

# Improving DeepFake Detection Using Dynamic Face Augmentation

Sowmen Das<sup>1†</sup> Arup Datta<sup>2†</sup> Md. Saiful Islam<sup>3†</sup> Md. Ruhul Amin<sup>4§</sup>

<sup>†</sup>Shahjalal University of Science and Technology, Bangladesh <sup>§</sup>Fordham University, USA

<sup>1</sup>sowmendipta@gmail.com  
<sup>3</sup>saif.acm@gmail.com

<sup>2</sup>arup.cse.sust@gmail.com  
<sup>4</sup>mamin17@fordham.edu

## Abstract

The creation of altered and manipulated faces has become more common due to the improvement of DeepFake generation methods. Simultaneously, we have seen detection models' development for differentiating between a manipulated and original face from image or video content. We have observed that most publicly available DeepFake detection datasets have limited variations, where a single face is used in many videos, resulting in an oversampled training dataset. Due to this, deep neural networks tend to overfit to the facial features instead of learning to detect manipulation features of DeepFake content. As a result, most detection architectures perform poorly when tested on unseen data. In this paper, we provide a quantitative analysis to investigate this problem and present a solution to prevent model overfitting due to the high volume of samples generated from a small number of actors. We introduce Face-Cutout, a data augmentation method for training Convolutional Neural Networks (CNN), to improve DeepFake detection. In this method, training images with various occlusions are dynamically generated using face landmark information irrespective of orientation. Unlike other general-purpose augmentation methods, it focuses on the facial information that is crucial for DeepFake detection. Our method achieves a reduction in LogLoss of 15.2% to 35.3% on different datasets, compared to other occlusion-based augmentation techniques. We show that Face-Cutout can be easily integrated with any CNN-based recognition model and improve detection performance.

## 1. Introduction

The term *DeepFake* has gained much attention in recent times. It is used to denote manipulated multimedia content, specifically video or images created or altered using deep learning techniques. Although the manipulation of digital media is not new, the use of deep neural architectures for this purpose has been gaining popularity with the increasing

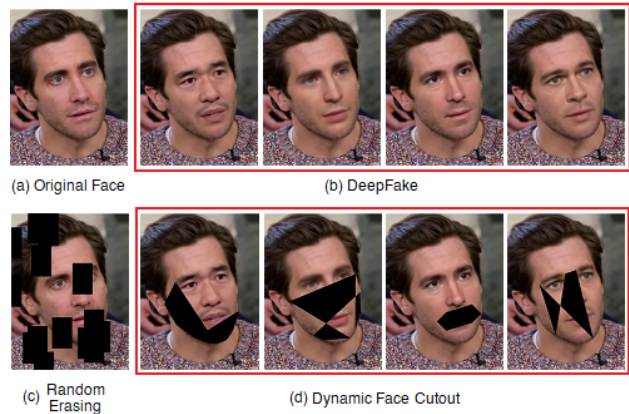


Figure 1: (a) & (b) Example of multiple DeepFake images created from an original face. (c) Random Erase augmentation. (d) Results of Dynamic Face-Cutout. Compared to Random-Erasing, our method augments a face based on the facial landmarks irrespective of its orientation and can identify manipulated facial features.

improvements in generative learning. The vast leaps in the development of Generative Adversarial Networks (GAN) [15] and Variational Autoencoders (VAE) [19] for generative modeling has enabled the creation of media forgeries almost unidentifiable by human observers. These forgeries include but are not limited to swapping of a person's face in a video with that of another [21], facial reenactment, i.e., transferring the facial expressions from one video to another [35], changing/altering physical or facial appearances [4] etc. While there are humorous applications of DeepFake videos, when created with malicious intent, they have the potential to harass individuals and spread misinformation. DeepFakes are also a cause of significant concern for the authenticity of news as well as privacy. In recent days, manipulated videos of political leaders have been heavily used to spread false information to bias public election outcomes.

In response to the increasing development of DeepFake generators, many actions are being taken to develop techniques for detecting manipulated content. Facebook re-

cently partnered with AWS and Microsoft to launch the DeepFake Detection Challenge (DFDC) to improve DeepFake detection performance.

The results from the competition showed that even though state-of-the-art image classification models like EfficientNet [34], ResNeXt [38], XceptionNet [7] etc. can extract manipulation features, their performance is unsatisfactory due to extremely high model overfitting. The best result from the competition is a LogLoss of 0.20336 achieved by *Selim Seferbekov* [10] on the public test dataset. His solution employed the EfficientNet-B7 architecture, an image classifier with the highest 88.5% Top-1 accuracy on the Imagenet [9] benchmark dataset.

All currently available public DeepFake datasets include both the real and manipulated videos and images. Datasets like UADFV [40], FaceForensics++ (FF++) [30] and CelebDF [25] were created using videos collected from online video streaming sites like YouTube. On the other hand, Google DFD [12] and the DFDC dataset was created using video clips recorded by professional actors. Fake videos were then generated from these recorded or collected videos using different DeepFake generators. From our analysis, we found that in most cases, a single source-video was used to generate multiple types of manipulations, as shown in Fig. 1. As a result, a *single* face was over-sampled to generate multiple fakes. Many DeepFake generators apply small changes to a source face to alter the facial expressions [27], making it challenging to differentiate manipulated and real faces. This lack of variation and oversampling causes deep neural networks to quickly overfit the data before learning the necessary features for DeepFake identification. Instead of learning manipulation features, the networks start to memorize the subjects' faces, resulting in poor performance. To solve this issue, we propose a data augmentation method that improves the existing *Random-Erasing* [42] augmentation.

