

Deepfake Detection: A Comparative Analysis

SOHAIL AHMED KHAN, SFI-MediaFutures, University of Bergen, Norway
DUC-TIEN DANG-NGUYEN, University of Bergen, Norway

This paper presents a comprehensive comparative analysis of supervised and self-supervised models for deepfake detection. We evaluate eight supervised deep learning architectures and two transformer-based models pre-trained using self-supervised strategies (DINO, CLIP) on four benchmarks (FakeAVCeleb, CelebDF-V2, DFDC, and FaceForensics++). Our analysis includes intra-dataset and inter-dataset evaluations, examining the best performing models, generalisation capabilities, and impact of augmentations. We also investigate the trade-off between model size and performance. Our main goal is to provide insights into the effectiveness of different deep learning architectures (transformers, CNNs), training strategies (supervised, self-supervised), and deepfake detection benchmarks. These insights can help guide the development of more accurate and reliable deepfake detection systems, which are crucial in mitigating the harmful impact of deepfakes on individuals and society.

Additional Key Words and Phrases: deepfakes; visual content verification; convolutional neural networks; transformers; video processing.

ACM Reference Format:

Sohail Ahmed Khan and Duc-Tien Dang-Nguyen. 2023. Deepfake Detection: A Comparative Analysis. 1, 1 (August 2023), 28 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Deepfakes, or deepfake media, are digital media that have been generated or modified using deep learning algorithms. They have gained notoriety in recent years due to their potential to manipulate and deceive using artificial intelligence (AI) techniques. While deepfakes can be used for harmless or even humorous purposes, they can also pose a serious threat when used for malicious purposes such as creating convincing fake media to manipulate public opinion, influence elections, or incite violence.

The research community has been working on proposing AI-based automated systems to detect deepfakes. However, one of the major challenges in detecting deepfakes is that the deepfake generation systems are constantly evolving and improving. With the availability of cheap compute resources, and open-source software, it is becoming easier (even for people with limited technical knowledge and expertise) to create realistic deepfakes that are harder to distinguish from the real content. In addition to that, the deepfake detection, and generation is like a cat-and-mouse game [12], where the researchers propose detection tools by exploiting certain shortcomings of the generation systems. Soon after the release of the detection systems, the generation techniques are reinforced and made undetectable for the previously proposed detection systems by overcoming the exploited vulnerabilities. For example, in [25] researchers proposed a deepfake detector which exploited eye blinking as a cue to detect deepfake media (they found that deepfake faces don't blink

Authors' addresses: [Sohail Ahmed Khan](mailto:sohail.khan@uib.no), sohail.khan@uib.no, SFI-MediaFutures, University of Bergen, Norway; [Duc-Tien Dang-Nguyen](mailto:ductien.dangnguyen@uib.no), University of Bergen, Norway, ductien.dangnguyen@uib.no.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.

XXXX-XXXX/2023/8-ART \$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

eyes). Soon after they released their study, the newer deepfake generation algorithms generated videos which blinked eyes, and thus making the detection system useless.

Another prominent problem of the available deepfake detection algorithms is the lack of generalisation capability. This means that the detection systems work excellently on detection deepfakes coming from the same data distribution as the training data used to train these systems. However, when exposed to deepfakes generated using different generation systems than the one which generated training samples, the detectors fail to achieve similar performance. Numerous studies have been proposed in the past which propose to employ different novel strategies to develop detection systems, however, all of the proposed studies have one thing in common, and that is, poor generalisation capability of the models on unseen data.

To address these issues, this study aims to provide insights into the challenge of detecting deepfake media by comparing multiple deep learning models and deepfake detection benchmarks. Specifically, we evaluate several different well-known image and video recognition architectures for their effectiveness in detecting deepfakes. Our primary objective is to identify which of these models perform well on unseen data as compared to other participating models.

To achieve this, we train all participating models on four deepfake detection datasets, including a newly released dataset, and evaluate them in both intra-dataset¹ and inter-dataset² configurations (see Figure 1). Additionally, we evaluate the difficulty level of each benchmark and investigate whether a more challenging benchmark leads to better generalisation performance on unseen data. To this end, we train participating models on all four datasets twice: first, without any image augmentations, and then with various image augmentations to improve their performance.

We also analyse self-supervised Vision Transformer (ViT) architectures pre-trained using two well-known strategies: **DINO** [5] and **CLIP** [31]. To study these models, we use self-supervised ViT-Base models as feature extractors and train a classification head on top of them. It is important to note that we only train the classification head and freeze the weights of the feature extractors to avoid backpropagating gradients through them.

Overall, our study aims to answer several questions, such as which model has the highest generalisation capability on unseen data, which dataset is most challenging for the models to learn, which dataset enables the models to achieve the best generalisation capability on unseen data, and which of the participating models and architectures are most successful for detecting deepfakes.

This next parts of this paper is organised as follows. In Section 2 we present a brief literature review on the topic of deepfake detection. Section 3 presents the proposed framework. In Section 4 we present the results and discussion of our findings, and finally Section 5 concludes this study by summarising our analysis, and presents future research direction.

2 LITERATURE REVIEW

Since recently quite a large number of research studies focused on deepfake media detection have been proposed. Most of the proposed studies employ CNN models to detect deepfake media. The proposed studies also employ different strategies e.g., novel augmentation techniques, ensemble models, behavioral features, multimodal features, temporal features along with spatial information, recurrent networks, transformer models etc to detect deepfake images/videos while trying to increase the models' generalisation capabilities. Below we present some well-known, as well as some of the recently proposed deepfake detection studies.

In one of the earliest studies on deepfake media detection, Afchar *et al.* proposed two different CNN models namely (1) Meso-4, and (2) MesoInception-4 [3]. Both of the proposed CNN networks

¹models trained and evaluated on the same dataset

²models trained on one dataset and evaluated on another dataset

were comprised of a very small number of layers which focused on mesoscopic image details. Authors tested their proposed models on one of the available deepfake detection benchmark. In addition to that, authors also collected a custom dataset and tested the models on it as well achieving excellent results on both participating datasets.

In [33] Sabir *et al.* proposed to detect deepfake media using a novel recurrent convolutional network. Authors used DenseNet CNN, and combined it with a gated recurrent neural network (RNN) to learn temporal features along with spatial features. The motivation was to detect inconsistencies within neighbouring frames of a video. Authors evaluated their model on the widely known FaceForensics++ [32] deepfake detection benchmark showing promising results.

Rossler *et al.* in [32] proposed a deepfake detection benchmark, called FaceForensics++. Along with the benchmark, authors proposed a simple CNN based deepfake detection technique using XceptionNet [7]. Authors trained and evaluated the simple XceptionNet on their FaceForensics++ deepfake detection benchmark. They reported excellent performance scores on high-quality version of the four subsets of the FaceForensics++ dataset [32], however, lost performance when evaluated on low-quality videos.

In [28] Nguyen *et al.* proposed to employ capsule networks for deepfake detection. The proposed technique was the first of its kind which proposed to employ capsule networks in contrast to most of the other techniques which proposed to employ CNN models at that time. The capsule networks based detection technique was evaluated on four different deepfake detection datasets comprising of a wide variety of fake videos and images. The authors reported excellent evaluation statistics of their proposed technique in comparison to other deepfake detection techniques.

Ciftci *et al.* in [8], developed a novel CNN and SVM based deepfake media detection model which was trained on biological signals (i.e., photoplethysmography or PPG signals). The CNN and SVM models make their individual predictions which were then fused together in order to get a final classification score. This deepfake detection model achieved promising results when tested on a number of different deepfake detection benchmarks including, CelebDF[26], Face Forensics, and Face Forensics++ [32] datasets.

Zhu *et al.* in [47] proposed a deepfake detection system which employed 3D face decomposition features in order to detect deepfakes. Authors showed that by merging the 3D identity texture and direct light features significantly improved the detection performance while also making the model to generalise well on unseen data when evaluated in cross-dataset setting. In this study, authors also employed the XceptionNet CNN architecture for feature extraction. Both a face cropped image, and its associated 3D detail were used to train their deepfake detection model. They also carried out an in-depth analysis of several different feature fusion strategies. The proposed model was trained on the FaceForensics++ [32] benchmark, and evaluated on (1) FaceForensics++, (2) Google Deepfake Detection Dataset, and (3) DFDC [11] dataset. Promising evaluation statistics were reported for all of the three participating datasets, while depicting the generalisation capability of the model when compared to the previously proposed deepfake detection systems.

In [23], Khan *et al.*, proposed to employ transformer architecture for the task of deepfake media detection. Authors proposed a novel video based model for deepfake detection which was trained on 3d face features as well as standard cropped face images. Authors also showed that their proposed model was capable of incrementally learning from new data without catastrophically forgetting what it was trained on earlier. Authors evaluated their models on different widely used deepfake detection benchmarks including FaceForensics++, DFDC, DFD and showed that their proposed models achieved excellent results on all of the participating datasets.

[41] introduce a Multi-modal Multi-scale TRansformer (M2TR) model, which processes patches of multiple sizes to identify local abnormalities in a given image at multiple different spatial levels. M2TR also utilises the frequency domain information along with RGB information using a

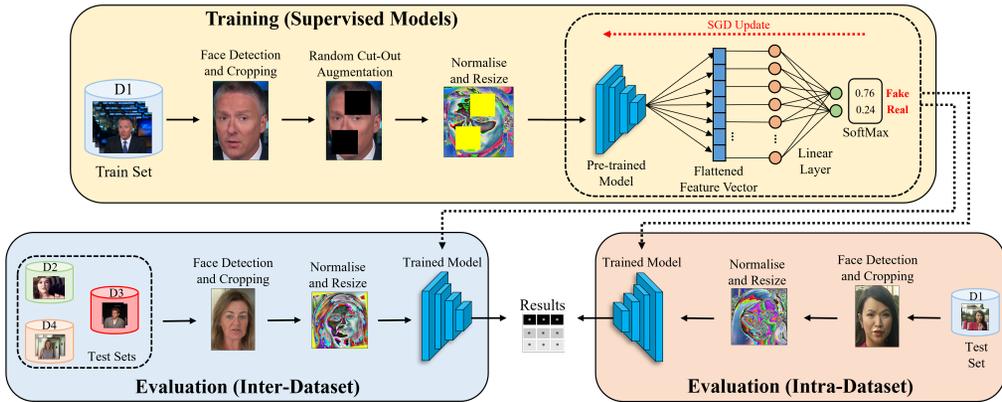


Fig. 1. The proposed framework. The process involves several steps, starting with the extraction and cropping of face frames from videos, followed by augmentation, normalisation, and resizing. The pre-trained models are then used as feature extractors, with a new classification head (linear layer) added on top for supervised models. During training, the weights of both the feature extractor and the classification head are updated for supervised models, while only the newly added classification head is updated for self-supervised models. The models are evaluated through both intra-dataset and inter-dataset evaluations to test their performance and generalisation capabilities. For image models, the input data is a single cropped face image, while for video models, it is a tensor containing eight consecutive cropped face images from a given video.

sophisticated cross-modality information fusion block to detect forgery related artifacts in a better way. Through extensive experiments authors establish the effectiveness of M2TR, and show their model outperforms SOTA Deepfake detection models by acceptable margins.

[9] propose a video deepfake detection model by employ a hybrid transformer architecture. Authors used an EfficientNet-B0 as feature extractor. The extracted features were then used to train two different types of Vision Transformer models in their study, e.g., (1) Efficient ViT, and (2) Convolutional Cross ViT. Through experimentation, authors established that the model comprising of EfficientNet-B0 feature extractor and Convolutional Cross ViT achieved the best performance scores as compared to other models that they tested.

