

TLDR

DMDF_Faces_V2 is a collection of deepfake detection datasets intended to provide you with a ready to use toolkit for deepfake detection. DMDF contains thousands of crops of over 50,000 real and deepfaked videos, see the split below:

	Real Videos:	Fake Videos:
DFDC	8601	17,093
VoxCeleb	13,680	0
Deeper_Forensics	2234	633
CelebDF	192	2985
Deepfacelab	0	253
Face_Forensics	777	782
Face2Face	0	790
FaceShifter	0	785
FaceSwap	0	798
NeuralTextures	0	655
GoogleDF	0	1438
Total:	25,484	26,212

For the full breakdown across face sizes and train/test/validation splits, please see the DMDF V2 Distribution doc available for free on the [DMDF github](#).

To download the full DMDF_Faces_V2 datasets use the command:

```
aws s3 sync s3://dmdf-v2 . --request-payer --region=us-east-1
```

When you download DMDF_Faces_V2 you will receive a directory pre-structured to accommodate a deepfake detector. The folder is subdivided into test, train, and validations sets, as well as fake and real within each set. Within each fake and real directory, the data is further subdivided into its original dataset. To set up a deepfake detector, point the data loader to the fake and real data, and the individual crops are already pre-processed.

What's new in V2? We are grateful for the hugely positive response to V2, and we hear your feedback loud and clear. Firstly, DMDF V2 is significantly smaller in download size, making it a much more easily accessible and downloadable tool for deepfake detection. Secondly, despite the decreased overall size, we have added dramatically more data to DMDF V2. DMDF V2 now contains over 50,000 pre-processed images, split almost exactly between original and manipulated videos. Lastly, we have begun doing demographic analysis of the large quantity of data we are providing, allowing users to more directly understand the racial, emotional, gender, and age distribution of the data they are using.

The dataset size is 350 GB(Compared to 2.7 TB for V1), and you must be authenticated with your Amazon Web Service account. By setting your region to US-east-1 you can avoid download costs, otherwise it is \$0.02 USD/GB, or about \$55 for the whole dataset. If you are in academia, or find this cost prohibitive, and are unable to set your region to US-east-1, we would be more than happy to subsidize your download. For more information please reach out to ryan@deepmedia.ai.

Abstract

With the creation of synthetic media such as deepfakes increasing in accessibility, quality, and ease, the challenge of accurately and reliably detecting these deepfakes grows dramatically. The variety of tools and synthetic media generation techniques, and the rapidly developing nature of deepfake generators make the work of deepfake detection difficult, but crucially important. A number of tools have been developed in recent years to begin the process of creating accurate and fast deepfake detectors, however all of these tools fall prey to the same fundamental problem. Though previous deepfake detection methods have achieved high accuracies, they have not yet been tested on the modern techniques for deepfake generation. In short, previous detectors have been trained on low quality deepfakes.

We at Deepmedia believe that to combat the threat of deepfakes, we need accurate detectors, trained on the most high quality synthetic media. Even moreso, we truly care about the ethical use of technology and protecting truth on digital platforms. For this reason, we have built and are releasing for public and academic use DMDF_V2, a dataset that is a collection of the highest quality and widest range images for deepfake detection. Combining real and fake sources from a number of datasets and deepfake generation techniques, and performing cropping and pre-processing of these images, we present DMDF as a highly powerful and ready to use toolkit for the creation of advanced and accurate deepfake detectors.

Version 2 improves upon the well received version 1 by dramatically reducing the overall size of the dataset while increasing the overall number of videos included significantly. We accomplished this by selecting two versatile and useful facial crops to include, rather than include every possible facial crop. Secondly, V2 contains a demographic analysis of the ages, races, genders, and emotional distributions of the faces included in DMDF. We believe this type of analysis is crucial toward lessening the dominance of White and Masculine faces in currently available facial data. Volume 1 of this dataset contains faces. Future versions will be coming soon, including additional data and other modalities of deepfake detection such as voice, text, and aerial imagery.

Introduction

To best understand how to use DMDF, it is first important to know its individual components. DMDF_V2 contains five different real datasets, and ten total fake datasets generated from these real images across a number of deepfake detection techniques. Let's begin by going over the real data.