Our proposed method, termed as *Face-Cutout*, is an erasing technique that repaints groups of pixels of different shapes on an image using face landmark information. Whereas, *Random-Erasing* employs the method of cutout by selecting rectangles of various sizes and replaces the pixels with random values. The size, shape, and number of those rectangles need to be manually adjusted. We have found in our experiments that plain *Random-Erasing* sometimes hampers performance for DeepFake detection as it does not consider the underlying pixel information of the erased segment. On the other hand, our proposed method uses the face landmark information and also the DeepFake locations to dynamically select the best cutout region, as can be seen from Fig. 1. It provides selective attention by occluding areas of the face that do not contain fake information. In addition, the proposed method can help deep learning models to learn the DeepFake features better by

generating random variations from an input image. To summarize the contributions of our paper,

- We provide a comprehensive analysis of popular DeepFake datasets to identify their shortcomings. We show the use of face clustering for evaluating the datasets and its use in data preprocessing to mitigate overfitting.
- We propose *Face-Cutout*, an erasing technique that uses the facial landmarks and underlying image information to dynamically determine cutout regions for augmentation.
- We show that our technique significantly improves DeepFake detection performance of existing CNN architectures by reducing model overfitting.
- We further present that dynamic *Face-Cutout* helps in the generalization of DeepFake detection models, independent of any specific dataset.

This paper is organized as follows. In Sec. 2, we talk about the existing work on deepfake generation and detection techniques. In Sec. 3 we discuss some of the currently available deepfake datasets and explain their limitations with the use of face clustering. In Sec. 4, we propose our augmentation method to overcome these limitations. We show the results of our experimentation in Sec. 5 and 6. Lastly, in Sec. 7 we conclude our work.

## 2. Related Works

Most state-of-the-art DeepFake generators utilize GANs, a two-part network consisting of a generator and a discriminator. The generator learns and replicates the features from the training data, and the discriminator tries to differentiate real from generated samples. Both the networks are trained simultaneously until the generator can produce samples well enough to *deceive* the discriminator. Multiple variations of GAN architectures [17, 6, 44, 16] are being used for DeepFake generation that are capable of performing image-to-image translation and generation of high resolution fake faces. An extensive study of DeepFake generation using GAN and Autoencoders is presented in [27].

On the other hand, to detect the DeepFake content, several different architectures have been proposed so far. Shallow networks were proposed in [1] to exploit the mesoscopic features in DeepFake videos having high compression artifacts. XceptionNet has been shown to perform very well in identifying facial forgeries in [30]. A Siamese approach is used in [45] that learns a difference function from both real and manipulated frames to encode deepfake features. In [28] multi-task learning was carried out to simultaneously classify and localize deepfakes using a Y-shaped encoder-decoder network. The use of capsule networks for forgery detection has been highlighted in [29]. In [8] a visual

attention network with supervised and unsupervised learning approach has been shown. They were able to gain high detection performance by integrating an attention module with existing image classifiers. In addition, many studies have been done to detect the various artifacts left by DeepFake generators, such as face warping artifacts [24], temporal and spatial inconsistencies [31], eye blinking [23], inconsistent head poses [40], etc. In [3] the authors conducted many experiments to determine which image features are most effective for generalizing fake face detection. A survey of multiple DeepFake detection architectures and their performance comparison has been presented in [36].

Although several studies have been done on designing architectures and identifying feature descriptors for fake face classification, there has been little analysis of the effects of data preprocessing and training specifics. The single preprocessing step used in all image-based detectors is the extraction of video frames and locating the facial regions. [5] studied the importance of preprocessing for face detection. They showed that current face detectors have a high false-positive rate, which can negatively affect classifiers’ training. To this end, they introduced a clustering-based approach to remove low confidence face detections. In addition to face detection, [31] adopted face tracking and face alignment using a Spatial Transformer Network before feeding into a CNN based feature extractor. In [22] a novel image representation has been presented that uses facial landmarks to generate a blending mask prior to training. This method tries to determine the boundary of manipulated facial regions created due to the blending of different faces.

Data augmentation plays a vital role in learning image features for convolutional architectures. Augmentations like Flipping, Rotation, Random crop, Jitter, Translation, Noise Injection, Color Transforms, etc. are being used widely for improving the performance of neural networks [33]. Random-Erasing is one such method where random patches from the input image are cut out or replaced with noise during training. It helps the network to learn or estimate features to match the neighboring information in the image. But there is a drawback to this approach in which randomly cutting out patches can result in the removal of essential object descriptors. For instance, if the upper half of an image containing the number ‘8’ is removed, it becomes a ‘6’. In such cases, random erasing is detrimental to the training process.

### 3. Data Exploration

A handful of DeepFake datasets have been published in recent years. Most of the datasets contain a collection of fake videos generated from real samples. Fig. 2 shows the ratio of unique face to the number of videos. Earlier datasets like UADFV and DF-TIMIT [20] had less than 500 unique videos and even fewer unique faces. The first large scale

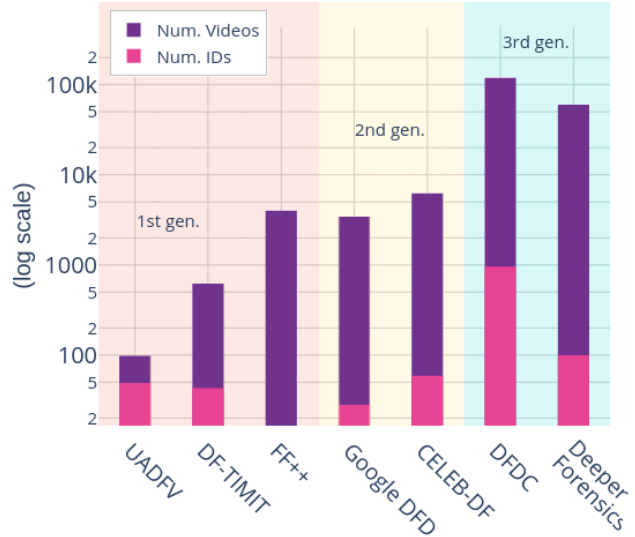


Figure 2: Comparison of current DeepFake datasets. Y-axis is shown in the log scale since the DFDC dataset is over an order of magnitude larger than any others. The number of videos and IDs is overlaid on the graph to show a comparison of the ratio of ID:Video. Datasets are divided into generations as given in [25].

face manipulation dataset was the FF+ dataset. Currently, DFDC is the largest dataset in terms of both the number of videos and faces. For our study, we chose to analyze the DFDC, FF++, and Celeb-DF datasets. These are the most popular datasets from their respective generations. We have explored each of those datasets, defined their shortcomings, and analyzed how it can affect DeepFake detection.