In [46], an Interpretable Spatial-Temporal Video Transformer (ISTVT) for deepfake detection was proposed. The proposed model incorporates a novel decomposed spatio-temporal self-attention as well as a self-subtract mechanism to learn forgery related spatial artifacts and temporal inconsistencies. ISTVT can be also visualise the discriminative regions for both spatial and temporal dimensions by using the relevance propagation algorithm [46]. Extensive experiments on large-scale datasets were conducted, showing a strong performance of ISTVT both in intra-dataset and inter-dataset deepfake detection establishing the effectiveness and robustness of proposed model.

Through this literature review it becomes apparent that the research community actively employs deep learning based models along with other techniques to try develop robust and efficient deepfake detectors. However, while carefully reading the research studies it also becomes noticeable that the models perform poorly on unseen data. Also, there is a lack of comparative studies which aim to identify which specific family of deep learning architectures is better than the others in detecting deepfakes. Also, it is a bit difficult to make sense of the capability of datasets in providing the generalisation capability to the models which can help them classify unseen data in a better way. To address this, in this study we propose to employ some of the most frequently used architectures (EfficientNets, Xception, Vision Transformers) in the literature of deepfake detection. We also

employ widely known datasets for experimentation, and try to find out the datasets offering best generalisation capabilities to the models. We also analyse somewhat understudied approaches for deepfake detection i.e., we train and evaluate the performance of self-supervised models and compare their performance with supervised models.

3 THE PROPOSED FRAMEWORK

The workflow followed in this study for training and evaluating the models is illustrated in Figure 1. On top we show the training pipeline where we start by extracting and cropping faces from videos. The cropped face frames are then augmented, normalised, and resized before being fed to the model for training. We load pre-trained models as feature extractors, i.e., we remove the last layer from the loaded models and add a new classification head (linear layer) on top. For supervised models, during training we update weights of both feature extractor as well as the classification head. In case of self-supervised models we only update weights of the newly added classification head, and keep the weights of feature extractor frozen.

For intra-dataset evaluation we evaluate models on the same dataset (test set) it was trained on, e.g., model trained on dataset D1 is evaluated on the test set of D1. The main goal of intra-dataset evaluation is to find out which model performs better than other participating models on each of the dataset. Besides this, intra-dataset evaluation provides an insight about the dataset which is more challenging for the models to learn, and which dataset is the least challenging to learn.

During inter-dataset evaluation, we evaluate models trained on one dataset on each of the other three datasets, e.g., model trained on dataset D1 is evaluated on D2, D3, and D4 datasets. The goal of inter-dataset evaluation is to study the generalisation capabilities of models as well as to find out how good the training dataset is at providing that generalisation capability.

The input data for training and evaluating models image models is a single face cropped image ([3, 224, 224]), whereas, input data for training and evaluating video models is a tensor containing 8 consecutive face cropped images ([8, 3, 224, 224]) from any given video.

3.1 Datasets

In this study we train and evaluate several different deep learning models on four deepfake detection datasets/benchmarks: FakeAVCeleb [22], CelebDF-V2 [26], DFDC [11], and FaceForensics++ [32]. All of the four datasets comprise of real and fake videos, where fake videos are generated using different deepfake generation methods. In upcoming sections, we present a brief description of these datasets.

FaceForensics++ [32] is one of the most widely studied deepfake detection benchmarks. FaceForensics++ comprises of 1000 real video sequences (mostly from YouTube) of mostly frontal faces and without any occlusions. These real videos were then manipulated using four different face manipulation methods: (1) FaceSwap [2], (2) Deepfakes [1], (3) Face2Face [39], and (4) NeuralTextures [38] to have four subsets. Each subset comprises of 1000 videos each. In total, the dataset contains 5000 videos, i.e., 1000 real and 4000 fake videos. FaceForensics++ offers 3 different qualities of data, (1) Raw, (2) High-Quality and (3) Low-Quality. In our study, we experimented the high-quality videos.

FaceSwap and Deepfakes subset contains videos generated using what is called, face swapping. As the name suggests, face of the target person is replaced with the face of source person and results in transferring the identity of the source person onto the target. Face2Face and NeuralTextures subsets are generated by a different process called, face re-enactment. In contrast to face swapping, face re-enactment swaps the faces of source and target, however, keeps the original identity of the target face (see Figure 2).

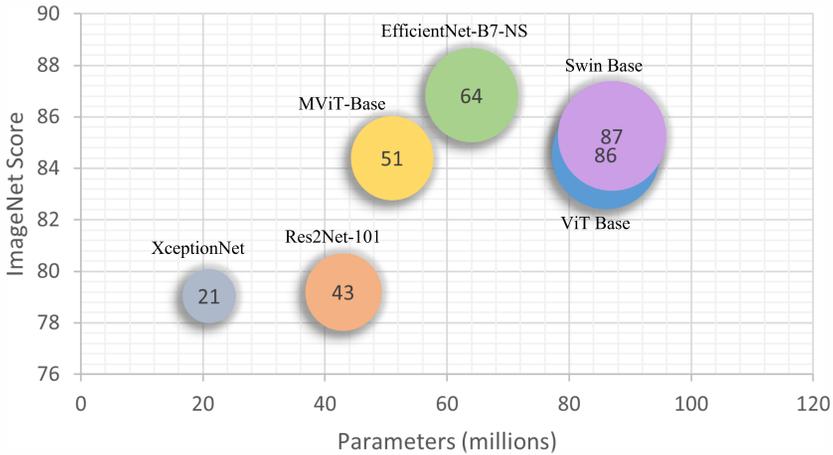


Fig. 2. Model size and its performance (Top-1 accuracy) on ImageNet [10].

Deepfake Detection Challenge (DFDC) [11] comprises of around 128k videos, out of which, around 104k are fake. Similar to the FaceForensics++ dataset, DFDC dataset also comprises of videos generated using more than one face manipulation algorithms. Five different methods were employed to generate fake videos, namely, (1) Deepfake Autoencoder[11], (2) MM/NN[18], (3) NTH[45], (4) FSGAN[29], and (5) StyleGAN[21]. In addition to these, a random selection of videos also underwent a simple sharpening post-processing operation which increases the videos' perceptual quality. Unlike FaceForensics++ dataset, the DFDC dataset also contains videos having undergone audio-swapping. However, in this study we do not use audio features to train and evaluate our models.

Since DFDC dataset is huge and we have limited resources, we only use a subset of the full dataset to train and evaluate our models. For training we use roughly around 19.5K (around 16.5K fake, and 3.1K real) videos from which we get 100k face cropped images (50k real, and 50k fake). We use 20k images as validation set. For testing the models we use 4000 image frames randomly selected from 3.5K (3.2K fake, and 0.3K real) videos.

CelebDF-V2 [26] contains 5639 fake, and 590 real videos. The real videos are collected from YouTube, and contain interview videos of 59 celebrities having diverse ethnic backgrounds, genders, age groups. CelebDF-V2 dataset comprises of fake videos generated using Encoder-Decoder models. Post-processing operations are also employed to circumvent color mismatch, temporal flickering, and inaccurate face masks.

FakeAVCeleb [22] is the most recently proposed deepfake detection dataset. FakeAVCeleb dataset contains 19500 fake, and 500 real videos. This dataset also contains audio modality, and manipulates audio as well as video content to generate deepfake videos. For video manipulation, FaceSwap[24], and FSGAN[29] algorithms are used. For audio manipulation, a real-time voice cloning tool called SV2TTS[19], and Wav2Lip[30] are used. The dataset is divided into 4 subsets, i.e., (1) FakeVideo/FakeAudio, (2) RealVideo/RealAudio, (3) FakeVideo/RealAudio, and (4) RealVideo/FakeAudio.

In this study, we only employ 2 of the mentioned subsets to train our models, i.e., (1) FakeVideo/FakeAudio, and (2) RealVideo/RealAudio.

Table 1. The amount of real/fake images used to train, validate, and test our image models.

Train/Test Data						
Dataset	Train		Validation		Test	
	Real	Fake	Real	Fake	Real	Fake
FakeAVCeleb [22]	47,099	45,912	9,301	9,301	2,000	2,000
CelebDF-V2 [26]	50,000	50,000	10,000	10,000	1,000	1,000
DFDC [11]	50000	50,000	10,000	10,000	2,000	2,000
FaceForensics++ [32]	50000	50,000	10,000	10,000	2,000	2,000

3.2 Dataset Preparation

Data preparation process was quite lengthy process as (1) the datasets are quite big, and (2) some of the chosen datasets do not come with helpful dataset preparation instructions, e.g., FakeAVCeleb does not comes with train/validation/test splits. We thus have to manually devise strategy in order to properly divide dataset into train/validation/test sets, while making sure that we do not encounter same identity in more than one splits.

Besides this, all of the datasets are unbalanced, i.e., contain more "fake" videos as compared to the "real" videos. We also made sure that we extract faces from videos in a way that the resulting face cropped image datasets are balanced, and contain at least a frame from all of the videos we chose to train/evaluate our models.

3.3 Preprocessing and Augmentations

We employ two different strategies to train our models in this study. First, we train models without using any image augmentations, and second, we train the models using different randomly selected image augmentations, including, horizontal flips, affine transformations, and random cut-out augmentations. All the face cropped images are then normalised according to the same strategy used in order to pre-train models on ImageNet dataset. We use ImgAug [20] library for the augmentations.

3.4 Models

We choose to experiment with six image recognition models trained using supervised strategy, three of them are CNNs and the rest three are transformer based models. We also evaluate two variants of transformer models trained using self-supervised strategies including (1) DINO [5], and (2) CLIP [31]. Besides the image classification models, we also train and evaluate two different video classification models, (1) ResNet-3D [16], which is a CNN model for video classification, and (2) TimeSformer [4], which is a transformer model for video classification.

We select models based on their performance on the ImageNet dataset [10], the number of parameters of the model, and for some models such as the Xception [7], we consider their previously reported performance on the task of deepfake detection.

3.4.1 Image Models. Deepfake detection task is typically considered as an image classification task, where a deep learning model is trained and evaluated on images separately (image-by-image), which is in contrast to the video based deepfake detection where the models are trained and evaluated on consecutive video frames to better detect the temporal inconsistencies present between different frames of the video along with spatial cues. Whereas the image based deepfake detection models only focus on learning to detect the spatial inconsistencies present in the images.

