dataset distribution across gender and the presence of irregular swaps.
2. Evaluation of the fairness of popular deepfake detection algorithms varying in size, architecture, and the methodology, trained and tested on gender annotated versions of the existing datasets.
3. Development of publicly available gender-balanced and annotated deepfake dataset, GBDF, from FaceForensics++ (FF++), Celeb-DF, and Deeper Forensics-1.0 consisting of $10,000$ live and fake videos generated using different identity and expression swapping deepfake generation techniques.
4. Cross-comparison of the performance differential of deepfake detectors trained on existing and our gender-balanced GBDF training set, across males and females.

This paper is organized as follows: Section 2 discuss the related work on deepfake detectors and gendered differences in facial analytics. Section 3 discusses the development of the GBDF dataset. Deepfake detection algorithms used in this study are discussed in Section 4. Evaluation metrics used for fairness analysis are discussed in Section 5. Results and discussion is detailed in Section 6. Conclusion and future research directions are discussed in Section 7.

## 2   Related Work

### 2.1   Deepfake Detection

In this section, we will discuss the existing countermeasure proposed for deep fake detection. Most of the existing methods are CNN-based classification baselines trained for deep fake detection [15,10,25,27].

In [20], Li and Lyu used VGG16, ResNet50, ResNet101, and ResNet152 based CNNs for the detection of the presence of artifacts from the facial regions and the surrounding areas for deep fake detection. Afchar et al. [1] proposed two

different CNN architectures composed of only a few layers in order to focus on the mesoscopic properties of the images: (a) a CNN comprised of 4 convolutional layers followed by a fully-connected layer (Meso-4), and (b) a modification of Meso-4 using a variant of the Inception module named MesoInception-4. In [28], an exhaustive analysis of different CNN-based deep fake detection methods by Rosslet et al. suggested efficacy of XceptionNet when evaluated on FaceForensics++. In [19], a face X-ray model has been proposed to detect forgery by detecting the blending boundary of a forged image using a two-class CNN model trained end-to-end.

Apart from the aforementioned CNN-based deep fake detection methods, spatial temporal information using Long Short-term Memory (LSTM) networks [6], facial and behavioral biometrics (i.e., facial expression, head, and body movement), and lipforensics [14] have been used for deep fake detection [13,2,3,27]. In [14], LipForensics that targets high-level semantic irregularities in mouth movements common in many generated deepfake videos, is used for deepfake detection. Studies have also been proposed for improving the performance of deepfake detectors across datasets and deep fake generation methods using techniques such as reinforcement learning [23] and fine-grained multi-attention network [36]. Readers are referred to the published survey in [31], [26] for detailed information on deep fake detection methods.

## 2.2   Gendered Differences in Facial Analytics

There is consensus in the published literature that face analytics-based computer vision applications obtain lower accuracy for females, who often have both a higher false match and a higher false non-match rate over males [8,4,18,5,17]. Examination of the fairness of the gender classification systems using commercial SDKs and deep learning-based CNNs suggest lower accuracy rates for females consistently [18,8]. 2019 Face Recognition Vendor test documents lower female accuracy rates across a broad range of algorithms and datasets[4]. Similarly, lower accuracy rates for females have been obtained for various in-house deep learning-based face recognition systems [5,4,30]. The cause and effect analysis suggests gendered hairstyles resulting in facial occlusion, make-up, and inherent lower variability between different female faces over males to be the factors contributing to lower performance for females [4,5]. The demographic balanced datasets have been proven to mitigate the performance differential of different facial analysis based applications across demographics [4,18,17].

## 3   GBDF: Gender Balanced DeepFake Dataset

The GBDF dataset is created using FF++($c23$ version), Celeb-DF, Deeper Forensics-1.0 and consist of $10,000$ videos with 5000 each for males and females.

---

[4]     https://www.nist.gov/system/files/documents/2019/11/20/frvt_report_2019_11_19_0.pdf

The FaceForensics++ [28] (FF++) is an automated benchmark for facial manipulation detection. It consists of several manipulated videos created using two different generation techniques: Identity Swapping (FaceSwap, FaceSwap-Kowalski, FaceShifter, Deep Fakes) and Expression swapping (Face2Face and NeuralTextures). The Celeb-DF [21] deep fake forensic dataset include 590 genuine videos from 59 celebrities as well as 5639 deep fake videos. Celeb-DF, in contrast to other datasets, has essentially no splicing borders, color mismatch, and inconsistencies in face orientation, among other evident deep fake visual artifacts. The deep fake videos in Celeb-DF are created using an encoder-decoder style model which results in better visual quality. The DeeperForensics-1.0 [16] is one of the largest deep fake datasets used for face forgery detection. It consists of $60,000$ videos that have around 17.6 million frames with substantial real-world perturbations. The dataset contains videos of 100 consented actors with 35 different perturbations. The real to fake videos ratio is 5:1 and the fake videos are generated by an end-to-end face-swapping framework.

**Gender Label Annotation**. As none of these existing deepfake datasets contain demographic information, we manually annotated ground truth gender labels for these datasets. To do so, we annotated each subject with the perceived gender male, female. Two graduate annotators were selected for the task of gender label annotation. For each subject, the annotators were presented with an average of 150 frames at various times in the video, which displayed the subject at different light angles and poses. The gender label was assigned to each video based on the consensus between the annotators. With the annotated gender labels, we evaluated the percentage of videos belonging to males, and females and those being irregular face-swaps. Recall that an irregular swap is defined as a swap where a person's face is swapped onto another person's face of a different gender. All the three datasets provided the IDs for pairs of swaps for all the manipulation methods, With the help of the available IDs which are unique for all the identities, we were able to segregate gender labels as well as irregular swaps. FaceForensics++ has 35.12% of irregular deepfakes. Irregular deepfakes were not found in Celeb-DF. DeeperForensics-1.0 dataset has negligible number of irregular swaps. To remain ethnically aware and to maintain demographic information, irregular swaps from FaceForensics++ and deeperforensics-1.0 datasets are not included in the GBDF dataset.