Table 1 shows a quantitative comparison of the selected three datasets. Here, the fake to real video ratio was calculated by averaging the fake video count generated from a single real source-video. Average videos per subject were calculated for every subject used either as a source or target for face swap or DeepFake. We can see a large imbalance in the data distribution from the count of videos per subject. For both the DFDC and Celeb-DF dataset, a single face appears in more than 120 videos on average. Considering the total number of frames in DFDC, a single face can be found in ~37.2k images.

#### 3.1. Face Clustering

The first step in face manipulation detection is to locate the face from a video frame. [30] showed that using only the facial region instead of the entire video frame improves detection performance. It also improves computation speed. Several different face detection and recognition models are used for this purpose, including, MTCNN [41], DLib [18], BlazeFace [2] etc. We used the MTCNN face detector to detect, and cluster faces from videos. We use facial clustering to aggregate similar faces and calculate the number of

Dataset	Frames per video	Real		Fake		Fakes per Real video	Unique Subjects	Videos per subject	No. of Clusters	Videos per cluster		
		#videos	#frames	#videos	#frames					min	max	avg
FF++*	~500	1,000	509.9k	4,000	1.7M	1:4	-	-	987	1	5	1
Celeb-DF	~382	590	225.4k	5,639	2.1M	1:10	59	1:182	45	5	57	14
DFDC	300	19,154	5.7M	99,992	29.9M	1:5	960	1:124	866	25	345	33

\* Subject data was not published

Table 1: Quantitative comparison and result of face clustering on various DeepFake datasets.

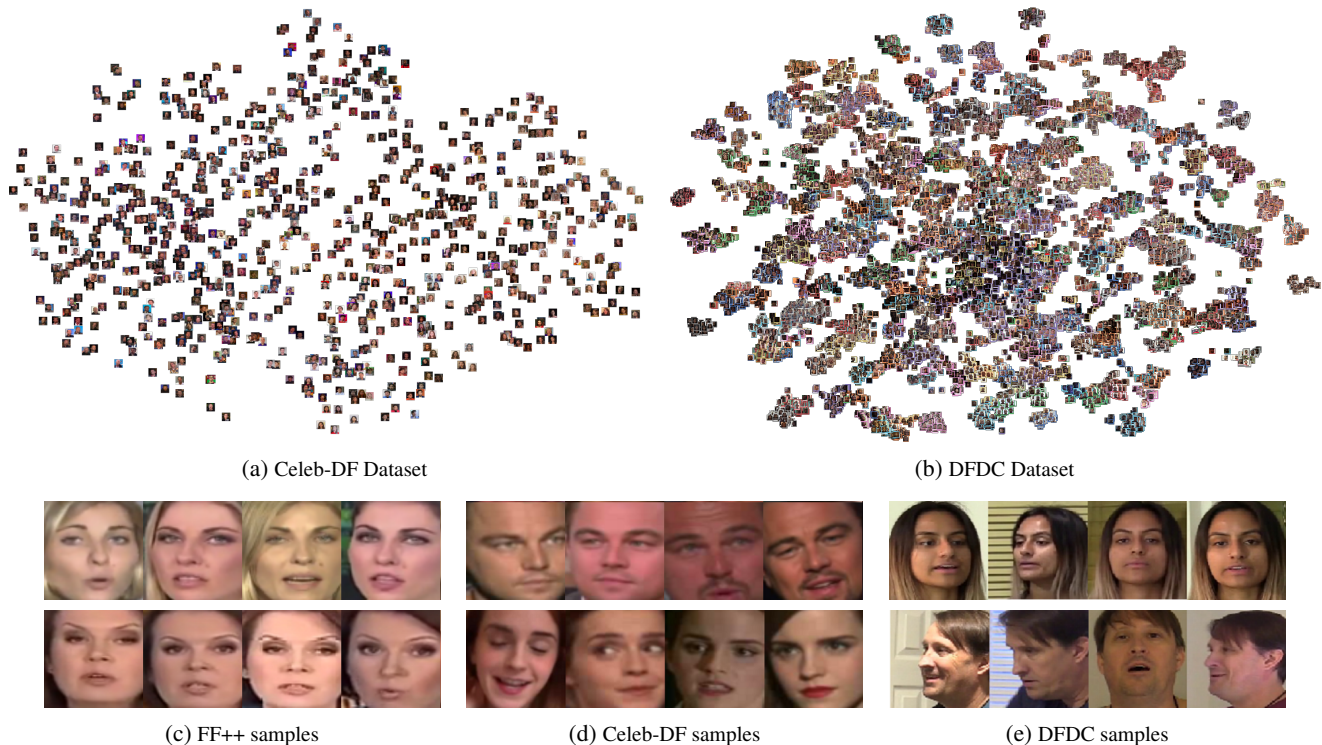


Figure 3: (a) & (b) shows the face clusters generated using DBSCAN for Celeb-DF and DFDC datasets, respectively. The density of the plots is representative of the number of videos in each dataset. (c), (d) & (e) are sample images from clusters for each dataset. All images in these examples are from separate videos.

unique videos per actor. This is done to get an insight into the actual distribution of faces in each dataset. It will also allow us to get a visual understanding of the data distribution and design our detection methods accordingly. Furthermore, we use the cluster information to prevent data leakage while separating validation and test data. It is further discussed in Sec. 5.1.