Real datasets used in DMDF_Faces_V2

The first dataset employed is Meta.AI's [Deepfake Detection Challenge Dataset](#) (DFDC). DFDC contains a wide and robust dataset of both real and fake videos. DFDC's greatest advantage is that it boasts an incredibly large dataset, across a wide number of paid actors. It is important to note however, that by the standard of deepfake technology, DFDC is already quite old, and should not be used as the sole metric of any detector.

The second dataset utilized is [Deeper Forensics](#), a dataset constructed specifically for the purpose of facial forgery detection. Deeper Forensics contains high quality videos, useful for both deepfake detection and generation, with actors spanning 26 different countries. Though it is powerful, deeper forensics contains only 100 faces, and should likely be combined with additional data to achieve the diversity important for deepfake detectors.

The third included dataset is [VoxCeleb](#), a massive audio visual dataset constructed for a wide range of facial and vocal analysis purposes. Containing over 2,000 hours of celebrity facial videos, VoxCeleb proved to be a crucially useful training resource for Deepfake detection. It is important to note however, that despite the massive amount of data provided, the video quality is not particularly high, and it is important to use VoxCeleb in conjunction with a higher quality real dataset such as Deeper Forensics. VoxCeleb also contains no synthetically manipulated videos, and thus a deepfake detector can never be trained on VoxCeleb alone.

The fourth dataset used by [CelebDF](#), a large scale dataset containing videos of celebrities intended for challenging deepfake forensic analysis. Though there are only 590 real videos included in the dataset, they are of a high quality, and there are many more synthetically manipulated videos included. For this reason, CelebDF ought to be thought of as the perfect complement for a dataset such as VoxCeleb, which has massive amounts of real data and no manipulated data.

Lastly, our final dataset is [Face Forensics ++](#). Though the individual video quality of Face Forensics is rather poor, it boasts a wide and robust range of content. Face forensics contains data from both youtube sources and outside Mp4 videos.

Next, it's important to understand the sources of synthetically manipulated data present in DMDF_Faces_V2

Fake datasets used in DMDF_Faces_V2

DFDC fakes are created using one of eight different facial modification algorithms. For this reason, training a detector on DFDC fakes can be useful and provide a test for the detector over a wide range of modification techniques. The fakes however, are perhaps where DFDC's age can be most felt, and it should never be taken as the modern standard for deepfake generation.

Deeper Forensics provides perhaps the best quality deepfakes of this version. These videos are high resolution, and employ a technique referred to as a DeepFake Variational Auto-Encoder (DF-VAE), providing for high quality face swapping.

CelebDF boasts a very impressive collection of high quality deepfake videos, using a technique of facial synthesis to provide thousands of compelling synthetically manipulated videos.

[Deepfacelab](#) is considered by many to be the premier deepfake generation tool publicly available. Not only are the videos incredibly compelling to the naked eye, but they are difficult and robust for even the most well trained detector. For this reason, despite the lack of a publicly available deepfacelab dataset, Deepmedia has collected approximately 250 deepfacelab clips to include in DMDF V2.

[Google Deepfakes Dataset](#) contains over 3000 high quality synthetically manipulated videos. Contained within the Face Forensics family, GoogleDF provides high quality videos generated from actors hired particularly for this purpose. For this reason, these videos are ideal for a deepfake detector, and these

crops should be considered of the highest quality available for deepfake detection purposes.

Face Forensics++ provides an incredibly useful tool for deepfake detection, one set of videos trained on 5 different deepfake generation techniques. To briefly summarize, the Face Forensics fakes are split into:

Face2Face which utilizes real time facial reenactment. **Face Shifter** which utilizes a high quality face swapping technique. **Face Swap** which utilizes facial remapping technology. **Neural Textures** which implements facial reconstruction using feature maps learned from real data. And lastly **Face Forensics** provides its own generated fakes using this real data.

To include a brief note on sampling, the combination of these datasets provides a rather huge completed collection. This remains the case even as we have only included 10% of the most massive datasets for this release, DFDC and Deeper Forensics. For this reason, using the entire dataset is almost certainly unnecessary. Rather, sampling 10% of the total dataset will provide a more useful and compact package.

Methods

What DMDF provides is pre-processed image crops generated from these datasets. It is these crops that will effectively be able to train a deepfake detector. Here, we quickly overview our pre-processing techniques.