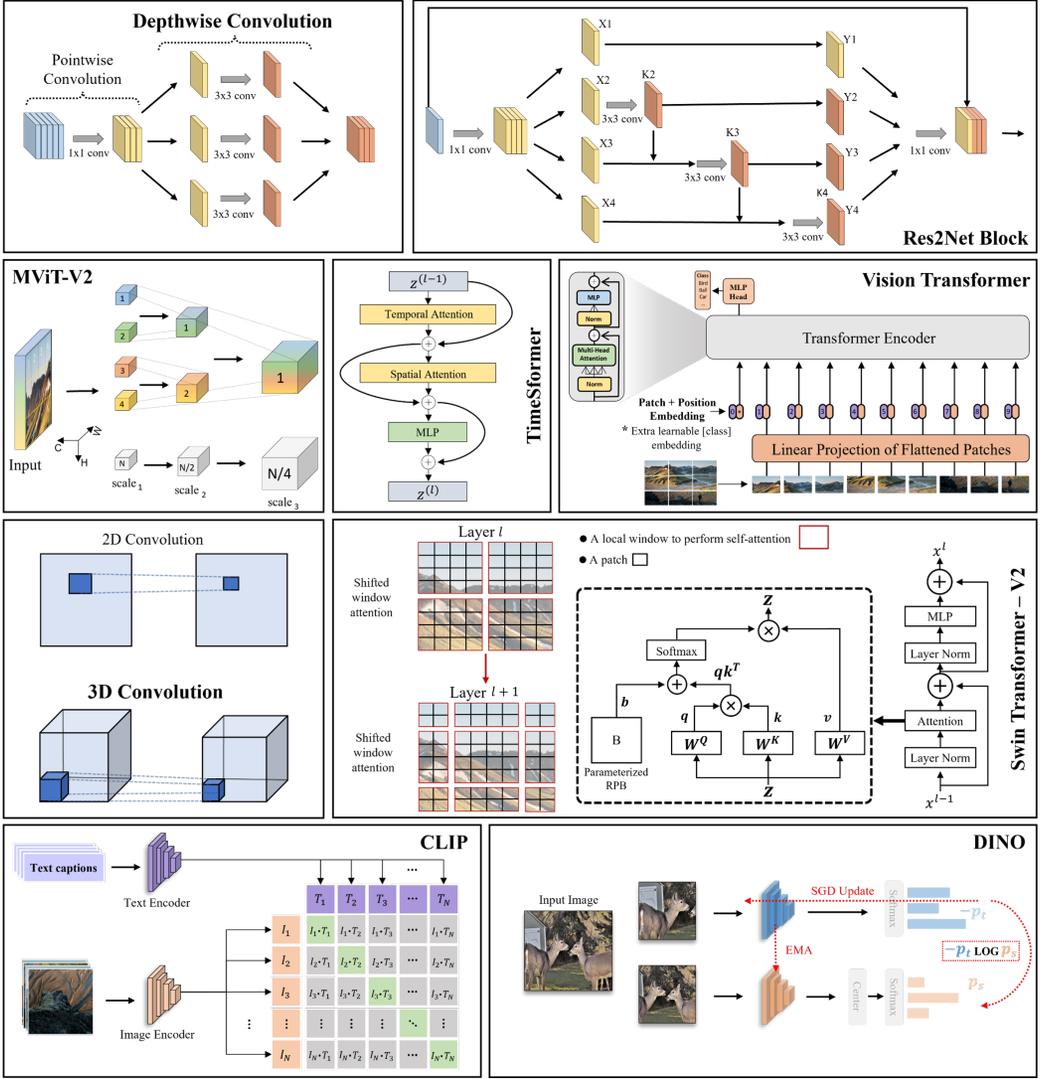


Fig. 3. Visual representation of the models used for analysis in this study. Due to space limitations, only basic, key concepts for each model are illustrated instead of the whole model. For optimal understanding of the essential components of each model, we recommend viewing this figure in color and at a higher magnification.

• **Xception** [7] is a convolutional neural network (CNN) architecture built upon the Inception architecture [36], but proposes to use depth-wise separable convolutions instead of the traditional Inception modules. Xception has a smaller number of trainable parameters as compared to some of the other widely used deep CNN models, however, it still shows comparable performance to other models having more parameters on the ImageNet benchmark [10]. Due to the smaller number of parameters, Xception is less prone to over-fitting on unseen data, as well as performs smaller number of computations, thus resulting in more efficient models. Depth-wise convolution is depicted in the top left corner of Figure 3. Besides good performance on ImageNet benchmark, Xception was also shown to achieve excellent results on deepfake detection task in studies conducted in the

past [32, 47]. Based on the results achieved by this architecture in the past on deepfake detection we opt to analyse this architecture in this study as well.

- **Res2Net-101** [15] is a convolutional neural network architecture. The main motivation behind Res2Net is to improve upon the popular ResNet architecture [17] by introducing a new type of building block called the "Res2Net Block", replacing the traditional bottleneck residual blocks of ResNet. The Res2Net architecture represents multi-scale features at a granular level and increases the range of receptive fields of each of the network layers. This results in a more efficient and powerful network that can achieve better performance on a wide range of computer vision tasks, such as, image classification, segmentation and object detection [15]. The proposed Res2Net block can be easily incorporated into other state-of-the-art backbone CNN models, e.g., ResNet[17], DLA[44], BigLittleNet[6], and ResNeXt[43]. Res2Net block is illustrated in the top right corner of Figure 3. We employ Res2Net-101 in this study to analyse whether multi-scale CNN features improve deepfake detection performance, and if so, does it also improve cross-dataset performance (generalisation capability)?

- **EfficientNet-B7** [37] is a convolutional neural network (CNN) architecture. The main idea behind the EfficientNet architecture is to increase the efficiency of convolutional neural networks by scaling the model's architecture, and parameters in a systematic manner. The authors proposed a new scaling technique that uniformly scales the depth, width, and resolution using a straightforward yet highly effective compound coefficient. In simple words, instead of arbitrarily scaling up model width, depth or resolution, the compound scaling strategy uniformly scales each dimension with a certain fixed set of scaling coefficients. Using this method the authors proposed seven different models of various scales [37]. The EfficientNet architecture achieves SoTA performance on a number of image classification benchmarks while being more computationally efficient than other architectures such as ResNet and Inception [37]. As is the case with Xception, a variant of EfficientNet architecture, specifically the EfficientNet-B7 architecture was also shown to perform excellently on deepfake detection task. The winning solution of the Google sponsored Deepfake Detection Challenge (DFDC) was also based on these EfficientNet-B7 models [11]. We thus choose to study this model in this paper.

- **Vision Transformer (ViT Base)** [40] is a class of neural network architectures based on the transformer architecture, which was initially designed for natural language processing tasks. In the context of computer vision, the Vision Transformer or simply ViT was the first transformer based architecture to be made available for image classification task [13]. It uses the self-attention mechanism to process visual data. The ViT uses a simple yet powerful approach, which is to divide the image into small patches and feed them into a transformer model at once. The small patches are then assigned positional embeddings in order to have an idea of the position of the image patch in the original image. A classification token is then inserted at the start of this input, which is then processed by the transformer encoder (similar to the encoder used in text related transformer models). The model learns to attend to different patches of the image at the same time to make predictions. By doing this, the network tends to better capture the context and relationships between different parts of the image, leading to comparable performance as compared to the SOTA CNN models on the ImageNet dataset after being trained on huge datasets, such as, the ImageNet-21k [10] or the JFT-300M [35] image datasets. ViT architecture is presented in Figure 3, second row on the right side. In this study, we train and evaluate the base version of vision transformer (ViT-Base) model on the task of deepfake detection and compare its performance with other participating models.

- **Swin Transformer (Swin Base)** [27] is a class of Vision Transformer models. It generates hierarchical feature maps by combining image patches in deeper layers. It is computationally efficient as compared to other vision transformer models, as it only performs self-attention within

each local window, resulting in linear computation complexity depending on the size of the input image. In contrast, vanilla Vision Transformers produce feature maps of a single low resolution and have quadratic computation complexity to the size of the input image, due to global self-attention computation. Swin Transformer achieves comparable performance when compared with other SoTA image classification models such as the EfficientNets[37]. Besides image classification, Swin Transformers also perform well on tasks such as image segmentation, object detection. Figure 3, third row on the right illustrates window generation, and attention calculation of Swin transformer. Because of the excellent performance swin transformer achieve on ImageNet, we use it for the task of deepfake detection, and try to study how it performs as compared to other participating models.

- **Multiscale Vision Transformer (MViT-V2 Base)** [14] is another class of vision transformer model. Unlike traditional vision transformers, the MViTs have multiple stages that vary in both channel capacity and resolution. These stages create a hierarchical pyramid of features, where initial shallow layers focus on capturing low-level visual information with high spatial resolution, while deeper layers extract complex, high-dimensional features at a coarser spatial resolution. This approach allows the network to capture the context and relationships between different parts of the image in a better way, which results in improved performance on a broad range of computer vision tasks including image classification, image segmentation. A broad overview of the architecture of MViT is shown on the left side, second row in Figure 3. Since MViTs are relatively new and achieve excellent performance on different vision tasks, we employ these in our study to analyse how well they perform on the task of deepfake detection.

- **DINO**[5] is a simple self-supervised training method, which is interpreted as a form of self-*DI*stillation with *NO* labels. The authors adapted self-supervised methods to train ViT (vision transformer) [13] architecture, and ViTs trained using supervised strategies. The authors make the following observations in their study, i.e., (1) self-supervised ViT features incorporate explicit information useful for computer vision tasks such as semantic segmentation, which does not come along as evidently with supervised ViTs, and also not with CNNs; (2) self-supervised ViT features are also shown to achieve excellent performance when tested as k-NN classifiers, attaining 78.3% top-1 on ImageNet with a ViT-small architecture. For more details about the strategy, please see in [5]. The DINO training strategy is shown in bottom right of Figure 3. Inspired from these findings, we also employ ViT-Base[13] architecture trained using DINO [5]. In our study, we use the ViT-Base as feature extractor, and add a classification head on top. We only train the added classification head on participating deepfake detection datasets, while freezing the weights of the ViT-Base feature extractor, i.e., we do not train the feature extractor, but only the classification head.

- **Contrastive Language-Image Pre-Training (CLIP)** [31] is a neural network that has been trained on a diverse set of (image, text) pairs in a self-supervised manner. It has the ability to infer the most suitable text excerpt for a given image using natural language, without explicit supervision for this task. It exhibits zero-shot capabilities similar to the ones exhibited by GPT-2/GPT-3. In CLIP's original research paper, authors show that it achieves performance scores equivalent to the original ResNet50 CNN model[17] when evaluated on ImageNet[10] in a "zero-shot" fashion, i.e., even though CLIP does not use any of the 1.28 million labelled examples from the original dataset it achieves comparable performance as a ResNet50 model trained on ImageNet in a supervised manner. CLIP is illustrated in the bottom left corner in Figure 3. For more details on CLIP, we refer readers to [31]. We employ a ViT-Base model trained using CLIP as a feature extractor for our study. Similar to DINO, we add a classification head on top of ViT-Base trained using CLIP. For our analysis, we only train the classification head, and keep the CLIP ViT-Base features frozen i.e., we do not update its weights during training.

3.4.2 Video Models. We studied two different video classification models, (1) ResNet-3D (a CNN based video classifier), and (2) TimeSformer (a transformer based video classification model). We study the performance of both these models on intra, as well as inter dataset performance on four well-known deepfake detection benchmarks. We choose to study video based models in addition to the image based detection models in order to find out whether temporal information help in the detection task. Below we briefly describe these models.

- **ResNet-3D** [16] is based on the same principles as the original ResNet architecture [17], but they are specifically designed to work with 3D data, such as videos and volumetric medical images. These models use 3D convolutions, instead of 2D layers, for feature extraction. In addition to that, ResNet-3D models generally use a large number of layers, which allows them to learn complex and abstract features in the data. ResNet-3D models have been utilised for a variety of computer vision tasks, including video classification, action recognition, and medical image segmentation. They have been shown to achieve SoTA performance on a number of different benchmarks, however, it is also worth noting that the ResNet-3D models are computationally costly, and need a large amount of data to train. For reference, we illustrate both 2D and 3D convolutions in Figure 3, on the left side of third row. We choose to employ ResNet-3D model for our study because, (1) it is widely studied in regards of video recognition, (2) pre-trained models are available, and (3) the available compute power which is not suitable for training bigger video recognition models using more video frames for training. We chose ResNet-3D model pre-trained on 8 frames per video to experiment in this study. In contrast, available video classification models are typically trained on 16/32 frames per video and tend to perform better than models trained using 8 frames per video.