For clustering, we extracted all faces found in each *real* video and encoded each face to a 128 dimension vector using a one-shot CNN based encoder [14]. Experiments showed that increasing the dimension size did not improve cluster estimation by any significant margin, but rather greatly increased the computation time. Then we used Density-Based Spatial Clustering of Applications with Noise (DBSCAN) [13] with a Euclidean distance param-

eter to congregate the images into groups. Using the labeled dataset, we assigned all fake videos to their respective source face cluster. The clusters for Celeb-DF and DFDC are shown in Fig. 3. The 128 dimensional embeddings have been reduced to 2 dimensions using Principal Component Analysis (PCA) for better visualization. Table 1 contains the result of clustering. We found a single cluster that included all outlier images that could not fit in any other group during our experiments. The data shown in the table was calculated discarding this outlier group. We can see that the algorithm identified 45 clusters in the Celeb-DF dataset, among the 59 subjects reported by the authors. The difference is due to the difficulties of differentiating among the actors of similar race and gender under different lighting conditions as well as low-resolution images extracted

from compressed videos. For FF++, we identified 987 clusters. This means almost all real videos in this dataset have unique faces. For the DFDC dataset, we identified 866 separate clusters compared to 960 reported subjects. The variety of faces in DFDC is meager compared to this dataset’s enormous size.

The results of face clustering further verify our claim concerning the data imbalance and variation. FF++ has the most variation in faces, while DFDC has the least. It is well known that using the same data repeatedly for training can cause models to overfit. Our proposed method, Face-Cutout, helps solve this problem by generating multiple samples through dynamic face augmentation.

### 4. Face-Cutout

Landmark positions on a face are the locations of eyes, ears, nose, mouth, and jawline. DLib can uniquely identify 68 positions on the face, shown in Fig. 4. Each landmark position is annotated with a point ranging from [0 – 67]. Points in [36 – 41] and [42 – 47] denote left and right eyes respectively, while [27 – 35] is used for the nose and [48 – 67] for the mouth. We use these positions to calculate polygons for Face-Cutout. Before training, we generate a pixel-wise difference mask by calculating the Structural Similarity Index (SSIM) [43] between the frame of a real, and its corresponding fake video, as shown in Fig. 5. This difference mask contains 1 for manipulated pixels and 0 for real ones. This mask would later be used during training for calculating the Face-Cutout polygons out of non-manipulated regions.

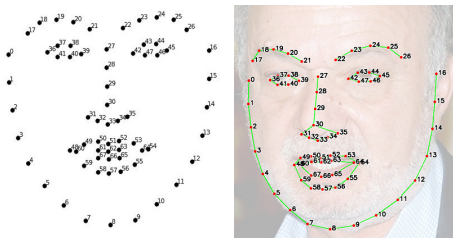


Figure 4: 68 landmark positions detected by DLib.

During training, Face-Cutout is performed with a certain probability. For a training image  $I$  in a mini-batch, the probability of it undergoing augmentation is  $p$ . We perform two types of cutout operations; 1) Sensory group removal and 2) Convex-hull removal. Sensory groups consist of the landmark groups of two eyes, nose and mouth. We randomly select one of these 3 groups and remove the maximum enclosing polygonal region  $I_c$  defined by the group’s points. Before removal,  $I_c$  is checked against the precalculated difference mask to evaluate the amount of fake artifact enclosed. From the difference mask, let  $C_o$  be the set of all pixels of value 1 contained within the cutout region and  $A$

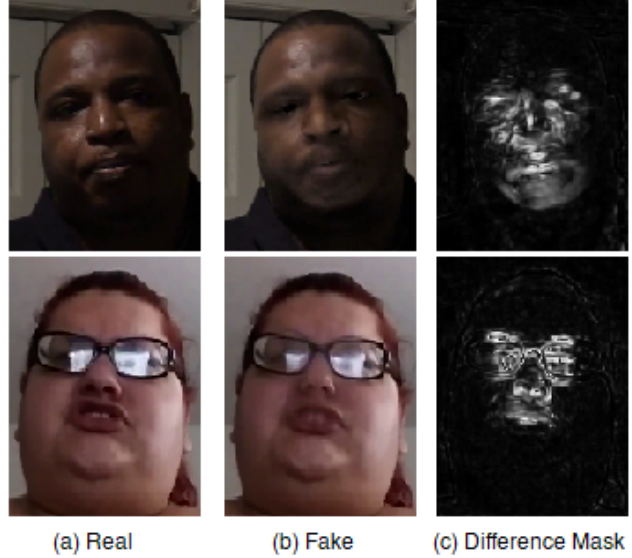


Figure 5: Face extracted from the frame of a real as well as its corresponding fake video. The difference mask shows the artifacts found by measuring the real and fake face’s pixel-wise difference.

the set of all 1 pixels in the entire image. Therefore, by definition  $C_o \subset A$ . The amount of envelop of the cutout region is denoted by  $\rho$  where,

$$\rho = \frac{|C_o|}{|A|} \tag{1}$$

The cutout region is selected to be augmented if  $\rho \leq \Gamma_h$  where  $\Gamma_h$  is a predefined threshold set to a default of 0.3. Since for real images, the difference mask does not contain any 1; therefore,  $|A|$  will always equal 0. So, they are augmented by the default polygon generated by the algorithm. Thus, the condition of  $\rho$  is only applied to fake images. By cutout, we mean that the pixel values of the selected region  $I_c$  are replaced with values from [0, 255].

For convex-hull removal, we select the landmark points [0 – 26] representing the facial boundary. We calculate the maximum polygon generated by points having the minimum envelop  $\rho$ . Points can be selected, 1) randomly from all boundary points, 2) as 8 or more contiguous points having the maximum polygonal area with minimum  $\rho$ , and 3) from one of 4 sub-polygons created by partitioning the outer polygon. Face-Cutout can also be combined with any existing image augmentation, like rotation, scaling, and color transforms. Fig. 6 shows some images generated using Face-Cutout.