The gender annotated version of the live and deepfake videos (excluding irregular swaps) from these deepfakes datasets are merged to create GBDF dataset. Deepfakes in the GBDF dataset are created using different Identity Swapping (i.e., FaceSwap, FaceSwap-Kowalski, FaceShifter, Encoder-decoder style and End-to-end Face Swapping techniques) and Expression swapping (i.e., Face2Face and NeuralTextures) deepfake generation techniques. The majority of the videos in GBDF are from Caucasians. The ratio of real to fake videos in the GBDF dataset is 1 : 4. The GBDF is further divided into gender-balanced and subject independent training and testing subsets in the ratio of 70 : 30. Figure 2 illustrates the comparison of deepfake videos among existing Deepfake datasets and our GBDF. The number of videos in GBDF is higher than many of the existing

deepfake datasets shown on the x axis. The GBDF dataset is publicly available
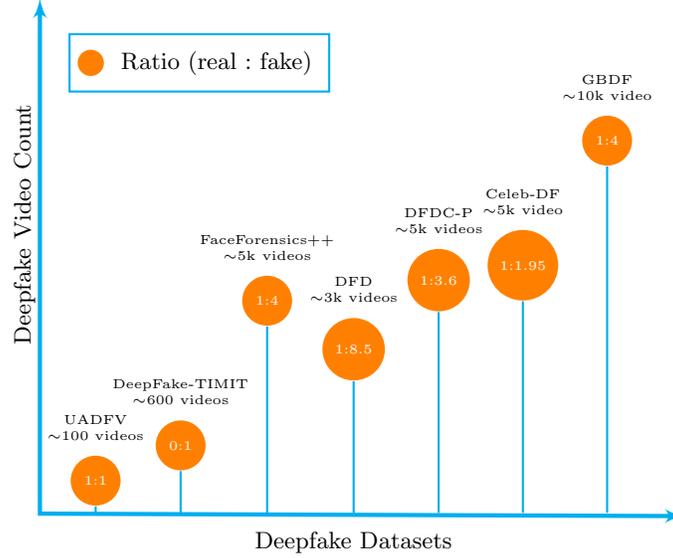at: `https://github.com/aakash4305/GBDF`



**Fig. 2.** Illustration of the number of videos in different deepfake datasets along with
our proposed GBDF dataset. The figure contains information about the real to fake
ratio of videos in the datasets along with deepfake video count. The no. of videos in
GBDF is higher than many of the existing deepfake datasets shown on the x-axis.

## 4   Deepfake Detection Algorithms Used

We investigated fairness of popular deepfake detection models of various sizes,
architectures and the underlying concept, across males and females. Specifically,
we trained MesoInception4[5], XceptionNet[6], EfficientNet V2-L[7], LipForensics[8]
and CNN-LSTM [9] based deepfake detectors.

  These models are trained on the popular FF++ dataset($c23$ version) and
our proposed GBDF training set. We used the sampling approach described
in [28] to choose 270 frames per video for training and 150 frames per video for
validation and testing of most of the models. The face images were detected and

---

[5] `https://github.com/HongguLiu/MesoNet-Pytorch`

[6] `https://github.com/i3p9/deepfake-detection-with-xception`

[7] `https://github.com/d-li14/efficientnetv2.pytorch`

[8] `https://github.com/ahaliassos/LipForensics`

[9] `https://github.com/oidelima/Deepfake-Detection`

aligned using MTCNN [35] algorithm. MTCNN utilizes a cascaded CNN based framework for joint face detection and alignment. The images are then resized to $256 \times 256$ for both training and evaluation.

For all the CNN-based models, we used a batch-normalization layer followed by the last fully connected layer of size 1024 and the final output layer for deep fake classification. The CNN models were trained using an Adam optimizer with an initial learning rate of 0.001 and a weight decay of 1e6. For CNN-LSTM model, we chose EfficientNet V2-L as the backbone CNN model due to its superior performance. The CNN network's output of 2048 feature vector is fed into the LSTM layer for deepfake detection. For LipForensics model, following authors implementation in [14], the network receives 25 grayscale, aligned mouth crops of size $88 \times 88$ as input for each video. The input is passed through pretrained ResNet-18 (pretrained for lipreading task with an initial 3-D convolutional layer) to obtain output embedding sensitive to mouth motion analysis. A multiscale temporal convolutional network (MS-TCN) was finetuned to detect fake videos based on semantically high-level anomalies in mouth motion, which was also pretrained for lipreading task. All the models were trained on 4 RTX 5000Ti GPUs with a batch size of 64.

## 5    Evaluation Metrics

Following the standard evaluation metrics adopted for deepfake detectors, we used partial AUC (pAUC) (at 10% False Positive Rate (FPR)) and Equal Error Rate (EER) for the evaluation of performance differences across males and females. Further, as deepfake detection is a binary classification task, we have also analyzed binary classification metrics for fairness evaluation across males and females. Similar to the bias evaluation study on gender classification by Buolamwini et al. [8], we follow the evaluation precedent established by the National Institute of Standards and Technology (NIST) and assessed the overall classification accuracy (ACC), along with the true positive rate (TPR), and false-positive rate (FPR) for males and females.