- **TimeSformer** [4] is a video recognition model based on the transformer architecture. TimeSformer utilises self-attention over space and time, instead of traditional convolutional layers, or the spatial attention as employed by ViT for image recognition. The TimeSformer model modifies the transformer architecture, generally used for image recognition, by directly learning the spatio-temporal features from a sequence of frame-level patches. This is accomplished by extending the self-attention mechanism from the image space to the 3D space-time volume. Similar to the Vision Transformer (ViT) model, the TimeSformer employs linear mapping and positional information to interpret ordering of the resulting sequence of features. In TimeSformer paper [4], authors experimented with different self-attention techniques. Out of different techniques, the "divided attention" technique which calculates temporal and spatial attention separately within each block, was found to perform better than other self-attention calculation techniques, and thus we choose to analyse the same architecture in this study. Divided space-time attention is illustrated in Figure 3, in the middle of second row. We opt to evaluate TimeSformer on the task of deepfake detection, and compare it with convolutional video classification network, ResNet-3D. We also chose 8 frame per video version of the TimeSformer model, same as the ResNet-3D model we described above.

3.5 Evaluation Metrics

In order to analyse the performance of our models in a comprehensive way, we employ multiple widely used classification metrics, e.g., (1) LogLoss, (2) AUC, and (3) Accuracy. Below we briefly describe the chosen evaluation metrics.

3.5.1 LogLoss. , also known as logarithmic loss or cross-entropy loss, is used to measure the classification performance of machine/deep learning models. LogLoss calculates the dissimilarity between the predicted probability score with the true label (0, 1 in case of binary classification). The LogLoss score is computed as the negative logarithm of the likelihood of the true labels given a set of predicted probabilities. The range of the LogLoss function is from 0 to infinity, with 0 representing the ideal outcome and higher values representing worse outcomes.

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (1)$$

Where L is the LogLoss, N is the total number of samples in the dataset, y_i is the true label of the i -th sample, p_i is the predicted probability for the i -th sample.

It is worth noting that Logloss is a widely used evaluation metric in machine learning competitions such as Kaggle competitions, as it gives a general idea of how good the predictions of the model are. We use LogLoss as one of the evaluation metrics in this study as other previously proposed deepfake detection research studies often use it as their evaluation metric, and thus would allow us to compare our results with them.

3.5.2 Area Under the Curve (AUC). is also a widely known metric used to evaluate classification models. AUC basically refers to calculating the entire two-dimensional area under the Receiver Operating Curve (ROC). AUC gives hints about how well a model has made a certain prediction. Quite understandably, the higher the area falling under the ROC, i.e., AUC, the better the performance of the model at discriminating between "real" and "fake" samples in our case. Most of the recently proposed deepfake detection studies employ AUC as the evaluation metric to study the performance of their models.

Note that the ROC curve is created by varying the threshold used to make predictions from 0 to 1, so the AUC provides a summary of the model's performance across all possible thresholds.

3.5.3 Accuracy. is also a widely used metric in the classification domain. Accuracy score is basically the measure of correct predictions made by a model in relation to all the predictions made by the model. Accuracy does not indicate how well a model has made a certain classification, as was the case with LogLoss and AUC. Accuracy score can be obtained by dividing the number of correct predictions by total predictions.

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (2)$$

Where TP is the number of true positives, TN refers to the number of true negatives, FP refers to the number of false positives, and FN refers to the number of false negatives.

It is worth noting that accuracy is the proportion of correctly classified samples out of the total number of samples. It is a common evaluation metric used in binary classification tasks, however, it can be misleading in cases where the classes (real, fake) are imbalanced, or if the cost associated with the false positives and false negatives is different. In such cases, other evaluation metrics like F1 score, precision, recall, or AUC may provide a more accurate evaluation of the classification model's performance [34]. In our study however, since we have balanced number of samples both for **real** and **fake** classes, we can use accuracy as one of the evaluation metric.

3.6 Implementation Details

We use PyTorch library to train and test our the models. To train our models we choose a batch size of 16 for image models, and 4 for video models. We choose a constant learning rate of 3×10^{-3} for both image and video models. We use CrossEntropyLoss as the loss function and SGD (Stochastic Gradient Descent) as the optimiser to train our models. We train our models for 5 epochs, and choose the model with lowest validation loss for further testing and evaluation. For evaluation we rely on Scikit-Learn library, i.e., to report on LogLoss, AUC, Accuracy scores.

We rely heavily on Ross Wightman's PyTorch Image Models repository [42] for model code implementations and pre-trained weights. We adapt some code from [5] to train linear classification

head on top of self-supervised feature extractors i.e., DINO and CLIP. For image augmentations we rely on *imgaug* [20] library

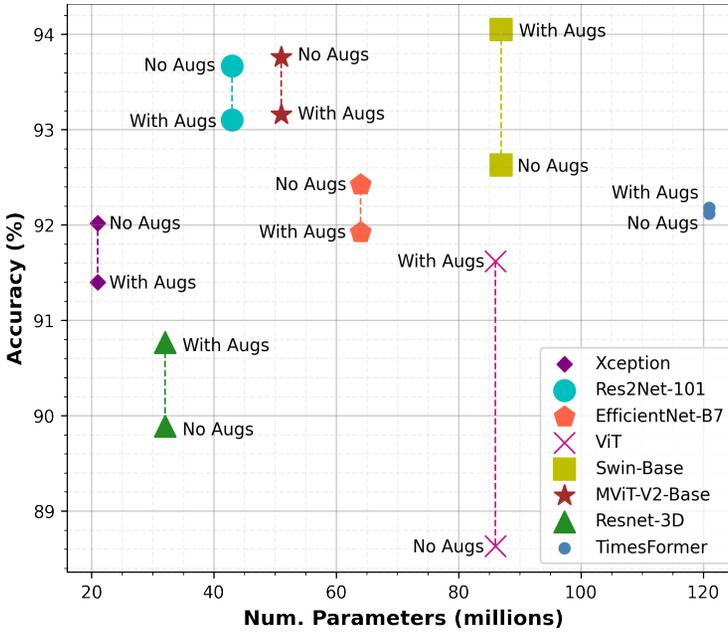


Fig. 4. Performance (accuracy) comparison of participating models on all datasets. The reported scores were achieved in an intra-dataset evaluation. Results in this figure are obtained by, (1) evaluating each model separately on each dataset, and (2) averaging the achieved scores i.e., add the 4 accuracy scores and divide by 4.

4 RESULTS

We conducted extensive experiments and evaluated six image deepfake detection models, as well as two video deepfake detection models on four different benchmarks (the details were discussed in Section 3). In addition to this, we also evaluate two vision transformer (ViT-Base) models pre-trained using self-supervised techniques mentioned in Section 3. We evaluated all models in both intra-dataset as well as inter-dataset settings. In the following sections, we report the performance of all the participating models both in an intra-dataset (trained and evaluated on same dataset), as well as inter-dataset (trained on one dataset and evaluated on the remaining datasets excluding the training dataset) settings.

In this section we refer to models as supervised models, and self-supervised models. Supervised models refer to eight models including six image models, and two video models. Self-supervised models refer to DINO, CLIP and a supervised ViT-Base (which is used as a feature extractor to compare with DINO, and CLIP based ViT-Base). Supervised models are trained end to end i.e., weights of feature extractor as well as the classification head are updated during training. In case of self-supervised models including DINO, CLIP, and a supervised ViT-Base, the weights of feature extractors are kept frozen during training and only the classification head is trained.

DINO, and CLIP are also ViT-Base models, however, the only difference is that both DINO and CLIP are pre-trained using self-supervised training strategies. The supervised ViT-Base is pre-trained using supervised training strategy. Through training a classification head on top of

these three models we aim to find out whether self-supervised features provide better feature representations as compared to supervised features.

Table 2. Intra-dataset comparison of image models. The table below presents scores achieved by image models when trained and evaluated on FakeAVCeleb [22] dataset. Best results are highlighted in yellow.

FakeAVCeleb						
Model	With Augs			No Augs		
	LogLoss	AUC	ACC	LogLoss	AUC	ACC
Xception	0.0047	100.00%	99.93%	0.0040	100.00%	99.85%
Res2Net-101	0.0008	100.00%	99.98%	0.0037	100.00%	99.93%
EfficientNet-B7	0.0132	100.00%	99.63%	0.0047	100.00%	99.83%
ViT	0.2073	99.29%	94.60%	0.3768	98.78%	92.43%
Swin	0.0033	100.00%	99.88%	0.0058	100.00%	99.83%
MViT	0.0008	100.00%	100.00%	0.0023	100.00%	99.95%
ResNet-3D	0.0041	100.00%	100.00%	0.0066	100.00%	100.00%
TimeSformer	0.0796	99.96%	97.50%	0.1238	99.94%	97.00%

4.1 FakeAVCeleb

FakeAVCeleb [22] is a newly proposed deepfake detection dataset containing four different categories of videos i.e., (1) FakeVideo/FakeAudio, (2) RealVideo/RealAudio, (3) FakeVideo/RealAudio, and (4) RealVideo/FakeAudio. Since we focus only on visual deepfakes in this study, and thus do not use the audios (real and fake) for training and evaluating our models. In fact, out of four subsets of the FakeAVCeleb dataset, we only use two for our experiments i.e., (1) FakeVideo/FakeAudio, (2) RealVideo/RealAudio.

Table 2 shows that all models perform pretty well in distinguishing between fake and real faces. We can see that all of the participating models achieved almost 99% AUC, and very low LogLoss score when tested in an intra-dataset configuration. The numbers in 2 suggest that FakeAVCeleb dataset is relatively easy and thus the models can accurately distinguish between real and fake samples.

In table 10 we report results achieved by all the models when trained on FakeAVCeleb, and evaluated on the remaining three datasets. When we look at the numbers in Table 10, it is apparent that almost all of the models perform poorly on all the other datasets. We can see that in terms of accuracy scores, the models are making random guesses. LogLoss, and AUC scores are also not remarkably good in inter-dataset evaluation.

In case of self-supervised models, the results are not as good as they are for the supervised models. That is because the self-supervised models are not trained in an end to end manner as we mentioned earlier. However, on FakeAVCeleb dataset, even though only the classification head is trained, still DINO, and ViT-Base (the supervised feature extractor) achieve good performance. However, DINO performs significantly better than the other two models as shown in Table 8. CLIP does not achieve good performance, and this might be because CLIP was initially pre-trained using both images and their associated text captions, however, in this study we use CLIP's image encoder only without any text captions. This might be a reason of bad performance by CLIP. We aim at

investigating this issue along with the inter-dataset analysis of self-supervised models in our future research.

From the results we can infer that FakeAVCeleb dataset is not challenging enough for the models to learn, and is fairly easy to distinguish between fake and real samples for both supervised and self-supervised models. In addition to that, this dataset does not enhance the models' ability to learn distinguishing features between real and fake faces, or in other words, it lacks at integrating the generalisation capability into the models, as is apparent from Table 9 and 10.

Table 3. Intra-dataset comparison of image models. The table below presents scores achieved by image models when trained and evaluated on CelebDF-V2 [26] dataset.

CelebDF						
Model	With Augs			No Augs		
	LogLoss	AUC	ACC	LogLoss	AUC	ACC
Xception	0.0712	99.73%	97.00%	0.0367	99.95%	98.55%
Res2Net-101	0.0237	100.00%	98.95%	0.0185	99.99%	99.45%
EfficientNet-B7	0.0433	99.95%	98.40%	0.0340	99.98%	98.75%
ViT	0.0336	99.96%	98.60%	0.0350	99.95%	98.60%
Swin	0.0340	99.94%	98.80%	0.0202	99.97%	99.40%
MViT	0.0075	100.00%	99.70%	0.0096	100.00%	99.70%
ResNet-3D	0.0748	99.68%	97.00%	0.1525	98.68%	95.00%
TimeSformer	0.0309	100.00%	98.00%	0.0220	99.96%	99.00%

4.2 CelebDF-V2

Table 3 presents the performance of supervised models when trained and evaluated on CelebDF-V2 [26] dataset. Same as it was the case with FakeAVCeleb dataset, almost all of the participating models achieve excellent scores i.e., more than 97% accuracy, and more than 99% AUC score, while having a very small LogLoss. We can thus infer that the models quite comfortably learnt to discriminate between real/fake samples of the CelebDF-V2 dataset, similar to FakeAVCeleb dataset.