### 5. Experimental Setup

We evaluate Face-Cutout on the three datasets, FF++, Celeb-DF, and DFDC. FF++ videos were used at 40% compression, and the other two datasets were used in their origi-

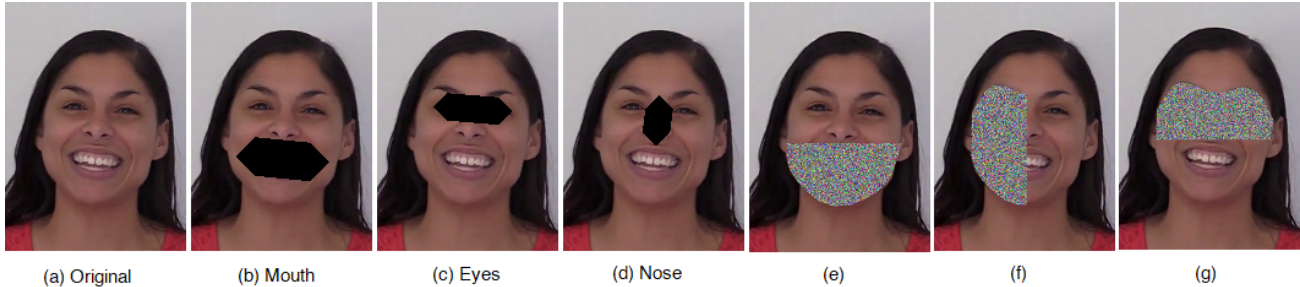


Figure 6: (a) The original face without any augmentations. (b), (c), (d) Three types of Sensory Group cutout of Mouth, Eyes and Nose. (e), (f), (g) Three random outputs for Convex-Hull Cutout. It also shows example of cutout fill with random values.

nal format. We also tested combining samples from all three datasets to present that our method reduces overfitting.

### 5.1. Test Set Selection

As explained earlier, DeepFake datasets are prone to overfitting due to a lack of face variation. While training and evaluating on these datasets, we need to prevent data leakage. For any machine learning system, the method for selecting the optimal validation data is a highly studied problem [39]. Train data should not be used for validation, as it would result in unrealistic model evaluation. For DeepFake data, leakage will occur if the same face is found in both train and test data.

To ensure data separation in our experiments, we created the validation and test set based on the face cluster data presented in Sec. 3.1. We separated the dataset such that all videos from the same cluster lie in the same set. We then trained our models with K-Fold Cross Validation with  $K = 10$  and used a single holdout set for the test. For the combined evaluation, we selected a subset of videos from each dataset based on clusters and used them together for the train and a separate set for the test. Since we do not have access to the private test data from Kaggle, our reported metrics might not correlate with the scores from the competition.

### 5.2. Model Selection

We selected two deep convolutional models; EfficientNet-B4 and XceptionNet. Both were initialized with pre-trained ImageNet weights. Additionally, EfficientNet-B4 was pre-trained using Noisy Student [37]. There are eight variants of the EfficientNet architecture based on their depth and number of parameters ranging from B0-B7. We chose B4 because of its lower parameter count and faster train time. The second model, XceptionNet, was introduced in [30, 11] as a baseline that achieves great results in forgery detection tasks.

### 5.3. Preprocessing

Each dataset consists of real and manipulated videos in *mp4* format. We used only facial regions from each frame for both training and inference. We picked every 10<sup>th</sup> frame from each video. First, we used MTCNN Face detector with thresholds (0.85, 0.95, 0.95) to get the face bounding boxes from each frame of the original videos. Face detection was only applied to the real faces to avoid any false positives. Since fake videos contain the same face in the same location for each frame, the face coordinates can be used to extract the fake face. After that, DLib was used on each of the extracted real faces to detect facial landmarks.

### 5.4. Training Setup

All the extracted images in each dataset were normalized using per-channel mean of (0.485, 0.456, 0.406), and standard deviation (0.229, 0.224, 0.225). Images were isotropically resized to  $224 \times 224$  with zero padding. Additional train augmentations were used, including Image Compression, Gaussian Noise, and Flipping, each with 10%-15% probability. Test and validation were done without any augmentation. The random generator was initialized with a seed of 777. We used Rectified Adam [26] optimizer with an initial learning rate of 0.001 and a weight decay of 0.0005. Learning rate scheduling was done using Reduction on Plateau by a factor of 0.25 and patience 2. All models were trained using Binary Cross-entropy Loss for 30 epochs and with Early stopping if no improvement was observed for consecutive 10 epochs. All experimentation was conducted with a training batch size of 40 on a system with an NVIDIA GTX 1080 Ti GPU and AMD Threadripper processor.

### 5.5. Evaluation Metrics

Since DeepFake datasets are heavily class imbalanced, accuracy is inefficient for measuring the model performance. We decided to use Area Under Curve (AUC) of ROC and Mean Average Precision (mAP) score for analyzing model performance. The AUC score summarizes the re-

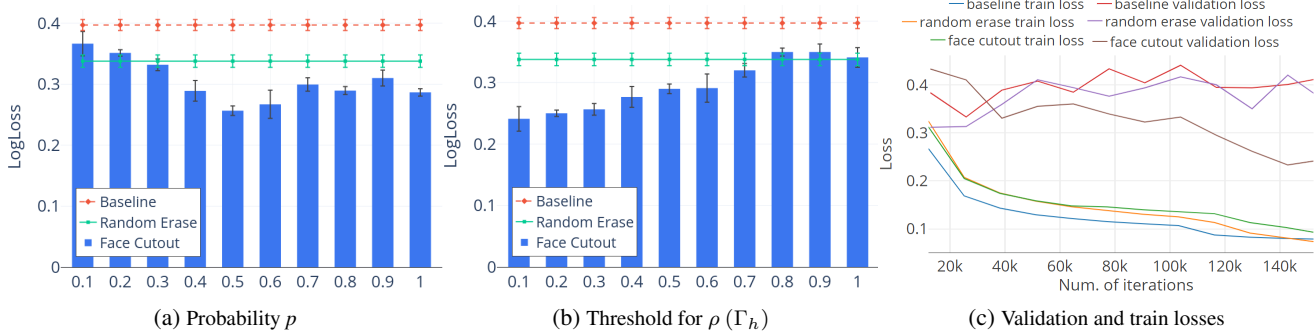


Figure 7: (a), (b) Test errors under different hyper-parameters. (c) Validation and train losses on optimized hyper-parameters.