## 6    Results and Analysis

In this section, we examine the fairness of the deepfake detectors, discussed in section 4, across males and females on FF++, Celeb-DF, GBDF, and an external DFDC-P [12] test sets. All the **evaluation metrics** (from section 5) are reported in the range $[0, 1]$.

### 6.1    Performance differential of deepfake detectors on FF++ test set

Table 1 shows the performance of the deepfake detectors across males and females when trained on FF++, GBDF, and tested on FF++. Similarly, Table 2 shows the corresponding ACC, TPR, and FPR values of these models. The top

**Table 1.** Evaluation of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on **FF++**. The metrics used are AUC, pAUC and EER. The performance differential (P.D) is also calculated as the absolute difference between EER of males and females.

| Models | Training Dataset | Overall | | | Male | | | Female | | | P.D↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AUC | pAUC | EER | AUC | pAUC | EER | AUC | pAUC | EER | |
| EfficientNet V2-L | FaceForensics++ | **0.991** | **0.979** | **0.024** | **0.995** | **0.986** | **0.019** | 0.987 | 0.972 | 0.029 | 0.010 |
| XceptionNet | | 0.985 | 0.969 | 0.037 | 0.987 | 0.975 | 0.029 | 0.983 | 0.963 | 0.045 | 0.016 |
| MesoInception-4 | | 0.857 | 0.832 | 0.229 | 0.863 | 0.837 | 0.221 | 0.851 | 0.827 | 0.237 | 0.016 |
| CNN-LSTM | | 0.987 | 0.972 | 0.032 | 0.991 | 0.979 | 0.024 | 0.983 | 0.967 | 0.039 | 0.015 |
| LipForensics | | 0.990 | 0.977 | 0.027 | 0.987 | 0.975 | 0.031 | **0.993** | **0.979** | **0.023** | **0.008** |
| EfficientNet V2-L | GBDF | 0.925 | 0.902 | 0.136 | 0.935 | 0.917 | 0.121 | 0.915 | 0.888 | 0.140 | 0.019 |
| XceptionNet | | 0.906 | 0.886 | 0.176 | 0.912 | 0.892 | 0.172 | 0.899 | 0.880 | 0.180 | **0.008** |
| MesoInception-4 | | 0.806 | 0.785 | 0.264 | 0.813 | 0.794 | 0.259 | 0.799 | 0.775 | 0.269 | 0.010 |
| CNN-LSTM | | 0.918 | 0.897 | 0.141 | 0.910 | 0.890 | 0.150 | 0.926 | 0.904 | 0.132 | 0.018 |
| LipForensics | | **0.932** | **0.917** | **0.122** | **0.937** | **0.921** | **0.117** | **0.928** | **0.913** | **0.128** | 0.011 |

performance results are highlighted in bold across various evaluation datasets. EfficientNet V2-L obtained the best results with an overall AUC of 0.991, EER of 0.024, and ACC of 0.975 when trained and tested on FF++.

When trained on FF++, the overall difference in the performance is 0.009 and 0.010 in terms of pAUC and EER, respectively, across males and females. Males outperformed females for the majority of the models despite having a lower percentage than females in FF++ training set. *The reason is 35.12% of the videos in FF++ are irregular deepfakes, it is not certain which gender-group-related features are dominant in irregular facial swaps.* The overall difference in ACC, TPR, and FPR is 0.006,0.0036, and 0.020, respectively, across males and females (see Table 2). The least performance differential is obtained by LipForensics model when trained and tested on FF++.

When trained on GBDF, the overall difference in the performance is 0.010 and 0.006 in terms of pAUC and EER, respectively, across males and females. The overall difference in ACC, TPR, and FPR was reduced to 0.011, 0.006 and 0.009, respectively, across males and females (see Table 2). XceptionNet model obtained the least performance differential when trained on GBDF and tested on FF++.

Therefore, the overall difference in EER and FPR was reduced to 0.04 and 0.011, respectively, when using GBDF over FF++ as the training set. Using GBDF as the training set, the highest bias mitigation is observed for XceptionNet with the EER difference reduced from 0.016 to 0.008 across gender. Most of the detectors obtained lower error rates when trained on FF++. This is obvious as the test bed is also FF++. The performance of most of the models dropped when trained using GBDF due to domain shift i.e., the data distribution change between the training (GBDF) and testing set (FF++). This is due to change in the image quality of real videos and deep fakes due to advances in sensor technology and the deep fake generation techniques. The GBDF dataset has an

**Table 2.** ACC, TPR and FPR of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on FF++. When trained on GBDF, the drop in the performance of the models is due to domain shift. The GBDF dataset consist of higher number of deepfake generation techniques over FF+.