Firstly, DeepMedia implements a technique of trajectory analysis of each processed video (patent pending). This allows us to save crucial data as we break down the video into individual crops. Secondly, we align each image crop such that it can be centrally observed by the detector. For each DeepFake and real video, DMDF includes two different image crops: ffhq-align-big and bounding-box-tight-v2. Additionally, the dataset includes 68pt landmarks for each video frame, the original audio samples, and video metadata. Also included is deepface json files, which contain demographic information for each video such as predictions for age, race, gender, and emotion of each speaker. Though we have collected a number of sizes and types of image crops, we have found “ffhq-align-big” cropping/alignment to be most effective. This both assists the detector, and allows for further customization of the crops by leaving plenty of space for further processing with DeepMedia provided facial landmarks if desired. Lastly, as many DeepFake videos contain multiple faces, many which have not undergone synthetic manipulation, we ensure that each face in the “FAKE” set has actually been synthetically altered, solving a problem which has been known to dramatically harm detector quality in published research.

To ensure the ready application for this dataset to train Deepfake Detectors, we have further subdivided the data into a ready to use structure. Firstly, we have pre-constructed a train, test, and validation set of images. Within this, we have subdivided the image crops into real and fake, each containing the respective crops for each dataset within. Lastly, the images are split into image size bins, and audio files are present wherever available. Through this combination of techniques and pre-processing methods, the image crops provided in DMDF serve as a ready to use and hugely practical deepfake detection toolkit.

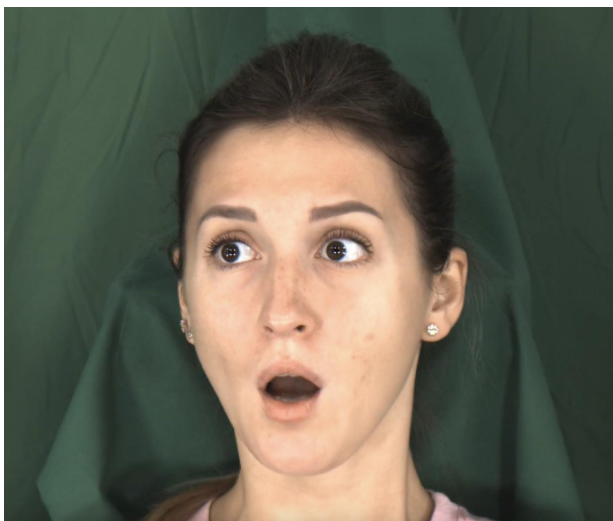
Results

Though future updates will contain detectors trained on the entirety of DMDF, for this release we hope to share our review of these datasets after many hours working with them. Ultimately, all of these datasets have utility in the construction of effective deepfake detectors, it is however important to use these datasets to their strengths.

GoogleDF provides one of the only sources of real videos intended directly for deepfake detection purposes. For this reason, the produced synthetically altered videos are of a high quality, and well positioned to be used by a deepfake detector. Anyone hoping to train a detector should include GoogleDF without hesitation, but note that it is likely not expansive enough to be used without the addition of other synthetic datasets.

Deeper Forensics provides a powerful and high quality deepfake detection dataset. It is high resolution, and though they proved useless for our detection purposes, they are further subdivided into emotion and camera angle. Further, the usage of their DeepFake Variational Auto-Encoder results in some impressive deepfakes that should be confidently sent to challenge a detector. Simply keep in mind its limitations, namely that it contains only 100 total faces, and Deeper Forensics provides a powerful tool for deepfake detection.

Attached below are two deeper forensics stills that illustrate the high quality of their source data:



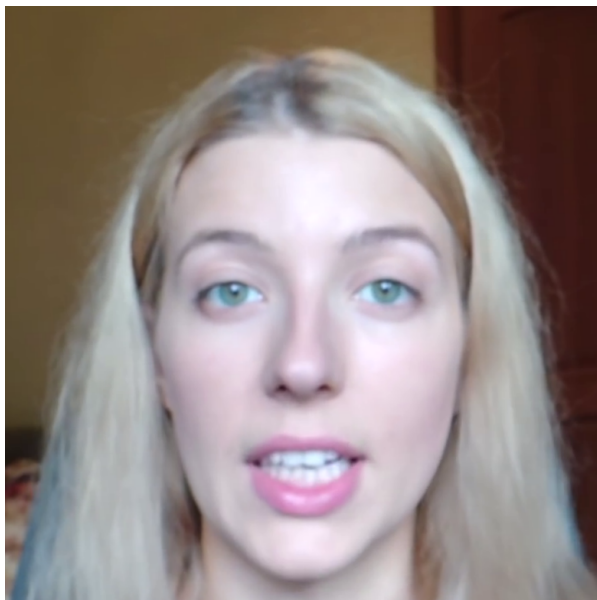
CelebDF has proved itself to be one of the most powerful sources of synthetically altered videos available. Using a rather simple facial swapping technique, the altered videos provided from CelebDF are compelling and rather extensive. CelebDF should be included without worry in any deepfake detection endeavor.