Model	DFDC			FF++ (c40)			Celeb-DF		
	LogLoss	AUC(%)	mAP(%)	LogLoss	AUC(%)	mAP(%)	LogLoss	AUC(%)	mAP(%)
EfficientNet-B4 Baseline	0.397	87.11	96.02	0.215	95.59	97.9	0.104	98.75	98.66
EfficientNet-B4 + Random Erase	0.3178	91.01	97.14	0.239	95.01	93.15	<b>0.048</b>	<b>99.54</b>	<b>99.69</b>
EfficientNet-B4 + Face-Cutout	<b>0.2566</b>	<b>92.71</b>	<b>98.59</b>	<b>0.178</b>	<b>98.77</b>	<b>99.03</b>	0.065	99.21	99.53
Xception Baseline	0.5598	78.61	88.51	0.247	89.91	95.21	0.199	98.17	98.05
Xception + Random Erase	0.5011	<b>82.07</b>	90.80	0.287	88.42	95.04	0.098	99.20	99.17
Xception + Face-Cutout	<b>0.4718</b>	81.99	<b>91.32</b>	<b>0.195</b>	<b>96.73</b>	<b>96.06</b>	<b>0.096</b>	<b>99.39</b>	<b>99.44</b>

Table 2: Test results of the deep learning models trained separately using DFDC, FaceForensics++ and Celeb-DF datasets.

lation between the False Positive Rate (FPR) and True Positive Rate (TPR) of our binary classifier. Moreover, since the datasets contain a larger number of true negatives, mAP is more indicative of how a detection model will perform over a real distribution of images [10]. We also measure LogLoss on videos proposed by [10] as a metric for ranking DeepFake models. It was also used for ranking submissions in the DFDC Competition. LogLoss was calculated using equation (2),

$$\text{LogLoss} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (2)$$

where  $n$  is the number of videos being predicted,  $\hat{y}_i$  the predicted probability of the video being fake and  $y_i$  the true label of the video, 1 if fake, and 0 if real. The final prediction for each video was made as a mean of each selected frame’s predicted probabilities.

## 6. Experimental Results

**Classification on Independent Datasets:** The results of evaluating Face-Cutout on the individual datasets are shown in Table 2. We set  $\Gamma_h = 0.3$  and cutout fill 0. Results indicate that models trained with Face-Cutout have significant improvement over baseline and Random-Erase. Moreover, Random-Erasing performs worse than baseline

in some cases, as seen from the results of FF++. Face-Cutout shows improvement in both EfficientNet and Xception models with an increase of 0.46% to 7.58% AUC(%) from baseline and improvement of 15.2% to 35.3% test LogLoss across models and datasets. Moreover, it improved LogLoss by 19.25% from Random-Erasing in DFDC and was almost on par in Celeb-DF. From Fig. 7c we can see that both the baseline model and Random-Erasing overfit to the DFDC dataset considerably. The decreasing validation loss for Face-Cutout shows its effectiveness in reducing model overfitting.

**Classification on Combined Dataset:** The combined dataset was kept small to balance the source videos since DFDC is magnitudes larger than the other two datasets. Results from Table 3 show that Face-Cutout performs equally well in the combined test data and outperforms Random-Erasing. It achieves an 11.3% improvement in LogLoss from Random-Erasing using EfficientNet.

**Impact of Hyper-Parameters:** There are two hyper-parameters for Face-Cutout, the threshold for  $\rho$  ( $\Gamma_h$ ) and the cutout probability  $p$ . We experimented on the DFDC dataset using EfficientNet-B4 to measure the impact of these hyper-parameters. When evaluating one parameter, the other one was fixed. From Fig. 7, we can see that a  $p$  of around 0.5 achieves the lowest LogLoss and improves baseline results

Model	LogLoss	AUC(%)	mAP(%)
EffNet-B4 Baseline	0.2719	92.99	96.22
EffNet-B4 + Random Erase	0.2698	95.00	98.71
EffNet-B4 + Face-Cutout	<b>0.2393</b>	<b>95.44</b>	<b>98.94</b>
Xception Baseline	0.3177	90.15	98.01
Xception + Random Erase	0.2713	95.02	98.59
Xception + Face-Cutout	<b>0.2586</b>	<b>95.66</b>	<b>98.76</b>

Table 3: Test results on the combined dataset.

by a factor of 0.14. With lower  $p$ , the results are close to the baseline as the augmentation isn’t as effective when applied to a small amount of data. The results deviate from the optimum at higher probabilities but are still better than both baseline and Random-Erasing. For the threshold  $\Gamma_h$ , we can see an increasing trend in LogLoss. The threshold decides how much fake artifact we allow inside the cutout region. With higher threshold values, essential visual and fake information gets removed from the images, and the results are almost similar to Random-Erasing. A threshold of 0.1 achieves the best score, but it has a higher error deviation. We chose 0.3 to allow more augmentations. For Random-Erasing, we used parameters as suggested in [42].

**Impact of Cutout Pixel Value:** We evaluate Face-Cutout by erasing pixels in the selected region using three types of values: 1) each pixel is assigned a random value between  $[0, 255]$ , (termed as F-R); 2) all pixels are assigned with 0, (termed as F-0); 3) all pixels are assigned with 255, (termed as F-255). Table 4a shows that all erasing schemes outperform the baseline. Moreover, F-R and F-0 performs equally well, and both are superior to F-255.

		Cutout Type	LogLoss
		Baseline (No Cutout)	0.3970
(a)	Fill Type	Mouth	0.3115
	F-R	Eyes	0.3203
	F-0	Nose	0.3101
	F-255	Randomized Sensory (R-S)	0.2801
		Convex-Hull	0.2985
		R-S + Convex-Hull	<b>0.2566</b>
		(b)	

Table 4: (a) Test results of different fill values. (b) Results of different cutout types evaluated individually.