| Models | Training Datasets | Overall | | | Male | | | Female | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | ACC | TPR | FPR | ACC | TPR | FPR | ACC | TPR | FPR |
| EfficientNet V2-L | | **0.975** | **0.952** | **0.091** | **0.979** | **0.955** | **0.058** | 0.971 | 0.949 | 0.119 |
| XceptionNet | | 0.969 | 0.942 | 0.128 | 0.971 | 0.947 | 0.109 | 0.967 | 0.937 | 0.139 |
| MesoInception-4 | FaceForensics++ | 0.825 | 0.805 | 0.256 | 0.834 | 0.813 | 0.245 | 0.816 | 0.797 | 0.267 |
| CNN-LSTM | | 0.971 | 0.945 | 0.115 | 0.976 | 0.954 | 0.093 | 0.966 | 0.936 | 0.137 |
| LipForensics | | 0.972 | 0.948 | 0.115 | 0.967 | 0.941 | 0.143 | **0.978** | **0.955** | **0.086** |
| EfficientNet V2-L | | 0.903 | 0.887 | 0.182 | 0.912 | 0.895 | 0.175 | 0.892 | 0.879 | 0.187 |
| XceptionNet | | 0.888 | 0.869 | 0.189 | 0.897 | 0.876 | 0.181 | 0.879 | 0.862 | 0.195 |
| MesoInception-4 | GBDF | 0.783 | 0.769 | 0.284 | 0.794 | 0.778 | 0.276 | 0.772 | 0.760 | 0.292 |
| CNN-LSTM | | 0.896 | 0.875 | 0.185 | 0.887 | 0.861 | 0.189 | 0.905 | 0.889 | 0.178 |
| LipForensics | | **0.912** | **0.896** | **0.176** | **0.919** | **0.901** | **0.169** | **0.905** | **0.891** | **0.183** |

additional number of deepfake generation techniques (based on encoder-decoder style and the end-to-end face swapping framework) over FF++.

## 6.2 Performance differential of deepfake detectors on Celeb-DF test set

**Table 3.** Evaluation of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on **Celeb-DF**. The metrics used are AUC, pAUC and EER. The performance differential (P.D) is also calculated as the absolute difference between EER of males and females.

| Models | Training Dataset | Overall | | | Male | | | Female | | | P.D↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AUC | pAUC | EER | AUC | pAUC | EER | AUC | pAUC | EER | |
| EfficientNet V2-L | | 0.658 | 0.635 | 0.379 | 0.667 | 0.645 | 0.372 | 0.649 | 0.625 | 0.386 | 0.014 |
| XceptionNet | | 0.651 | 0.629 | 0.383 | 0.657 | 0.634 | 0.379 | 0.645 | 0.623 | 0.390 | **0.011** |
| MesoInception-4 | FaceForensics++ | 0.544 | 0.519 | 0.459 | 0.558 | 0.528 | 0.442 | 0.530 | 0.510 | 0.476 | 0.034 |
| CNN-LSTM | | 0.675 | 0.656 | 0.359 | 0.686 | 0.662 | 0.348 | 0.664 | 0.650 | 0.370 | 0.022 |
| LipForensics | | **0.821** | **0.795** | **0.254** | **0.829** | **0.805** | **0.242** | **0.813** | **0.785** | **0.266** | 0.024 |
| EfficientNet V2-L | | 0.861 | 0.844 | 0.235 | 0.869 | 0.853 | 0.228 | 0.853 | 0.835 | 0.242 | 0.014 |
| XceptionNet | | 0.864 | 0.847 | 0.233 | 0.872 | 0.855 | 0.226 | 0.856 | 0.839 | 0.240 | 0.014 |
| MesoInception-4 | GBDF | 0.742 | 0.725 | 0.298 | 0.755 | 0.735 | 0.292 | 0.730 | 0.715 | 0.305 | 0.013 |
| CNN-LSTM | | 0.887 | 0.869 | 0.215 | 0.898 | 0.875 | 0.209 | 0.876 | 0.863 | 0.221 | **0.012** |
| LipForensics | | **0.908** | **0.885** | **0.175** | **0.917** | **0.896** | **0.163** | **0.900** | **0.874** | **0.187** | 0.024 |

Table 3 shows the performance differential of the deepfake detectors when trained on FF++, GBDF, and tested on Celeb-DF. Similarly, Table 4 shows

the corresponding ACC, TPR, and FPR values for these models. The top performance results are highlighted in bold across various evaluation datasets. The LipForensics model obtained the best results with an overall AUC of 0.908, EER of 0.175, and ACC of 0.889 when trained on GBDF and tested on Celeb-DF.

When trained on FF++, the overall difference in the performance is 0.0162 and 0.021 in terms of pAUC and EER, respectively, across males and females. The overall difference in ACC, TPR and FPR is 0.019,0.02 and 0.021, respectively, across males and females (see Table 4). The least performance differential is obtained by XceptionNet when trained on FF++ and tested on Celeb-DF.

When trained on GBDF, the overall difference in the performance is 0.017 and 0.015 in terms of pAUC and EER, respectively, across males and females. The overall difference in ACC, TPR and FPR is 0.018,0.01 and 0.018, respectively, across males and females (see Table 4). The least performance differential is obtained by CNN-LSTM when trained on GBDF and tested on Celeb-DF.

Therefore, the difference in AUC, EER, TPR, and FPR is reduced to 0.001, 0.006, 0.01, and 0.003, respectively, when using GBDF as a training set over FF++. Using GBDF, the highest bias mitigation is observed for MesoInceptionNet-4 model with the EER difference reduced from 0.034 to 0.013 across gender. The overall performance of all the models increased when trained on GBDF over FF++ because of the presence of higher number of deepfake generation techniques. **It is worth noting that the training and testing subset of GBDF and Celeb-DF, respectively, has no subject overlap**. This experiment points out the **merit** of using a demographically balanced dataset for deepfake detection.