In order to find out how helpful the dataset is in making models learn robust features making them better at generalisation, we also conduct extensive inter-dataset evaluation of all the participating models trained on CelebDF-V2. We report the results achieved by the models in Table 11. However, similar to the results achieved by the models trained on FakeAVCeleb dataset and evaluated on the remaining datasets, the models trained on CelebDF-V2 and evaluated in an inter-dataset setting also seem to perform poorly. This might be a result of CelebDF-V2 not being a very challenging dataset for the models to discriminate, and they can almost classify every real/fake sample in a perfect manner. However, this also makes the models less powerful against unseen data, as can be seen by the performance scores reported in table 11.

The fact that CelebDF-V2 is also not a challenging dataset is also supported by the results achieved by the self-supervised models as presented in Table 8. It is apparent from the numbers that only training a classification head on top of frozen feature extractors, models still achieve good results. However, in this case as well, CLIP does not achieve good performance as compared to DINO and supervised ViT.

Table 4. Intra-dataset comparison of image models. The table below presents scores achieved by image models when trained and evaluated on FaceForensics++ [32] dataset.

FaceForensics++						
Model	With Augs			No Augs		
	LogLoss	AUC	ACC	LogLoss	AUC	ACC
Xception	0.2342	96.96%	91.05%	0.2957	95.85%	89.03%
Res2Net-101	0.2165	97.87%	93.48%	0.3213	97.30%	91.85%
EfficientNet-B7	0.3111	96.92%	90.33%	0.3737	94.02%	86.95%
ViT	0.2445	97.27%	92.18%	0.3571	94.04%	85.15%
Swin	0.1573	98.58%	94.90%	0.2191	97.60%	92.18%
MViT	0.1828	98.34%	94.10%	0.1918	97.63%	93.00%
ResNet-3D	0.3224	96.42%	90.36%	0.3085	96.19%	91.07%
TimeSformer	0.2807	97.10%	90.00%	0.2451	96.76%	90.71%

4.3 FaceForensics++

Table 4 reports the performance metrics of all the supervised models when trained and evaluated on the FaceForensics++ [32] dataset. We can see that the results are not as good as they were in the case of previous two datasets, FakeAVCeleb and CelebDF-V2. Almost none of the models achieved more than 95% accuracy, and the LogLoss scores are also not as great as they were for the previous two datasets. We can thus infer by looking at the metrics that it is a relatively challenging dataset for the models to distinguish between real/fake samples. Self-supervised models also are not able to achieve excellent results on FaceForensics++ dataset as is apparent from the numbers in Table 8, affirming that it is indeed challenging to properly distinguish between fake and real faces coming from FaceForensics++ dataset. What we would now like to see now is, whether a more challenging dataset means better generalisation capability?

We thus evaluate all of the supervised models trained on FaceForensics++ dataset in an inter-dataset evaluation setting and report the results in Table 12. In this case, it can be clearly seen that the models perform in a somewhat acceptable manner even on unseen data from other datasets. For example, MViT trained on FaceForensics++ and evaluated on FakeAVCeleb dataset achieves more than 80% accuracy, and more than 90% AUC score. This also supports our claim that FakeAVCeleb as well as CelebDF-V2 datasets are not very challenging, and that the models can easily learn to distinguish real/fake videos coming from these datasets.

Furthermore, not only on the FakeAVCeleb dataset, we can also see somewhat better performance from all of the participating supervised models on the other two datasets, i.e., CelebDF-V2 and DFDC. However, it must be noted that how all the models suffer when tested using unseen data, i.e., lack of generalisation, which is a big problem the current, more sophisticated deepfake detection systems suffer from. The results in Table 12 somehow support the statement that more challenging datasets mean better generalisation capability. But we have to further re-enforce this statement after analysing the metrics of models trained using DFDC [11] dataset in the section below.

4.4 DFDC

DFDC is one of the biggest and most challenging deepfake detection benchmarks. This is apparent by the results we present in Table 5. Only one of the supervised models (i.e., Res2Net-101) managed

Table 5. Intra-dataset comparison of image models. The table below presents scores achieved by image models when trained and evaluated on DFDC [11] dataset.

DFDC						
Model	With Augs			No Augs		
	LogLoss	AUC	ACC	LogLoss	AUC	ACC
Xception	0.5613	88.75%	77.63%	0.5120	91.68%	80.65%
Res2Net-101	0.5570	90.64%	79.98%	0.5691	91.78%	83.45%
EfficientNet-B7	0.5542	89.97%	79.30%	0.4263	93.30%	84.15%
ViT	0.4696	91.89%	81.08%	0.5709	89.44%	78.35%
Swin	0.5602	90.89%	82.60%	0.6650	87.77%	79.05%
MViT	0.6079	88.41%	78.90%	0.5491	90.65%	82.40%
ResNet-3D	0.5865	85.64%	75.75%	0.6739	84.69%	73.50%
TimeSformer	0.4870	91.18%	83.25%	0.6176	92.30%	81.75%

to achieve more than 84% accuracy score, 93% AUC score on the DFDC dataset. Self-supervised models also achieve relatively low scores when trained and evaluated on DFDC, as apparent from Table 8. This establishes that DFDC is the most challenging dataset out of all the four datasets in this study.

In Table 13 we present inter-dataset evaluation scores achieved by the supervised models trained on DFDC dataset. It is apparent from the results that the models trained using DFDC dataset still achieve acceptable performance on unseen data, as compared to the scores achieved by the models which were trained on FakeAVCeleb, and CelebDF-V2. Also, by looking at the results we can *somewhat* affirm the statement that models trained using more challenging datasets seem to achieve better results. We say *somewhat*, because in the scope of this study, even though DFDC is more challenging to learn for the models, a better generalisation is offered by FaceForensics++ which is relatively less challenging to learn for the models.

4.5 Discussion

In Figure 4 we present a comparison of all the participating models on the basis of achieved accuracies while evaluated in an intra-dataset setting (i.e., models are trained and tested on same dataset). From the figure, it is apparent that there is not a lot of performance difference between the participating models. In most cases, the models achieves around 92% to 94% accuracies when tested in an intra-dataset evaluation setting. The figure also shows that the image augmentations do not prove to be significantly helpful in all cases, for example, we can see that for XceptionNet, Res2Net-101, MViT-V2-Base, and EfficientNet-B7, the models trained without image augmentations achieved better scores when compared to models trained using image augmentations. The difference between the accuracies achieved by the models when trained with and without image augmentations is not big except for the ViT. The ViT trained with image augmentations achieved 91.62% while the ViT trained without image augmentations achieved 88.63% accuracy. However, it is apparent from the Figure 4 that all of the transformer models perform better when trained using augmentations. Besides this, video models also perform when trained using augmentations. Another reason to no disregard image augmentations is that the best performing model Swin-Base achieved top most accuracy while being trained using image augmentations.

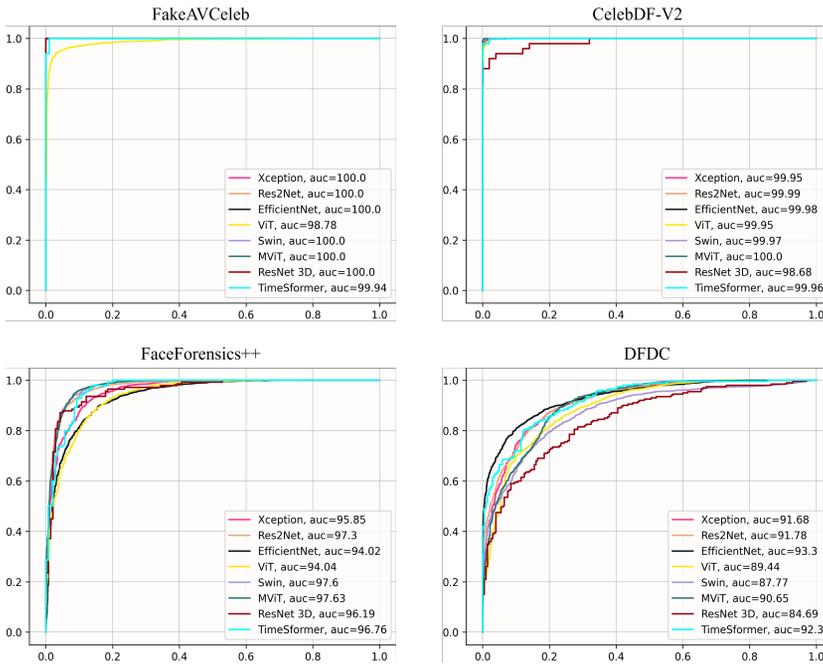


Fig. 5. ROC curves of each of the model when evaluated on each of the 4 different participating datasets in an intra-dataset evaluation setting.

Also, it can be noted that the transformer models (Swin-Base, and MVit-V2-Base) outperform the CNN based models. We can also see that the Res2Net-101 model also achieves excellent performance scores in intra-dataset evaluation setting, while having almost half the amount of parameters (43 million parameters) as compared to the best performing Swin-Base (87 million parameters) model. From Figure 4, and Table 6 we make one useful observation, i.e., models with multi-scale feature processing capabilities (Res2Net, MVit-V2, and Swin Transformer) are the best performing models out of all.

Moving towards inter-dataset analysis, we illustrate results achieved by supervised models when evaluated in an inter-dataset setting in Figure 7. From the figure it is apparent that in inter-dataset evaluation, models perform in a significantly poor manner, as compared to intra-dataset evaluation. However, this is understandable i.e., detection models loose performance when tested on data coming from a different distribution. However, their is a useful finding in Figure 7, i.e., for all datasets, the best performing models are transformers.

In case of self-supervised models, by analysing the results in Table 7 and 8, we can infer that self-supervised features (DINO) indeed provide better feature representations as compared to supervised feature representations. DINO seems to outperform both supervised ViT as well as CLIP as is also apparent from the ROC curvers illustrated in Figure 8.

We also visualise the TSNE plots for all the participating models in Figure 6, to get an idea about how the models separate real faces from those of the fake ones. Also, it gives us an idea about how the models group together faces coming from same datasets near to each other as compared to the faces coming from a different dataset. The TSNE plots also help us visualise which datasets are more challenging than the others. For example, if we look at the TSNE plots in Figure 6, we can see that the models tend to separate the easier datasets (FakeAVCeleb, and CelebDF-V2) in a better way,



Fig. 6. TSNE visualisations of the participating detection models. We chose the best performing models on all datasets (with/without image augmentations).

as compared to how they separate the more challenging datasets (FaceForensics++, and DFDC). We can see that the Res2Net-101 model does this separation in a better way as compared to the other models i.e., even better than the best performing Swin-Base model (we don't know how/why). The plots also support our findings that FaceForensics++ and DFDC are indeed challenging datasets, and that the models are not as good at separating the fake, and real faces coming from these datasets as they can separate the fake, and real faces coming from FakeAVCeleb and CelebDF-V2 datasets.

Another finding which is quite apparent is that the image models do this separation job relatively better than those of the video models. This is understandable, as we pointed out above that the video models need relatively larger amounts of training data (in our case we train both the image/video models on the same amount of data). Also, the video models we use are not the newest, most powerful models out there, however, we choose these because of the lack of compute resources, and for the sack of experimentation.