**Performance of Cutout Types:** We perform two types of Face-Cutout: 1) Sensory group and 2) Convex-Hull. We evaluated the impact of these two groups on the model performance. Table 4b shows that all cutout methods outper-

form the baseline. The use of any single sensory group cutout on all data gives approximately similar results, but choosing it randomly from the three groups improves the LogLoss by 0.03. Using convex-hull cutout on its own performs on par with the sensory group. The best result is achieved using randomized sensory and convex-hull augmentations together.

**CAM Visualization:** To visualize the effect and impact of our augmentation on training, and to verify whether it actually helps the models to better identify fake regions instead of memorizing faces, we visualize the GradCAM [32] output of the EfficientNet model in Fig. 8. The baseline model identified both the sample images correctly as fake with a 98% accuracy. But the CAM output shows the model highlighting the entire face, including arbitrary parts of the image. This means the model has been overfitted. However, the output produced by the same model when trained with Face-Cutout using the same parameters highlights only the fake portions of the face as seen in the difference mask. So, it has successfully reduced the model overfitting.

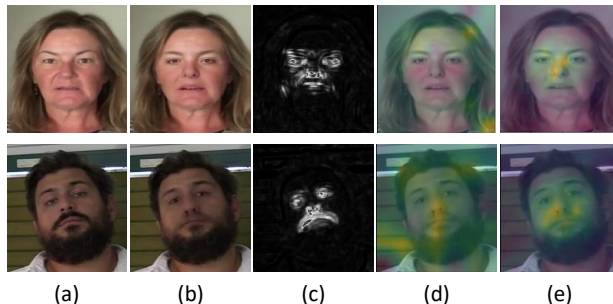


Figure 8: (a) Real face, (b) DeepFake, (c) SSIM difference mask showing fake pixels, (d) GradCAM output of a baseline model, (e) GradCAM output of Face-Cutout trained model.

## 7. Conclusion

In this paper, we showed the use of face clustering to identify the short-comings of DeepFake datasets and proposed Face-Cutout, an improvement over the Random-Erasing data augmentation for training convolutional neural networks. Our analysis provides significant directions to evaluate a DeepFake dataset, including preprocessing methods to mitigate overfitting. We hope these analyses will also help to create better datasets for DeepFake research. Our proposed augmentation policy enables the networks to overcome the problem of overfitting. It also helps the models to perform equally well on unseen data. We have shown experimentation results on 3 different datasets and using 2 separate models to prove that our scheme is not only independent of any single type of dataset but also it is not limited to any specific architecture. We have significantly



improved the DeepFake detection task using existing architectures, demonstrating our proposed method’s generalizability. Our data augmentation technique can be introduced into any existing DeepFake detection pipeline without any additional significant modifications. In the future, we wish to explore the use of this augmentation policy on more diverse face manipulation and forgery datasets. We also wish to explore its effects on video-based classification networks that use the spatial and temporal information of consecutive frames for DeepFake detection.

## References

- [1] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Mesonet: a compact facial video forgery detection network, 2018.
- [2] Valentin Bazarevsky, Yury Kartynnik, Andrey Vakunov, Karthik Raveendran, and Matthias Grundmann. Blazeface: Sub-millisecond neural face detection on mobile gpus, 2019.
- [3] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. What makes fake images detectable? understanding properties that generalize, 2020.
- [4] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now, 2018.
- [5] Polychronis Charitidis, Giorgos Kordopatis-Zilos, Symeon Papadopoulos, and Ioannis Kompatsiaris. A face preprocessing approach for improved deepfake detection, 2020.
- [6] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation, 2017.
- [7] François Chollet. Xception: Deep learning with depthwise separable convolutions, 2016.
- [8] Hao Dang, Feng Liu, Joel Stehouwer, Xiaoming Liu, and Anil Jain. On the detection of digital face manipulation, 2019.
- [9] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [10] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge dataset, 2020.
- [11] Brian Dolhansky, Russ Howes, Ben Pflaum, Nicole Baram, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) preview dataset, 2019.
- [12] Nicholas Dufour, Andrew Gully, Per Karlsson, Alexey Victor Vorbyov, Thomas Leung, Jeremiah Childs, and Christoph Bregler. Deepfakes detection dataset by google & jigsaw, 2019.
- [13] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 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*, KDD’96, page 226–231. AAAI Press, 1996.
- [14] Adam Geitgey, 2018. [https://github.com/ageitgey/face\\_recognition](https://github.com/ageitgey/face_recognition).
- [15] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [16] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks, 2016.
- [17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks, 2018.
- [18] Davis E. King. Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research*, 10(60):1755–1758, 2009.
- [19] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- [20] Pavel Korshunov and Sebastien Marcel. Deepfakes: a new threat to face recognition? assessment and detection, 2018.
- [21] Iryna Korshunova, Wenzhe Shi, Joni Dambre, and Lucas Theis. Fast face-swap using convolutional neural networks, 2016.
- [22] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection, 2020.
- [23] Yuezun Li, Ming-Ching Chang, and Siwei Lyu. In icu oculi: Exposing ai generated fake face videos by detecting eye blinking, 2018.
- [24] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts, 2018.
- [25] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics, 2019.
- [26] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond, 2019.
- [27] Yisroel Mirsky and Wenke Lee. The creation and detection of deepfakes: A survey, 2020.
- [28] Huy H. Nguyen, Fuming Fang, Junichi Yamagishi, and Isao Echizen. Multi-task learning for detecting and segmenting manipulated facial images and videos, 2019.
- [29] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. Use of a capsule network to detect fake images and videos, 2019.
- [30] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images, 2019.
- [31] Ekraam Sabir, Jiaxin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. Recurrent convolutional strategies for face manipulation detection in videos, 2019.
- [32] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2):336–359, Oct 2019.
- [33] Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6:1–48, 2019.
- [34] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks, 2019.