DFDC epitomizes a robust and extensive approach for deepfake detection. Boasting a practically unlimited number of real and altered images, DFDC can effectively provide the bulk of a training set. Furthermore, the usage of 8 different facial modification algorithms ensures that a detector will not be singly useful, but rather have a wide range of ability to detect media altered by any number of techniques. In the age of deepfakes however, three years is an eternity, and any detector trained on simply DFDC will be left powerless in the face of the best modern deepfakes. For this reason, consider DFDC as a supplement to fill any additional desired training data.

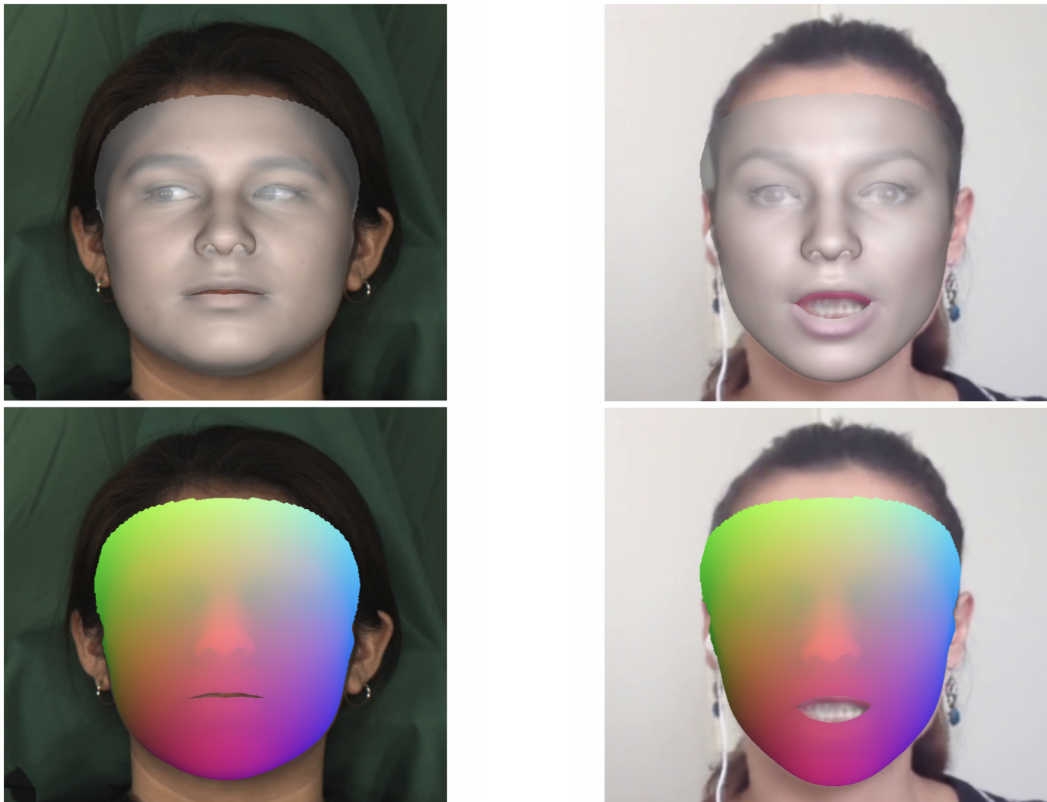
VoxCeleb serves as a perfect complement to DFDC, in that it similarly contains an extensive collection of real videos, but each video is not of the highest available video quality. Combined with a dataset such as CelebDF or GoogleDF, VoxCeleb remains a powerful and useful source of near endless real videos.

Face Forensics is positioned uniquely amongst these other datasets. It is relatively low quality, which is to be expected given its age. Nonetheless, Face Forensics contains an incredibly powerful resource, the same images processed through 5 different facial manipulation techniques. For this reason, consider Face Forensics as a practical and useful test for any detector, to test for accuracy and its ability to recognize synthetic media of various different approaches.