We also present the ROC curves of the participating models evaluated in an intra-dataset setting in Figure 5. The AUC scores achieved by the models also show that FakeAVCeleb and CelebDF-V2 datasets are easier to learn for the models, as compared to FaceForensics++ and DFDC datasets. This also suggests that while training the models for deepfake detection, it will give better generalisation capability to the models when they are trained on challenging datasets, rather than the easier ones.

$$\text{Score} = \frac{s_1 + s_2 + s_3 + s_4}{4} \quad (3)$$

Table 6. This table compares the performance of all the participating (supervised) models. We present scores after averaging the scores (LogLoss, AUC, Accuracy) achieved by each model when evaluated in an intra-dataset setting as given in Equation 3.

Performance Comparison of All Models on All Datasets						
Model	With Augs			No Augs		
	LogLoss	AUC	ACC	LogLoss	AUC	ACC
Xception	0.2179	96.36%	91.40%	0.2121	96.87%	92.02%
Res2Net-101	0.1395	97.12%	93.10%	0.2282	97.27%	93.67%
EfficientNet-B7	0.2305	96.71%	91.92%	0.2097	96.83%	92.42%
ViT	0.2388	97.10%	91.62%	0.3350	95.55%	88.63%
Swin	0.1494	97.35%	94.05%	0.2275	96.34%	92.63%
MViT	0.1998	96.68%	93.16%	0.1882	97.07%	93.76%
Resnet-3D	0.2470	95.44%	90.77%	0.1620	94.89%	89.89%
TimeSformer	0.2196	97.06%	92.18%	0.2521	97.24%	92.12%

where s_1 refers to score (LogLoss, AUC, ACC) achieved by a model when trained and evaluated on dataset1, s_2 refers to score (LogLoss, AUC, ACC) achieved by a model when trained and evaluated on dataset2 and so on. The scores reported in Tables 6 and 7 for each model are calculated using this equation.

Table 7. This table compares the performance of the self-supervised models. We present scores after averaging the scores (LogLoss, AUC, Accuracy) achieved by each model when evaluated in an intra-dataset setting, as given in Equation 3. In this table, **Supervised** refer to ViT-Base model pre-trained using supervised training scheme. **DINO** refers to ViT-Base model pre-trained using self-supervised scheme proposed in [5], and **CLIP** refers to ViT-Base pre-trained using self-supervised scheme proposed in [31]. All of these ViT-Base models are used as feature extractors, where we only train a classification head on top of each of the feature extractor, and freeze the weights of feature extractors.

Performance Comparison of All Models on All Datasets						
Model	With Augs			No Augs		
	LogLoss	AUC	ACC	LogLoss	AUC	ACC
Supervised	0.4516	87.00%	78.54%	0.4191	88.84%	81.59%
Dino	0.9924	91.43%	84.80%	0.7932	92.31%	85.54%
CLIP	1.0244	66.17%	62.26%	1.0513	69.40%	63.98%

5 CONCLUSIONS

In conclusion, this paper investigates the performance of various image and video classification architectures (supervised, self-supervised) on the task of deepfake detection when trained and evaluated on four different datasets. We aimed at identifying which models perform better than other participating models, which model generalises well on unseen data as compared to the other models. Through experimentation and analysis of the achieved results we conclude that models possessing the capability of processing multi-scale features (Res2Net-101, MViT-V2, and SWIN Transformer) achieve better overall performance in intra-dataset comparison. For inter-dataset

Table 8. This table compares the performance of all the participating (self-supervised) models when evaluated in an intra-dataset setting. The statistics of this table are illustrated in Figure 8.

Performance Comparison of All Models on Individual Datasets							Dataset
Model	With Augs			No Augs			
	LogLoss	AUC	ACC	LogLoss	AUC	ACC	
Supervised	0.4105	90.19%	82.50%	0.3727	91.77%	85.50%	FakeAVCeleb
Dino	0.1444	99.00%	95.33%	0.0801	99.64%	97.25%	
CLIP	1.2851	63.93%	60.65%	1.3111	65.20%	61.13%	
Supervised	0.2941	95.52%	88.05%	0.2237	97.18%	91.80%	CelebDF-V2
Dino	0.3655	97.31%	90.90%	0.3930	97.10%	88.90%	
CLIP	0.7811	77.00%	69.50%	0.6860	82.10%	74.10%	
Supervised	0.5182	83.11%	74.95%	0.4971	85.47%	77.43%	FaceForensics++
Dino	1.1758	88.67%	80.60%	1.1186	89.48%	81.85%	
CLIP	0.9634	62.46%	59.05%	1.1070	62.61%	59.73%	
Supervised	0.5836	79.19%	68.65%	0.5829	80.93%	72.63%	DFDC
Dino	2.2839	80.72%	72.38%	1.5812	83.03%	74.15%	
CLIP	1.0681	61.30%	59.83%	1.1011	67.69%	60.95%	

comparison, or in other words, the generalisation capability comparison, we infer from the results that transformer models perform better than that of the CNN models. It is also apparent from the results obtained through both inter-dataset as well as intra-dataset comparisons, the image augmentations do not always help achieve better performance scores.

Through intra-dataset comparisons we establish that the DFDC dataset is the most challenging dataset to learn for the models, whereas FaceForensics++ dataset is ranked as second most challenging dataset. However, through inter-dataset evaluation, we establish that the FaceForensics++ dataset offers the best generalisation capabilities to the models, as compared to other datasets. DFDC ranks second in providing the generalisation capabilities. The remaining two datasets FakeAVCeleb, and CelebDF-V2 appear to be fairly easy for the models to learn and achieve excellent performance in intra-dataset comparison. However, they do not provide the models with any generalisation capability, i.e., models trained using these datasets perform poorly when evaluated on other datasets.

In addition to analysing supervised image/video recognition models, we also explore the performance of self-supervised models for deepfake detection in an intra-dataset setting. Through our experiments we find the ViT-Base model which was pre-trained using DINO [5] to achieve better performance as compared to the supervised ViT-Base, and self-supervised CLIP ViT-Base. We also find in these experiments these models achieve better performance scores when trained without using image augmentations.

All in all, we present a detailed analysis of the performance achieved by several different deepfake detection architectures on four different deepfake detection benchmarks. We carry out extensive experiments and provide detailed results along with useful visualisations to help understand the overall contributions of this paper to the reader. We regard this study as an entry point for researchers exploring the research field of deepfake detection who are trying to make sense of

different architectures and datasets to develop their own solutions. We are confident that this study provides useful insights into the problem of deepfake detection.

In future work, we aim at analysing even more diverse set of architectures, and newer datasets. In addition to that we plan to focus more towards self-supervised training strategies to train models, as well as try to incorporate knowledge distillation, domain adaptation strategies to help make models better at classifying unseen samples correctly.

ACKNOWLEDGMENTS

This research was supported by industry partners and the Research Council of Norway with funding to MediaFutures: Research Centre for Responsible Media Technology and Innovation, through the Centres for Research-based Innovation scheme, project number 309339.

REFERENCES

- [1] Accessed: 2022-12-02. Deepfakes github. <https://github.com/deepfakes/faceswap>.
- [2] Accessed: 2022-12-02. FaceSwap github. <https://github.com/MarekKowalski/FaceSwap/>.
- [3] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. 2018. MesoNet: a Compact Facial Video Forgery Detection Network. In *Proceedings of IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE. <https://arxiv.org/abs/1809.00888>
- [4] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is Space-Time Attention All You Need for Video Understanding?. In *International Conference on Machine Learning*.
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Herv'e J'egou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging Properties in Self-Supervised Vision Transformers. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)* (2021), 9630–9640.
- [6] Chun-Fu Chen, Quanfu Fan, Neil Rohit Mallinar, Tom Sercu, and Rogério Schmidt Feris. 2018. Big-Little Net: An Efficient Multi-Scale Feature Representation for Visual and Speech Recognition. *ArXiv abs/1807.03848* (2018).
- [7] François Chollet. 2017. Xception: Deep Learning with Depthwise Separable Convolutions. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017), 1800–1807.
- [8] Umur Aybars Ciftci, Ilke Demir, and Lijun Yin. 2020. FakeCatcher: Detection of Synthetic Portrait Videos using Biological Signals. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*. IEEE. <https://arxiv.org/abs/1901.02212>
- [9] Davide Alessandro Cocomini, Nicola Messina, Claudio Gennaro, and F. Falchi. 2021. Combining EfficientNet and Vision Transformers for Video Deepfake Detection. In *International Conference on Image Analysis and Processing*.
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition* (2009), 248–255.
- [11] Brian Dolhansky, Joanna Bitton, Ben Pfau, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. 2020. The DeepFake Detection Challenge (DFDC) Dataset. *arXiv: Computer Vision and Pattern Recognition* (2020).
- [12] Luke Dormehl. 2021. Inside the rapidly escalating war between deepfakes and deepfake detectors. <https://www.digitaltrends.com/cool-tech/deepfake-detection-war/>.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *ArXiv abs/2010.11929* (2021).
- [14] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. 2021. Multiscale Vision Transformers. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)* (2021), 6804–6815.
- [15] Shanghua Gao, Ming-Ming Cheng, Kai Zhao, Xinyu Zhang, Ming-Hsuan Yang, and Philip H. S. Torr. 2021. Res2Net: A New Multi-Scale Backbone Architecture. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43 (2021), 652–662.
- [16] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. 2017. Learning Spatio-Temporal Features with 3D Residual Networks for Action Recognition. *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)* (2017), 3154–3160.
- [17] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 770–778.
- [18] Dong Huang and Fernando De la Torre. 2012. Facial Action Transfer with Personalized Bilinear Regression. In *European Conference on Computer Vision*.
- [19] Ye Jia, Yu Zhang, Ron J. Weiss, Quan Wang, Jonathan Shen, Fei Ren, Z. Chen, Patrick Nguyen, Ruoming Pang, Ignacio Lopez-Moreno, and Yonghui Wu. 2018. Transfer Learning from Speaker Verification to Multispeaker Text-To-Speech

- Synthesis. In *Neural Information Processing Systems*.
- [20] Alexander B. Jung. 2018. imgaug. <https://github.com/aleju/imgaug>. [Online; accessed 30-Oct-2018].
- [21] Tero Karras, Samuli Laine, and Timo Aila. 2018. A Style-Based Generator Architecture for Generative Adversarial Networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018), 4396–4405.
- [22] Hasam Khalid, Shahroz Tariq, and Simon S. Woo. 2021. FakeAVCeleb: A Novel Audio-Video Multimodal Deepfake Dataset. *ArXiv abs/2108.05080* (2021).
- [23] Sohail Ahmed Khan and Hang Dai. 2021. Video Transformer for Deepfake Detection with Incremental Learning. In *Proceedings of the 29th ACM International Conference on Multimedia*.
- [24] Iryna Korschunova, Wenzhe Shi, Joni Dambre, and Lucas Theis. 2016. Fast Face-Swap Using Convolutional Neural Networks. *2017 IEEE International Conference on Computer Vision (ICCV)* (2016), 3697–3705.
- [25] Yuezun Li, Ming-Ching Chang, and Siwei Lyu. 2018. In icu oculi: Exposing ai created fake videos by detecting eye blinking. In *2018 IEEE international workshop on information forensics and security (WIFS)*. IEEE, 1–7.
- [26] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. 2020. Celeb-DF: A Large-Scale Challenging Dataset for DeepFake Forensics. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 3204–3213.
- [27] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)* (2021), 9992–10002.
- [28] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. 2019. Capsule-Forensics: Using Capsule Networks to Detect Forged Images and Videos. In *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. <https://ieeexplore.ieee.org/document/8682602>
- [29] Yuval Nirkin, Yosi Keller, and Tal Hassner. 2019. FSGAN: Subject Agnostic Face Swapping and Reenactment. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 7183–7192.
- [30] Prajwal K R, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C. V. Jawahar. 2020. A Lip Sync Expert Is All You Need for Speech to Lip Generation In the Wild. *Proceedings of the 28th ACM International Conference on Multimedia* (2020).
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *International Conference on Machine Learning*.
- [32] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. 2019. FaceForensics++: Learning to Detect Manipulated Facial Images. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 1–11.
- [33] Ekraam Sabir, Jiaxin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. 2019. Recurrent Convolutional Strategies for Face Manipulation Detection in Videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. IEEE, 80–87. https://openaccess.thecvf.com/content_CVPRW_2019/html/Media_Forensics/Sabir_Recurrent_Convolutional_Strategies_for_Face_Manipulation_Detection_in_Videos_CVPRW_2019_paper.html
- [34] Soner Yıldırım. 2021. When You Should Not Use Accuracy to Evaluate Your Machine Learning Model. <https://tinyurl.com/3mjwxx9w>. [Online; accessed 20-Nov-2022].
- [35] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Kumar Gupta. 2017. Revisiting Unreasonable Effectiveness of Data in Deep Learning Era. *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), 843–852.
- [36] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, D. Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2014. Going deeper with convolutions. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014), 1–9.
- [37] Mingxing Tan and Quoc V. Le. 2019. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *ArXiv abs/1905.11946* (2019).
- [38] Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred Neural Rendering: Image Synthesis using Neural Textures. *arXiv: Computer Vision and Pattern Recognition* (2019).
- [39] Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2Face: Real-Time Face Capture and Reenactment of RGB Videos. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 2387–2395.
- [40] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. *ArXiv abs/1706.03762* (2017).
- [41] Junke Wang, Zuxuan Wu, Jingjing Chen, and Yu-Gang Jiang. 2021. M2TR: Multi-modal Multi-scale Transformers for Deepfake Detection. *Proceedings of the 2022 International Conference on Multimedia Retrieval* (2021).
- [42] Ross Wightman. 2019. PyTorch Image Models. <https://github.com/rwightman/pytorch-image-models>. <https://doi.org/10.5281/zenodo.4414861>