**Table 4.** ACC, TPR and FPR of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on Celeb-DF.

| Models | Training Datasets | Overall | | | Male | | | Female | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | ACC | TPR | FPR | ACC | TPR | FPR | ACC | TPR | FPR |
| EfficientNet V2-L | FaceForensics++ | 0.637 | 0.604 | 0.385 | 0.650 | 0.614 | 0.372 | 0.626 | 0.594 | 0.398 |
| XceptionNet | | 0.629 | 0.602 | 0.395 | 0.635 | 0.609 | 0.383 | 0.623 | 0.595 | 0.402 |
| MesoInception-4 | | 0.525 | 0.502 | 0.437 | 0.534 | 0.518 | 0.422 | 0.516 | 0.486 | 0.455 |
| CNN-LSTM | | 0.652 | 0.609 | 0.367 | 0.664 | 0.615 | 0.359 | 0.640 | 0.600 | 0.379 |
| LipForensics | | **0.798** | **0.774** | **0.275** | **0.807** | **0.785** | **0.271** | **0.791** | **0.763** | **0.282** |
| EfficientNet V2-L | GBDF | 0.843 | 0.825 | 0.242 | 0.849 | 0.833 | 0.239 | 0.837 | 0.817 | 0.246 |
| XceptionNet | | 0.847 | 0.825 | 0.240 | 0.854 | 0.834 | 0.232 | 0.840 | 0.816 | 0.251 |
| MesoInception-4 | | 0.718 | 0.701 | 0.324 | 0.733 | 0.712 | 0.309 | 0.703 | 0.690 | 0.340 |
| CNN-LSTM | | 0.863 | 0.849 | 0.225 | 0.876 | 0.854 | 0.211 | 0.850 | 0.844 | 0.235 |
| LipForensics | | **0.889** | **0.866** | **0.187** | **0.895** | **0.878** | **0.183** | **0.883** | **0.854** | **0.193** |

**Table 5.** Evaluation of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on **GBDF**. The metrics used are AUC, pAUC and EER. The performance differential (P.D) is calculated as the absolute difference between EER of males and females.

| Models | Training Dataset | Overall | | | Male | | | Female | | | P.D↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AUC | pAUC | EER | AUC | pAUC | EER | AUC | pAUC | EER | |
| EfficientNet V2-L | | 0.904 | 0.889 | 0.179 | 0.912 | 0.897 | 0.171 | 0.896 | 0.879 | 0.187 | 0.016 |
| XceptionNet | | 0.889 | 0.868 | 0.217 | 0.902 | 0.885 | 0.206 | 0.876 | 0.850 | 0.228 | 0.022 |
| MesoInception-4 | FaceForensics++ | 0.769 | 0.747 | 0.286 | 0.759 | 0.742 | 0.295 | 0.779 | 0.750 | 0.277 | 0.018 |
| CNN-LSTM | | 0.909 | 0.888 | 0.177 | 0.917 | 0.898 | 0.161 | 0.901 | 0.877 | 0.192 | 0.031 |
| LipForensics | | **0.942** | **0.926** | **0.109** | **0.938** | **0.922** | **0.113** | **0.947** | **0.929** | **0.105** | **0.008** |
| EfficientNet V2-L | | 0.967 | 0.943 | 0.052 | 0.972 | 0.948 | 0.050 | 0.962 | 0.938 | 0.054 | **0.004** |
| XceptionNet | | 0.972 | 0.952 | 0.046 | 0.979 | 0.956 | 0.043 | 0.965 | 0.948 | 0.049 | 0.006 |
| MesoInception-4 | GBDF | 0.819 | 0.800 | 0.256 | 0.828 | 0.805 | 0.250 | 0.811 | 0.795 | 0.264 | 0.014 |
| CNN-LSTM | | 0.975 | 0.957 | 0.044 | 0.983 | 0.964 | 0.038 | 0.967 | 0.950 | 0.050 | 0.012 |
| LipForensics | | **0.978** | **0.954** | **0.039** | **0.982** | **0.958** | **0.036** | **0.974** | **0.950** | **0.042** | 0.006 |

### 6.3  Performance differential of deepfake detectors on GBDF and DFDC-P test sets

Table  5 shows the performance of the deepfake detectors across males and females when trained on FF++, GBDF, and tested on GBDF subject independent test set. Similarly, Table 6 shows the ACC, TPR, and FPR values associated with these models. The LipForensics model obtained the best results with an overall AUC of 0.978, EER of 0.039, and ACC of 0.967 when trained and tested on GBDF.

When trained on FF++, the overall difference in the performance is 0.012 and 0.0092 in terms of pAUC and EER, respectively, across males and females. The overall difference in ACC, TPR and FPR is 0.010, 0.0112 and 0.012, respectively, across males and females (see Table 6). The least performance differential is obtained by EfficientNet V2-L when trained on FF++ and tested on GBDF.

When trained on GBDF, the overall difference in the performance is 0.010 and 0.008 in terms of pAUC and EER, respectively, across males and females. The overall difference in ACC, TPR and FPR is 0.008,0.010 and 0.009, respectively, across males and females (see Table 6). The least performance differential is obtained by the LipForensics model when trained and tested on GBDF.