- [35] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2387–2395, 2016.
- [36] Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection, 2020.
- [37] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V. Le. Self-training with noisy student improves imagenet classification, 2019.
- [38] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks, 2016.
- [39] Yun Xu and Royston Goodacre. On splitting training and validation set: A comparative study of cross-validation, bootstrap and systematic sampling for estimating the generalization performance of supervised learning. *Journal of Analysis and Testing*, 2, 10 2018.
- [40] Xin Yang, Yuezun Li, and Siwei Lyu. Exposing deep fakes using inconsistent head poses, 2018.
- [41] 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, Oct 2016.
- [42] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation, 2017.
- [43] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- [44] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017.
- [45] Y. Zhuang and C. Hsu. Detecting generated image based on a coupled network with two-step pairwise learning. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 3212–3216, 2019.

## Appendix

### A1. Implementation Details

The basic algorithm for Face-Cutout is similar for both sensory group and convex-hull cutout. This is shown in Algorithm 1. In both cases, after selecting the specified point group, random polygons are generated by shuffling the points, and each time  $\rho$  is calculated to select the polygon with the minimum overlap. The main distinguishing factor is the selection of the polygon region for either sensory group or convex-hull cutout.

The calculation of  $\rho$  is dependent on the original image since that is used to determine the difference mask. For real faces, the original image is not present. If the original image is not supplied, then this calculation is omitted.

After selecting the points for the specified group for sensory group cutout, we draw a line using the terminal vertices. The line is drawn between the points: 1) 36, 45 for eyes, 2) 27, 33 for nose, and 3) 48, 54 for the mouth when using DLib. Our implementation used the MTCNN landmark detector that also generates the same terminal points for each landmark. We draw the line using `cv2.line` function of `python-opencv` package. Next we used `scipy.ndimage.binary-dilation` from the `scipy` package for 5 iterations with various weights to expand the line and generate a polygon. Each time  $\rho$  is measured for the resulting polygon, and we select the polygon with minimum  $\rho$ .

For convex-hull cutout, we use three polygon selections.

1. Firstly, 8 to 15 points are selected randomly from the 27 boundary points. A polygon is drawn using the selected points using the `skimage.draw.polygon` function from `skimage` package. For each polygon, we calculate it’s area using the equation 3,

$$\mathbf{A} = \frac{1}{2} \left| \sum_{i=1}^{n-1} x_i y_{i+1} + x_n y_1 - \sum_{i=1}^{n-1} x_{i+1} y_i - x_1 y_n \right| \quad (3)$$

where,  $(x, y)$  are the coordinates of each vertex point. We select the polygon with the maximum area having  $\rho \leq \Gamma_h$ .

2. We select a number  $i$ , randomly between 8 and 15. We select  $i$  consecutive points starting from 0 till 26. For example, for  $i = 11$  on first iteration we select points 0 – 10, on second iteration 1 – 11, on third iteration 2 – 12 and so on. Each time we draw the polygon using the selected points and calculate  $\rho$ , and it’s area. For all polygons that satisfy  $\rho \leq \Gamma_h$ , we select the one with the maximum area.

---

**Algorithm 1: Face-Cutout**

---

**Input** : Fake Image  $I$   
Original Image  $O$  (Optional)  
Erasing probability  $p$   
Landmark Coordinates  $L$   
Cutout threshold  $\Gamma_h$   
Cutout Fill  $F$

**Output:** Augmented Image  $I^*$

```
1 if  $O$  is not empty then
2   | Difference Mask  $M \leftarrow \text{diff}(I, O)$ 
3 end
4  $p_1 \leftarrow \text{Rand}(0, 1)$ ;
5  $I^* \leftarrow I$ ;
6 if  $p_1 \geq p$  then
7   | return  $I^*$ ;
8 else
9   |  $\text{choice} \leftarrow \text{Rand}[0, 1, 2, 3]$ 
10  | if  $\text{choice} = 0$  then  $\text{points} \leftarrow L[36 - 47]$  ;
11  | else if  $\text{choice} = 1$  then  $\text{points} \leftarrow L[48 - 67]$  ;
12  | else if  $\text{choice} = 2$  then  $\text{points} \leftarrow L[27 - 35]$  ;
13  | else  $\text{points} \leftarrow L[0 - 26]$  ;
14  | for  $i \leftarrow 0$  to 5 do
15  |   |  $I_c \leftarrow$  pixels inside the polygon created by
16  |   | randomly shuffling the vertices in
17  |   |  $\text{points}$ ;
18  |   | if  $O$  is empty then
19  |   |   |  $I(I_c) \leftarrow F$ ;
20  |   |   |  $I^* \leftarrow I$ ;
21  |   |   | return  $I^*$ ;
22  |   | end
23  |   |  $|C_0| \leftarrow$  count of 1's in  $I_c$ ;
24  |   |  $|A| \leftarrow$  count of 1's in  $M$ ;
25  |   |  $\rho \leftarrow \frac{|C_0|}{|A|}$ ;
26  |   | if  $\rho \leq \Gamma_h$  then
27  |   |   |  $I(I_c) \leftarrow F$ ;
28  |   |   |  $I^* \leftarrow I$ ;
29  |   |   | return  $I^*$ ;
30  |   | end
31  | end
32  | return  $I^*$ ;
33 end
```

---

3. We draw the polygon using all 27 boundary points. We find the centroid of this polygon using `skimage.measure.centroid` function. Then we split the entire polygon laterally and vertically through this centroid, which gives us four sub polygons. From each of these sub-polygons, we select the one with minimum  $\rho$ .

## A2. Data Split

For training and testing on each of the datasets, we used the number of videos given in Table 5. Since DFDC is magnitudes larger than FF++ and Celeb-DF, for selecting the combined dataset, we used a smaller number of videos from DFDC.

Dataset		Individual data split		Combined data split	
		Real	Fake	Real	Fake
FF++	Train	789	2,411	600	2,000
	Test	100	300	200	173
Celeb-DF	Train	449	4,436	300	2,000
	Test	60	367	209	237
DFDC	Train	16,148	22,510	5,000	7,000
	Test	2,180	10,371	2,000	4,000

Table 5: Number of videos in the train and test set for individual testing and combined data testing.