- [43] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2016. Aggregated Residual Transformations for Deep Neural Networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)*, 5987–5995.
- [44] Fisher Yu, Dequan Wang, and Trevor Darrell. 2017. Deep Layer Aggregation. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (2017)*, 2403–2412.
- [45] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor S. Lempitsky. 2019. Few-Shot Adversarial Learning of Realistic Neural Talking Head Models. *2019 IEEE/CVF International Conference on Computer Vision (ICCV) (2019)*, 9458–9467.
- [46] Cairong Zhao, Chutian Wang, Guosheng Hu, Haonan Chen, Chun Liu, and Jinhui Tang. 2023. ISTVT: Interpretable Spatial-Temporal Video Transformer for Deepfake Detection. *IEEE Transactions on Information Forensics and Security* 18 (2023), 1335–1348.
- [47] Xiangyu Zhu, Hao Wang, Hongyan Fei, Zhen Lei, and S. Li. 2021. Face Forgery Detection by 3D Decomposition. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2928–2938.

A INTER-DATASET EVALUATION SCORES

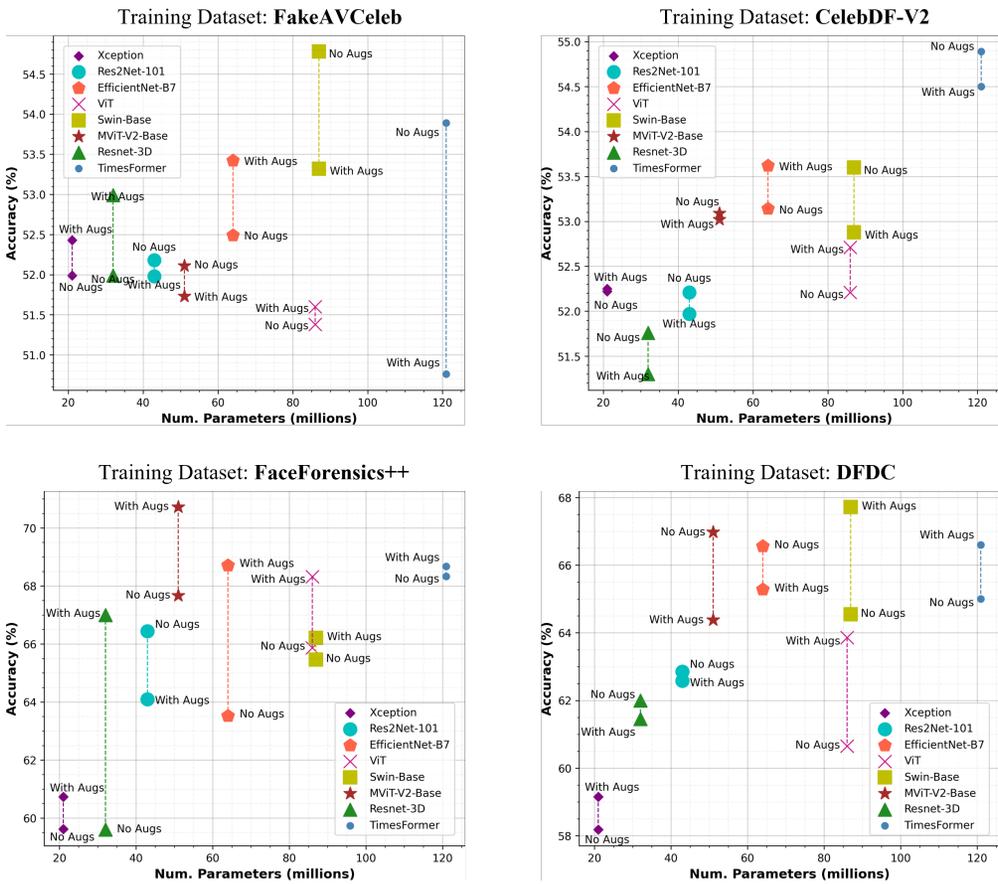


Fig. 7. Performance (accuracy) comparison of participating models evaluated using inter-dataset scheme. Results in this figure are obtained by, (1) evaluating each model trained on one dataset on each of the remaining datasets, and (2) averaging the achieved scores, i.e., add the 3 accuracy scores and divide by 3.

Table 9. This table compares the performance of all the participating (supervised) models evaluated in an inter-dataset setting. We present scores after averaging the scores (LogLoss, AUC, Accuracy) achieved by each model on each of the datasets. Figure 7 illustrate the statistics of this table.

Inter-Dataset Evaluation							Training Dataset
Model	With Augs			No Augs			
	LogLoss	AUC	ACC	LogLoss	AUC	ACC	
Xception	7.4484	57.69%	52.43%	6.8728	55.41%	51.99%	FakeAVCeleb
Res2Net-101	8.5574	58.77%	51.98%	7.0666	60.64%	52.18%	
EfficientNet-B7	8.5664	62.32%	53.42%	10.7718	60.05%	52.49%	
ViT	6.7348	61.01%	51.60%	9.1672	58.45%	51.38%	
Swin	5.1077	62.54%	53.32%	4.3274	64.88%	54.78%	
MViT	4.7564	58.78%	51.73%	4.2891	59.38%	52.11%	
ResNet-3D	4.4308	57.61%	52.99%	3.8206	60.09%	51.99%	
TimeSformer	4.7334	61.55%	50.76%	4.7759	63.95%	53.89%	
Xception	3.9439	65.06%	52.25%	4.8776	66.40%	52.22%	
Res2Net-101	5.4266	65.90%	51.97%	5.6891	66.21%	52.21%	
EfficientNet-B7	5.9514	66.99%	53.62%	8.9668	67.13%	53.14%	
ViT	5.4921	68.52%	52.71%	8.9981	66.36%	52.21%	
Swin	5.6007	70.06%	52.88%	4.8405	70.56%	53.60%	
MViT	4.8723	70.71%	53.02%	4.6419	67.20%	53.09%	
ResNet-3D	6.8365	61.57%	51.30%	5.0504	64.52%	51.76%	
TimeSformer	4.5629	69.04%	54.50%	4.9391	69.43%	54.89%	
Xception	1.0701	69.92%	60.73%	1.1262	67.78%	59.62%	FaceForensics++
Res2Net-101	1.0165	73.46%	64.09%	1.2360	73.61%	66.44%	
EfficientNet-B7	0.8792	79.51%	68.71%	1.0068	69.80%	63.52%	
ViT	0.7899	78.45%	68.32%	0.8301	73.24%	65.87%	
Swin	0.8517	77.94%	66.21%	0.8482	78.03%	65.46%	
MViT	0.8407	79.75%	70.72%	0.7292	75.85%	67.67%	
ResNet-3D	1.0639	74.47%	67.00%	1.3331	66.61%	59.50%	
TimeSformer	1.0665	75.59%	68.67%	0.8492	77.03%	68.33%	
Xception	1.2959	63.62%	59.15%	1.6780	64.18%	58.18%	
Res2Net-101	2.0224	67.80%	62.58%	1.7396	69.50%	62.85%	
EfficientNet-B7	1.0388	71.32%	65.28%	1.2764	72.28%	66.56%	
ViT	1.2198	70.45%	63.86%	1.2498	64.71%	60.65%	
Swin	1.2423	73.49%	67.72%	1.3802	69.49%	64.55%	
MViT	1.2329	72.37%	64.38%	1.2254	72.68%	66.98%	
ResNet-3D	1.1354	66.69%	61.45%	1.1354	66.27%	62.00%	
TimeSformer	1.1421	70.66%	66.60%	1.6584	71.77%	65.00%	

Table 10. Inter-dataset evaluation scores of models trained on FakeAVCeleb [22] dataset and evaluated on the remaining three datasets.

Training Dataset: FakeAVCeleb							Evaluation Dataset
Model	With Augs			No Augs			
	LogLoss	AUC	ACC	LogLoss	AUC	ACC	
Xception	3.1366	50.52%	56.10%	3.9659	44.63%	54.25%	CelebDF-V2
Res2Net-101	3.9007	57.61%	54.60%	3.2127	56.73%	54.50%	
EfficientNet-B7	5.7925	59.31%	54.50%	12.3786	54.30%	51.65%	
ViT	4.7035	56.51%	51.60%	7.0900	52.18%	46.65%	
Swin	3.5148	61.13%	57.70%	2.8360	62.54%	61.45%	
MViT	4.5526	65.37%	54.00%	3.6492	58.18%	55.05%	
ResNet-3D	2.4752	53.88%	51.00%	1.7475	57.36%	49.00%	
TimeSformer	3.9086	51.60%	48.00%	3.2374	58.52%	55.00%	
Xception	10.3539	62.90%	50.38%	8.9711	63.06%	50.33%	FaceForensics++
Res2Net-101	11.6456	59.23%	50.18%	9.7934	58.54%	50.53%	
EfficientNet-B7	10.5412	63.80%	52.45%	10.4825	63.13%	52.05%	
ViT	9.3036	61.70%	51.10%	12.3768	58.44%	50.93%	
Swin	6.1675	62.97%	50.70%	5.5166	64.87%	51.23%	
MViT	4.8833	56.51%	50.65%	4.7116	63.40%	50.85%	
ResNet-3D	6.7343	51.30%	50.71%	5.9796	53.99%	50.71%	
TimeSformer	6.0010	62.61%	51.79%	6.1889	62.60%	51.43%	
Xception	8.8546	59.65%	50.80%	7.6813	58.53%	51.40%	DFDC
Res2Net-101	10.1260	59.48%	51.15%	8.1937	66.66%	51.50%	
EfficientNet-B7	9.3656	63.86%	53.30%	9.4543	62.71%	53.78%	
ViT	6.1972	64.81%	52.10%	8.0348	64.71%	56.58%	
Swin	5.6410	63.51%	51.55%	4.6297	67.25%	51.65%	
MViT	4.8333	54.46%	50.55%	4.5065	56.55%	50.43%	
ResNet-3D	4.0828	67.65%	57.25%	3.7347	68.91%	56.25%	
TimeSformer	4.2907	70.43%	52.50%	4.9015	70.74%	55.25%	

Table 11. Inter-dataset evaluation scores of models trained on CelebDF-V2 [26] dataset and evaluated on the remaining three datasets.