Therefore, the difference in ACC, EER, TPR, and FPR decreased by 0.002, 0.001, 0.0012, and 0.003, respectively, when using balanced GBDF as training and testing sets. Using balanced GBDF as a training and testing set, the highest bias mitigation is observed for CNN-LSTM and XceptionNet models. For CNN-LSTM, the difference in EER across gender reduced from 0.031 to 0.012 when trained with FF++ over GBDF as the training set (the test set is GBDF). Similarly, for XceptionNet, the difference in EER across gender reduced from 0.022 to 0.006 when trained with FF++ over GBDF as the training set (the test set is GBDF). Recall that the subjects do not overlap between the training

**Table 6.** ACC, TPR and FPR of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on GBDF test set.

| Models | Training Datasets | Overall | | | Male | | | Female | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | ACC | TPR | FPR | ACC | TPR | FPR | ACC | TPR | FPR |
| EfficientNet V2-L | | 0.880 | 0.858 | 0.205 | 0.895 | 0.869 | 0.194 | 0.865 | 0.847 | 0.220 |
| XceptionNet | | 0.865 | 0.841 | 0.222 | 0.878 | 0.854 | 0.209 | 0.852 | 0.828 | 0.235 |
| MesoInception-4 | FaceForensics++ | 0.745 | 0.721 | 0.288 | 0.735 | 0.715 | 0.295 | 0.754 | 0.727 | 0.275 |
| CNN-LSTM | | 0.883 | 0.868 | 0.195 | 0.894 | 0.884 | 0.187 | 0.872 | 0.854 | 0.209 |
| LipForensics | | **0.925** | **0.905** | **0.178** | **0.920** | **0.899** | **0.180** | **0.930** | **0.910** | **0.175** |
| EfficientNet V2-L | | 0.948 | 0.935 | 0.154 | 0.951 | 0.939 | 0.149 | 0.945 | 0.930 | 0.159 |
| XceptionNet | | 0.955 | 0.941 | 0.144 | 0.958 | 0.947 | 0.139 | 0.952 | 0.935 | 0.147 |
| MesoInception-4 | GBDF | 0.802 | 0.778 | 0.282 | 0.808 | 0.788 | 0.275 | 0.796 | 0.770 | 0.287 |
| CNN-LSTM | | 0.953 | 0.939 | 0.148 | 0.959 | 0.941 | 0.144 | 0.949 | 0.935 | 0.154 |
| LipForensics | | **0.967** | **0.949** | **0.142** | **0.971** | **0.953** | **0.135** | **0.965** | **0.946** | **0.145** |

and testing set of GBDF. Further, the samples in the training and testing set of GBDF are from three different deepfake datasets.

**Table 7.** Evaluation of the DeepFake Detectors Across Males and Females when trained on FF++, GBDF and tested on **DFDC-P**. The metrics used are AUC, pAUC and EER. The performance differential (P.D) is calculated as the absolute difference between EER of males and females.

| Models | Training Dataset | Overall | | | Male | | | Female | | | P.D↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AUC | pAUC | EER | AUC | pAUC | EER | AUC | pAUC | EER | |
| EfficientNet V2-L | | 0.659 | 0.634 | 0.378 | 0.665 | 0.641 | 0.374 | 0.653 | 0.626 | 0.384 | **0.010** |
| XceptionNet | | 0.642 | 0.624 | 0.391 | 0.649 | 0.632 | 0.384 | 0.635 | 0.617 | 0.398 | 0.014 |
| MesoInception-4 | FaceForensics++ | 0.619 | 0.597 | 0.421 | 0.609 | 0.591 | 0.432 | 0.630 | 0.605 | 0.410 | 0.022 |
| CNN-LSTM | | 0.667 | 0.648 | 0.372 | 0.675 | 0.661 | 0.356 | 0.659 | 0.635 | 0.382 | 0.026 |
| LipForensics | | **0.718** | **0.705** | **0.312** | **0.724** | **0.710** | **0.306** | **0.712** | **0.701** | **0.319** | 0.013 |
| EfficientNet V2-L | | 0.684 | 0.662 | 0.349 | 0.689 | 0.665 | 0.345 | 0.680 | 0.661 | 0.354 | **0.009** |
| XceptionNet | | 0.668 | 0.652 | 0.368 | 0.675 | 0.657 | 0.363 | 0.663 | 0.649 | 0.374 | 0.011 |
| MesoInception-4 | GBDF | 0.615 | 0.592 | 0.427 | 0.621 | 0.596 | 0.419 | 0.608 | 0.585 | 0.432 | 0.013 |
| CNN-LSTM | | 0.689 | 0.665 | 0.343 | 0.683 | 0.658 | 0.350 | 0.694 | 0.674 | 0.334 | 0.016 |
| LipForensics | | **0.732** | **0.721** | **0.299** | **0.736** | **0.724** | **0.292** | **0.727** | **0.716** | **0.307** | 0.015 |

Table  7 shows the performance of the deepfake detectors across males and females when trained on FF++, GBDF, and tested on DFDC-P. **Note that DFDC-P dataset has not been used in the creation of the GBDF dataset**. As the original DFDC-P test set does not contain subject IDs, the subset of DFDC training set is manually annotated with gender labels and used as a test set for this study. Overall, low performance is obtained for all the models on DFDC dataset. This is because DFDC consist of low quality videos that are diverse across gender, skin-tone and age-group. The LipForensics model

obtained the best results with an overall AUC of 0.732, EER of 0.299 when tested on DFDC-P.

When trained on FF++, the overall difference in the performance is 0.010 and 0.0082 in terms of pAUC and EER, respectively, across males and females. The least performance differential is obtained by EfficientNet V2-L when trained on FF++ and tested on DFDC-P.

When trained on GBDF, the overall difference in the performance is 0.004 and 0.007 in terms of pAUC and EER, respectively, across males and females. The least performance differential is obtained by EfficientNet V2-L when trained on GBDF and tested on DFDC-P. *These results suggest that using our gender-balanced GBDF training set, bias is mitigated across gender even on an external DFDC-P dataset, not used in the creation of GBDF.*



(a) Real Images

(b) Fake Images

**Fig. 3.** Grad-CAM visualization of the EfficientNet V2-L based deepfake detector on randomly selected live and fake samples from males and females. The distinctive image regions used by the CNN model for deepfake detection differs across gender.