Below is an example of a base face forensics source image(left) and a face forensics deepfake(right). Observe the lower image quality compared to a data source such as deeper forensics.



Following our processing of the individual datasets, we processed certain images through a 3D facial alignment algorithm(3DFFA) designed by [Jianzhu Guo](#). Upon processing these images, it became clear that the real images generated significantly more accurate facial alignment masks than their fake counterparts. This kind of facial alignment analysis can prove as a useful tool in complement with our machine learning based deepfake detectors. See below a 3D facial morph as well as heat maps for a real image from Deeper Forensics and a deepfake generated using Neural Textures.



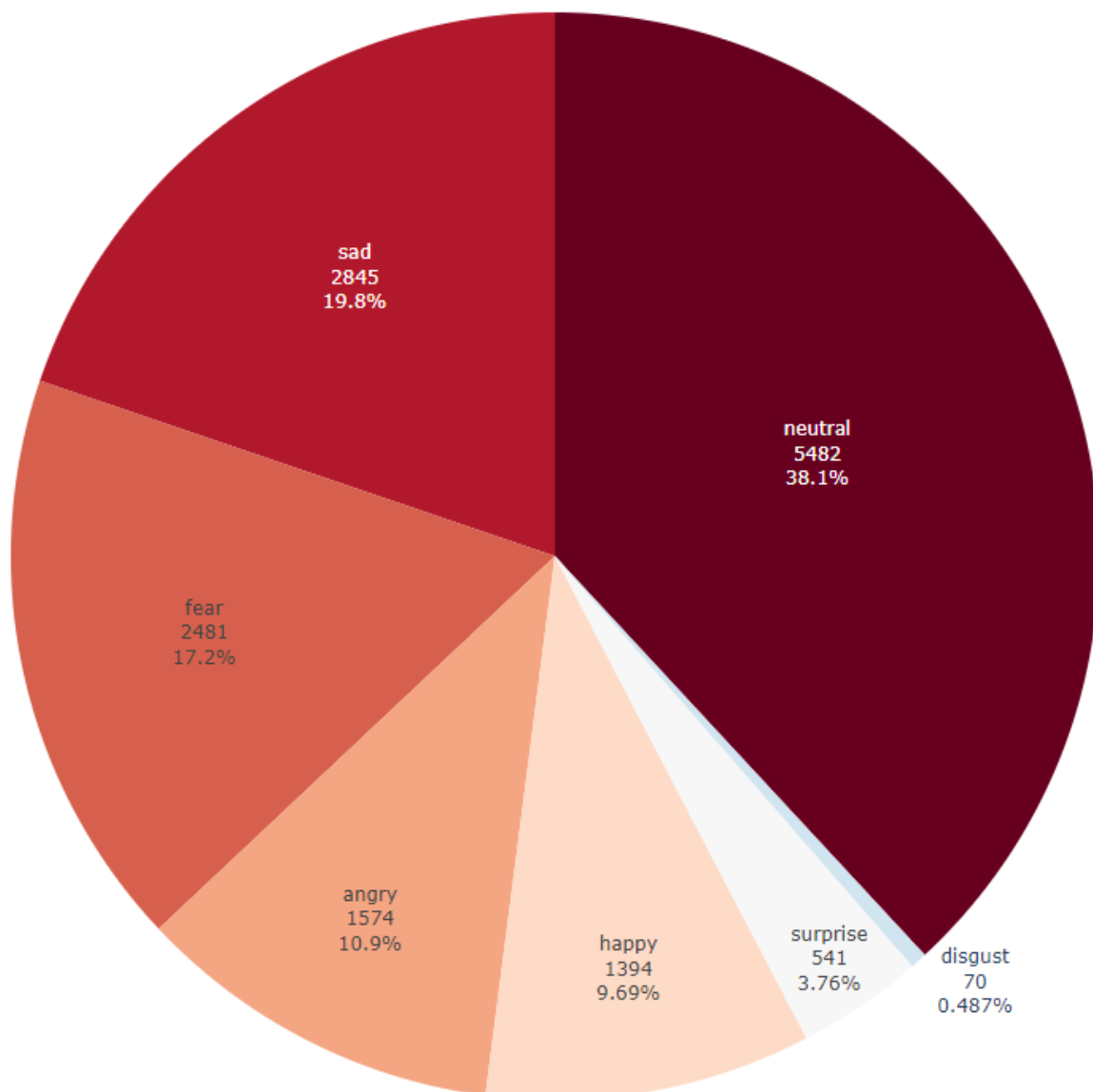
Demographic Results:

In addition to creating powerful and easily accessible deepfake detection resources, Deepmedia is committed to working toward a more equitable and diverse use of data. Simply put, currently available facial and synthetically altered facial data lean heavily white and masculine. There are a variety of reasons this occurs, but one thing is abundantly clear, our deepfake detectors must be able to function independent of age, race, and gender. The only way to accomplish this is to be intentional and critical about the data we select when training machine learning networks. This is certainly not going to happen overnight, but we at Deepmedia hope that this demographic report can make clear the need for more diverse data, and provide the awareness of this disparity within our facial datasets.

For the training set of each dataset, we have constructed visual and textual resources that outline the split of gender, age, race, and emotional distributions for each dataset included in DMDF. It is important to note that we do not use all of the data from each dataset, and therefore this is not necessarily indicative of the demographic split of each collection of data. However, what data was selected was done so agnostically of gender, age, race, and emotion. Therefore, the following results should be seen as at least somewhat representative of the distributions of the original datasets.

To view a spreadsheet containing the exact distribution of age, race, gender, and emotion, please see the DMDF V2 Demographic Information document [here](#).

Attached below is a sample static image chart for emotional distribution, the full collection of interactable data visualizations for Age, Race, Gender, and Emotion in HTML format can be found in the visualizations folder of the DMDF github.



Applications

We see the applicability of DMDF to be threefold.

Firstly, a robust and powerful deepfake detection dataset will be crucially useful to academia and DoD/IC. Creating a dataset containing powerful and high quality deepfakes will advance the work of academia and make clear the threat of synthetic media to national and international security.

Secondarily, we envision this dataset to be crucial in the creation of more advanced deepfake detectors. Applying a wide range of quality deepfakes will serve to raise the stakes on what our current detectors are capable of, and emphasize that detectors trained on easy mode simply are not good enough.

Lastly, DMDF provides a useful and powerful tool for comparing and better understanding the difference in deepfake generators. Understanding their strengths and weaknesses will be crucial in creating detectors most effective at processing this high quality synthetic media.

Discussion

Having worked with thousands of deepfakes, and dozens of different approaches to their generation and detection, one thing is made abundantly clear. Deepfakes are getting better, and they are doing so faster than our detectors can keep up. Previous detectors simply cannot keep up with the advancements to deepfake technology, and they are showing no signs of stopping. As these deepfakes improve, they pose all the more potent threats to truth online. It is for this reason that it is crucial to prioritize and collaborate on the work of deepfake detection.

At DeepMedia, we are committed to continue and advance the work of deepfake detection, and that involves creating modern detectors trained on the highest quality synthetic media, and made rapidly adaptable for the improving quality of deepfakes that are generated. We believe that we can only do this with the help of academia and governments, and we hope that DMDF provides individuals with the tools necessary to begin to take on this important and difficult work.

It is also clear that while our datasets are improving in quality, the racial and gender disparity present in these datasets is not improving in turn. Overwhelmingly, deepfake detection resources feature masculine, white presenting individuals. We hope that performing the kind of demographic analysis presented in this paper will serve as a first

step in facilitating the creation of more diverse datasets, and we are committed to help in any way to generate more diverse and representative datasets.

Next Steps

DeepMedia is actively working to train detectors on all of DMDF, providing the community with benchmark data, accuracy, and loss metrics for the best detectors currently available. Firstly, we will be training a Cross Effective Vision Based Transformer(CE-VIT), and observing performance across DMDF.

Secondarily, we are constantly searching for the most modern and advanced deepfake generators available, and will be adding many additional deepfakes and source videos in the following version. We are most interested in the deepfake generators used in the wild, as they prove the most relevant and challenging test to our detectors. These techniques include but are not limited to Thin-Plate-Spline, FSM, SimSwap, and DeepFaceLab. We will be adding proprietary “second stage” DeepFake processors for the generation of 4K-quality deepfakes to be added in future versions of DMDF.

Lastly, we hope to add additional modalities beyond facial imaging to DMDF such as synthetic audio manipulations and satellite imagery where applicable to our pre-processing, allowing for a new generation of vocal deepfake detectors.