Training Dataset: CelebDF-V2							Evaluation Dataset
Model	With Augs			No Augs			
	LogLoss	AUC	ACC	LogLoss	AUC	ACC	
Xception	4.7313	65.82%	51.68%	5.1136	67.77%	51.78%	FakeAVCeleb
Res2Net-101	5.7429	69.08%	52.30%	4.3332	71.01%	52.83%	
EfficientNet-B7	7.4940	63.86%	52.05%	9.6846	65.93%	51.45%	
ViT	5.0347	69.00%	53.08%	9.6735	61.89%	52.18%	
Swin	6.0084	67.63%	52.20%	4.8922	68.28%	52.70%	
MViT	5.1980	72.43%	52.75%	5.3953	61.24%	51.55%	
ResNet-3D	6.0756	62.79%	50.50%	4.9703	61.58%	50.50%	
TimeSformer	4.8465	69.73%	53.00%	5.8829	68.77%	54.00%	
Xception	4.2473	63.26%	53.53%	5.6357	63.68%	53.58%	FaceForensics++
Res2Net-101	6.3947	64.79%	53.33%	6.9000	63.59%	52.90%	
EfficientNet-B7	6.3164	65.07%	54.80%	8.9065	66.31%	53.98%	
ViT	6.1010	68.14%	53.53%	9.9676	65.50%	53.50%	
Swin	6.0278	70.13%	54.23%	5.5408	68.45%	54.30%	
MViT	5.2175	70.01%	53.15%	4.6160	67.88%	54.08%	
ResNet-3D	7.1877	60.00%	52.14%	5.6544	66.01%	54.29%	
TimeSformer	4.8228	68.84%	57.50%	5.0219	67.55%	56.43%	
Xception	2.8532	66.11%	51.55%	3.8835	67.74%	51.30%	DFDC
Res2Net-101	4.1424	63.83%	50.28%	5.8342	64.01%	50.90%	
EfficientNet-B7	4.0438	72.05%	54.00%	8.3092	69.17%	54.00%	
ViT	5.3405	68.41%	51.53%	7.3534	71.69%	50.95%	
Swin	4.7659	72.42%	52.20%	4.0886	74.95%	53.80%	
MViT	4.2014	69.69%	53.15%	3.9144	72.48%	53.65%	
ResNet-3D	7.2461	61.91%	51.25%	4.5265	65.97%	50.50%	
TimeSformer	4.0195	68.56%	53.00%	3.9124	71.98%	54.25%	

Table 12. Inter-dataset evaluation scores of models trained on FaceForensics++ [32] dataset and evaluated on the remaining three datasets.

Training Dataset: FaceForensics++							Evaluation Dataset
Model	With Augs			No Augs			
	LogLoss	AUC	ACC	LogLoss	AUC	ACC	
Xception	0.8691	79.62%	65.88%	0.7795	76.14%	66.93%	FakeAVCeleb
Res2Net-101	0.7693	83.01%	71.28%	0.6527	85.48%	76.83%	
EfficientNet-B7	0.5782	89.59%	77.05%	0.7375	77.88%	70.08%	
ViT	0.6648	83.05%	70.65%	0.7419	76.40%	69.23%	
Swin	0.5880	87.72%	72.95%	0.6373	89.10%	71.15%	
MViT	0.3654	92.96%	84.65%	0.4047	90.25%	81.90%	
ResNet-3D	0.7903	83.55%	68.00%	1.1338	73.34%	62.50%	
TimeFormer	0.9135	79.33%	75.00%	0.7900	76.65%	70.50%	
Xception	1.0426	65.92%	61.60%	1.2566	62.39%	58.65%	CelebDF-V2
Res2Net-101	1.0751	67.85%	62.40%	1.4218	65.46%	59.80%	
EfficientNet-B7	0.7759	78.46%	69.95%	1.0103	67.24%	61.25%	
ViT	0.5915	82.44%	74.10%	0.8504	75.11%	65.40%	
Swin	0.7136	74.58%	67.05%	0.7879	70.94%	63.75%	
MViT	0.9791	76.66%	65.35%	0.7912	68.69%	62.70%	
ResNet-3D	1.1992	66.12%	65.00%	1.5866	59.44%	55.00%	
TimeFormer	1.1745	73.68%	63.00%	0.7446	80.40%	71.00%	
Xception	1.2988	64.22%	54.70%	1.3424	64.81%	53.28%	DFDC
Res2Net-101	1.2052	69.51%	58.60%	1.6336	69.89%	62.70%	
EfficientNet-B7	1.2835	70.49%	59.13%	1.2726	64.29%	59.23%	
ViT	1.1135	69.87%	60.20%	0.8981	68.20%	62.98%	
Swin	1.2534	71.53%	58.63%	1.1194	74.04%	61.48%	
MViT	1.1775	69.63%	62.15%	0.9917	68.61%	58.40%	
ResNet-3D	1.2023	73.75%	68.00%	1.2788	67.04%	61.00%	
TimeFormer	1.1116	73.77%	68.00%	1.0129	74.04%	63.50%	

Table 13. Inter-dataset evaluation scores of models trained on DFDC [11] dataset and evaluated on the remaining three datasets.

Training Dataset: DFDC							Evaluation Dataset
Model	With Augs			No Augs			
	LogLoss	AUC	ACC	LogLoss	AUC	ACC	
Xception	1.4046	58.38%	55.25%	1.8346	60.31%	53.63%	FakeAVCeleb
Res2Net-101	2.0891	59.23%	56.33%	1.6953	59.77%	55.43%	
EfficientNet-B7	1.0800	65.31%	61.63%	1.0920	71.87%	65.40%	
ViT	1.2515	59.31%	56.00%	1.1361	60.12%	57.43%	
Swin	1.2053	67.81%	62.90%	1.2668	63.48%	60.25%	
MViT	1.2121	63.46%	60.05%	1.2139	65.75%	61.30%	
ResNet-3D	1.1114	63.19%	54.50%	1.2748	62.02%	56.50%	
TimeFormer	1.1582	65.34%	62.00%	1.6968	67.80%	61.00%	
Xception	1.1784	67.90%	61.25%	1.7465	64.95%	58.20%	CelebDF-V2
Res2Net-101	1.2293	83.01%	74.95%	1.1859	83.57%	72.35%	
EfficientNet-B7	0.8278	79.82%	70.15%	1.2972	74.27%	68.45%	
ViT	0.7301	85.62%	76.45%	0.9351	73.81%	67.25%	
Swin	0.8411	84.36%	76.60%	1.1246	80.34%	73.20%	
MViT	0.8548	87.83%	71.55%	0.7711	84.75%	76.75%	
ResNet-3D	0.7638	79.60%	72.00%	0.7806	77.88%	72.00%	
TimeFormer	1.0558	76.48%	71.00%	1.3635	79.60%	74.00%	
Xception	1.3048	64.59%	60.95%	1.4530	67.29%	62.70%	FaceForensics++
Res2Net-101	2.7490	61.15%	56.48%	2.3375	65.15%	60.78%	
EfficientNet-B7	1.2085	68.82%	64.08%	1.4401	70.71%	65.83%	
ViT	1.6779	66.43%	59.13%	1.6781	60.19%	57.28%	
Swin	1.6806	68.32%	63.65%	1.7493	64.66%	60.20%	
MViT	1.6317	65.82%	61.55%	1.6911	67.54%	62.88%	
ResNet-3D	1.5308	57.27%	57.86%	1.3507	58.91%	57.50%	
TimeFormer	1.2122	70.17%	66.79%	1.9148	67.90%	60.00%	

B ANALYSIS OF RESULTS ACHIEVED BY SELF-SUPERVISED MODELS

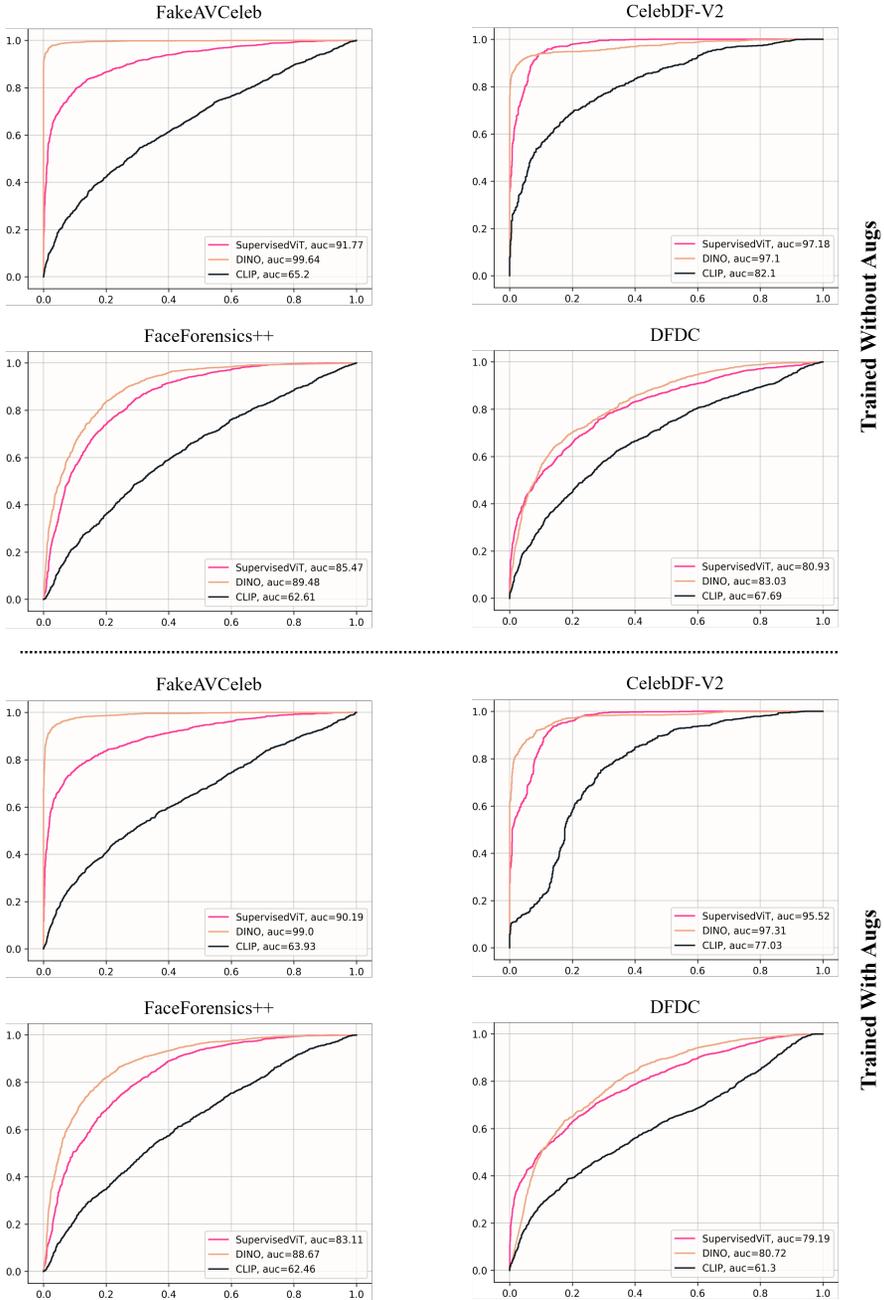


Fig. 8. ROC curves of self-supervised models trained and evaluated on each dataset using the intra-dataset evaluation scheme.