Finally, we also used Explainable AI (XAI) based Gradient weighted Class Activation Mapping (Grad-CAM) [29] visualization to understand the distinctive image regions used by the CNN models in detecting deepfakes across gender. GRAD-CAM uses the gradients of any target concept to generate a coarse localization map that highlights distinctive image regions used for making a decision/prediction [29]. Figure 3 shows the GRAD-CAM visualization of the EfficientNet V2-L-based deepfake detector for live and fake images for males and females. This detector was trained on GBDF dataset. The highly activated region is shown by the red zone on the map, followed by green and blue zones. It can be seen that the highly activated region is the cheek for females and the

ocular region for males. For fake images, the mouth and cheek region for males and the complete face region for females are the most activated region. These results were consistent across the datasets depending on the deepfake generation technique. Therefore, different image regions were used by the deepfake detector for live and fake classification across gender.

In **summary**, males outperformed females for most of the models, with the disparity of about 0.034 in terms of EER in the range $[0, 1]$ for MesoInception-4 model. The shallow MesoInception-4 model demonstrated high performance differential across gender for most of the experiments. The LipForensics model, on the other hand, obtained least disparity across gender for most of the experiments. This is because it uses mouth crops for mouth motion analysis. Thus the impact of gendered differences in facial images attributed to bias are mitigated to a major extent. When trained on FF+, males outperformed females for the majority of the models despite having a lower percentage than females. As large number of the videos in FaceForensics++ are irregular deepfakes, it is not certain which gender-group-related features are dominant in irregular facial swaps. The gender-balanced GBDF training set reduced the performance difference over FF++, with the highest being from 0.031 to 0.012 in terms of EER across males and females when tested on GBDF test set. The advantage of using GBDF training set towards gender fair deepfake classification is also noticed for an external DFDC-P set. The grad-CAM visualization suggests the distinctive image regions used by the CNN model for deepfake classification differs across gender. As these automated deepfake detection systems are used at the mass-level for audit of the social media data, even a small reduction in the bias across demographics would positively impact millions of people belonging to the deprived sub-group.

## 7   Conclusion and Future Research Directions

With the volume of deepfake videos showing staggering growth, there is a growing reliance on automated systems to combat deepfakes. For the massive rollout of this high-impact technology, it becomes vital to understand all the societal aspects including demographic disparities. In this work, we thoroughly examined the fairness of the deepfake detectors on gender-aware deepfake datasets. On manual annotation of gender labels, we found that current deepfake datasets have a highly skewed distribution across gender and contain irregular swaps. The popular deepfake detectors have exhibited disparities in the performance across gender when evaluated on gender-aware datasets, with mostly males outperforming females. This suggest an additional threat imposed by deep fake technology on female subjects, primarily due to the performance differential of SOTA deep fake detectors.

However, using our gender-balanced GBDF dataset, the unequal performance of the deepfake detectors across gender is mitigated to some extent. Our work echoes the importance of benchmarking demographically balanced and labeled deepfake datasets to facilitate intersectional subgroup-based audits of existing

deepfake detectors along with the cause and effect analysis. As a part of future work, fairness of the deepfake detectors will also be evaluated across race. Further, the fairness-aware deepfake detectors will be developed for increased demographic transparency and accountability of these high-impact systems.

## 8    Acknowledgement

## References

1. Afchar, D., Nozick, V., Yamagishi, J., Echizen, I.: Mesonet: a compact facial video forgery detection network. In: 2018 IEEE International Workshop on Information Forensics and Security (WIFS). pp. 1–7 (2018). https://doi.org/10.1109/WIFS.2018.8630761
2. Agarwal, S., Farid, H., El-Gaaly, T., Lim, S.N.: Detecting deep-fake videos from appearance and behavior. In: 2020 IEEE International Workshop on Information Forensics and Security (WIFS). pp. 1–6 (2020). https://doi.org/10.1109/WIFS49906.2020.9360904
3. Agarwal, S., Farid, H., Gu, Y., He, M., Nagano, K., Li, H.: Protecting world leaders against deep fakes. In: CVPR Workshops (2019)
4. Albiero, V., Zhang, K., Bowyer, K.W.: How does gender balance in training data affect face recognition accuracy? In: 2020 ieee international joint conference on biometrics (ijcb). pp. 1–10. IEEE (2020)
5. Albiero, V., Zhang, K., King, M.C., Bowyer, K.W.: Gendered differences in face recognition accuracy explained by hairstyles, makeup, and facial morphology. Trans. Info. For. Sec. **17**, 127–137 (jan 2022). https://doi.org/10.1109/TIFS.2021.3135750, `https://doi.org/10.1109/TIFS.2021.3135750`
6. Amerini, I., Caldelli, R.: Exploiting prediction error inconsistencies through lstm-based classifiers to detect deepfake videos. In: Proceedings of the 2020 ACM Workshop on Information Hiding and Multimedia Security. p. 97–102. IH&MMSec '20, Association for Computing Machinery, New York, NY, USA (2020). https://doi.org/10.1145/3369412.3395070, `https://doi.org/10.1145/3369412.3395070`
7. Biometrics, I.J.S.: Iso/iec wd tr 22116. In: information technology – biometrics – identifying and mitigating the differential impact of demographic factors in biometric systems (unpublished)
8. Buolamwini, J., Gebru, T.: Gender shades: Intersectional accuracy disparities in commercial gender classification. In: Conference on fairness, accountability and transparency. pp. 77–91. PMLR (2018)
9. Cellan-Jones, R.: Deepfake videos 'double in nine months' (Oct 2019), `https://www.bbc.com/news/technology-49961089`

10. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1800–1807. IEEE Computer Society, Los Alamitos, CA, USA (jul 2017). https://doi.org/10.1109/CVPR.2017.195, `https://doi.ieeecomputersociety.org/10.1109/CVPR.2017.195`

11. Citron, D.: How deepfakes undermine truth and threaten democracy, `https://www.ted.com/talks/danielle_citron_how_deepfakes_undermine_truth_and_threaten_democracy?language=en`

12. Dolhansky, B., Bitton, J., Pflaum, B., Lu, J., Howes, R., Wang, M., Ferrer, C.C.: The deepfake detection challenge (dfdc) dataset (2020). https://doi.org/10.48550/ARXIV.2006.07397, `https://arxiv.org/abs/2006.07397`

13. Dong, X., Bao, J., Chen, D., Zhang, W., Yu, N., Chen, D., Wen, F., Guo, B.: Identity-driven deepfake detection. ArXiv **abs/2012.03930** (2020)

14. Haliassos, A., Vougioukas, K., Petridis, S., Pantic, M.: Lips don't lie: A generalisable and robust approach to face forgery detection. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5037–5047 (2021). https://doi.org/10.1109/CVPR46437.2021.00500

15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 770–778 (2016)

16. Jiang, L., Li, R., Wu, W., Qian, C., Loy, C.C.: DeeperForensics-1.0: A large-scale dataset for real-world face forgery detection. In: CVPR (2020)

17. Karkkainen, K., Joo, J.: 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. pp. 1548–1558 (2021)

18. Krishnan, A., Almadan, A., Rattani, A.: Understanding fairness of gender classification algorithms across gender-race groups. In: 2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA). pp. 1028–1035 (2020). https://doi.org/10.1109/ICMLA51294.2020.00167

19. Li, L., Bao, J., Zhang, T., Yang, H., Chen, D., Wen, F., Guo, B.: Face x-ray for more general face forgery detection. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5000–5009 (2020). https://doi.org/10.1109/CVPR42600.2020.00505

20. Li, Y., Lyu, S.: Exposing deepfake videos by detecting face warping artifacts. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (June 2019)

21. Li, Y., Yang, X., Sun, P., Qi, H., Lyu, S.: Celeb-df: A large-scale challenging dataset for deepfake forensics. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3204–3213 (2020). https://doi.org/10.1109/CVPR42600.2020.00327

22. Matern, F., Riess, C., Stamminger, M.: Exploiting visual artifacts to expose deepfakes and face manipulations. 2019 IEEE Winter Applications of Computer Vision Workshops (WACVW) pp. 83–92 (2019)

23. Nadimpalli, A.V., Rattani, A.: On improving cross-dataset generalization of deepfake detectors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 91–99 (2022)

24. Nadimpalli, A.V., Reddy, N., Ramachandran, S., Rattani, A.: Harnessing unlabeled data to improve generalization of biometric gender and age classifiers. 2021 IEEE Symposium Series on Computational Intelligence (SSCI) pp. 1–7 (2021)

25. Nguyen, H.H., Fang, F., Yamagishi, J., Echizen, I.: Multi-task learning for detecting and segmenting manipulated facial images and videos. 2019 IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS) pp. 1–8 (2019)

26. Nguyen, T.T., Nguyen, C.M., Nguyen, D., Nguyen, D.T., Nahavandi, S.: Deep learning for deepfakes creation and detection. ArXiv **abs/1909.11573** (2019)

27. Ramachandran, S., Nadimpalli, A.V., Rattani, A.: An experimental evaluation on deepfake detection using deep face recognition. In: 2021 IEEE International Carnahan Conference on Security Technology (ICCST). pp. 1–6 (2021). https://doi.org/10.1109/ICCST49569.2021.9717407

28. Rössler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Niessner, M.: Faceforensics++: Learning to detect manipulated facial images. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 1–11 (2019). https://doi.org/10.1109/ICCV.2019.00009

29. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 618–626 (2017). https://doi.org/10.1109/ICCV.2017.74

30. Singh, R., Majumdar, P., Mittal, S., Vatsa, M.: Anatomizing bias in facial analysis. arXiv preprint arXiv:2112.06522 (2021)

31. Tolosana, R., Vera-Rodríguez, R., Fierrez, J., Morales, A., Ortega-Garcia, J.: Deepfakes and beyond: A survey of face manipulation and fake detection. Inf. Fusion **64**, 131–148 (2020)

32. Trinh, L., Liu, Y.: An examination of fairness of ai models for deepfake detection. ArXiv **abs/2105.00558** (2021)

33. Verdoliva, L.: Media forensics and deepfakes: An overview. IEEE Journal of Selected Topics in Signal Processing **14**, 910–932 (2020)

34. Wang, M., Deng, W., Hu, J., Tao, X., Huang, Y.: Racial faces in the wild: Reducing racial bias by information maximization adaptation network. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 692–702 (2019)

35. Zhang, K., Zhang, Z., Li, Z., Qiao, Y.: Joint face detection and alignment using multitask cascaded convolutional networks. IEEE Signal Processing Letters **23**(10), 1499–1503 (2016). https://doi.org/10.1109/LSP.2016.2603342

36. Zhao, H., Wei, T., Zhou, W., Zhang, W., Chen, D., Yu, N.: Multi-attentional deepfake detection. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2185–2194 (2021). https://doi.org/10.1109/CVPR46437.2021.00222