# DeepFake MNIST+: A DeepFake Facial Animation Dataset

Jiajun Huang  
The University of Sydney

Xueyu Wang  
The University of Sydney

Bo Du  
Wuhan University

Pei Du  
AntGroup

Chang Xu  
The University of Sydney

## Abstract

*The DeepFakes, which are the facial manipulation techniques, is the emerging threat to digital society. Various DeepFake detection methods and datasets are proposed for detecting such data, especially for face-swapping. However, recent researches less consider facial animation, which is also important in the DeepFake attack side. It tries to animate a face image with actions provided by a driving video, which also leads to a concern about the security of recent payment systems that reply on liveness detection to authenticate real users via recognising a sequence of user facial actions. However, our experiments show that the existed datasets are not sufficient to develop reliable detection methods. While the current liveness detector cannot defend such videos as the attack. As a response, we propose a new human face animation dataset, called DeepFake MNIST+[1], generated by a SOTA image animation generator. It includes 10,000 facial animation videos in ten different actions, which can spoof the recent liveness detectors. A baseline detection method and a comprehensive analysis of the method is also included in this paper. In addition, we analyze the proposed dataset's properties and reveal the difficulty and importance of detecting animation datasets under different types of motion and compression quality.*

## 1. Introduction

DeepFake[2] has become a critical topic for our digital society. With DeepFake, we can now easily change the identity or expression of the face in an image or a piece of video with another person's identity or expression. The mainstream DeepFake techniques are based on Deep Neural Networks (DNNs), especially Generative Adversarial Networks (GANs) [17], to produce visually plausible images or videos which are hard to be discriminated by humans. There is growing concern about DeepFake, as malicious people could use the techniques to palm off the victims and illude presence and activities, even if they never did before.

A number of DeepFake methods have been developed to manipulate the attributes of human face in images or videos. For example, the swapping methods [6, 5, 4, 8, 3, 25] mostly focus on the identity of the face and try to replace the face in one image/video with the face from others. The deepfake method [5], a famous swapping algorithm, trains identity-dependent two auto-encoders to swap the faces of two identities. Besides face identity, many other face attributes have been studied in the literature. For example, [34, 33] modified the expression in one face image/video, and facial image animation [11, 37, 31, 32, 12], as a compose of expression manipulation, is increasingly being important within the DeepFakes. Given a face source image and the driving video, DeepFake can now generate a new video where the source face performs the same action as the driving video. For instance, Siarohin et al. [31] use an encoder to capture optical flow information from the videos, embedding the information with source images, and generate videos. Zakharov et al. [39] pass the identity embedding to the image generator to produce manipulated face with the given landmark. Burkov [12] extracts the identity and pose information separately and generate videos with embedding.

Given the growing anxiety on the high-quality generation by DeepFake and its potential negative social impacts, it becomes especially urgent to study the defense techniques against DeepFake. Recently a few datasets have been created for the study of DeepFake detection methods. UADFV [38] and Celeb-DF [27] collect youtube videos to generate face-swapping videos, and DFDC [16] captures 48,190 videos with paid actors and generates a large scale swapping dataset with over 104,500 videos. By analyzing these datasets, Rossler et al. [30] and Seferbekov [7] suggested the importance of CNN architectures on detecting the swapped faces. Li et al. [26] further found the manipulate boundary between the face and head can be an effective clue for the detection. However, all these works are mostly about detecting the face identity change, which is a limitation of existing deepfake datasets. There is rare dataset or deepfake detection work about facial animation, though it occupies a significant part in DeepFake attack side.

Recently, facial animation by DeepFake has been deployed on mobile devices, producing a large number of fake

---

[1] https://github.com/huangjiadidi/DeepFakeMnist

[2] DeepFake not only indicates the facial modification methods but also is the name of one algorithm called deepfake [5].

videos and broadcasts through the Internet [1]. These fake videos thus challenge the security of many intelligent systems in our daily life. For example, the face recognition based payment systems usually rely on liveness detection to verify whether users are the real people by requiring them to do a sequence of specific actions in videos. However, our experiments show that DeepFake detectors trained on existing deepfake datasets consisting of face identity changes are not applicable for detecting facial animation videos. Further, we observe that the SOTA liveness detector in public cannot defend the animation data with specific actions as they claimed.

In this paper, we propose a new dataset, called DeepFake MNIST+, a human face animation video dataset. The DeepFake MNIST+ dataset is developed as a response to the wide use of the MNIST dataset, but for a different deepfake detection problem. We create this dataset to provide the basis for learning and practicing how to develop, evaluate, and use deep neural networks for fake face animation detection. The dataset contains 10,000 face animation videos in ten different actions, *plus* 10,000 real face videos to enable a supervised detector training. These fake facial animation are of higher fidelity and able to spoof the popular liveness detectors on the market (as of the time of this manuscript submission). Given the gap of different deepfake types, the detectors trained on existing DeepFake datasets with face identity change cannot well detect fake animations in the proposed dataset. We establish deepfake detection baselines on the DeepFake MNIST+ dataset and carefully evaluate their performance in different scenarios.

We also present comprehensive analysis related to the properties of proposed dataset. We explore the impact of motion type and compression quality of generated videos. As observed from Figure 5, the actions with large movements will challenge existing detectors. But by taking these difficult DeepFake data into account, the learned detector can enjoy a much better generalization across other types of DeepFake data. The relevant discussion and empirical studies in this paper therefore shed new light on both DeepFake generation and detection research.

## 2. Related Works

### 2.1. Image Animation

The interest in facial manipulation methods has rapidly increased recently. One approach for manipulation is based on 3D modeling. Zollhofer et al., [42] build a 3D morphable model for the source face, to perform realistic animation for given actions. Suwajanakorn [32] attempt to model lip to forgery talking. The deepfake [5] introduce a DNN based face-swapping method, replacing faces within two identities with two encoders. Although it requires plenty of videos/images of both identities to achieve better results, the promising results show the potential of manipulation with DNNs. The recent researches focus on identity-independent swapping methods. Li et al. [25] implement two encoders to extract attribution and identity information and embed them in GAN to generate high fidelity swapped results. Several DNN based expression manipulation and animation are also proposed. Thies et al. [33] consider facial reenactment as a domain transfer problem using Pix2Pix architecture [22] to produce results. Siarohin et al. [31] extract the motion information of driving videos with optimal flow estimation to generate high-quality animation results.

### 2.2. DeepFake and Liveness Detection

An increasing number of DeepFake detection methods are proposed as a response to the huge concern from society. The approaches could be separated into three categories. The first type is trying to detect the unnatural section of the manipulated videos, such as swapping boundary [26], inconsistent head angles between the face and head [38]. The second type detects the synthesis signal of GAN to distinguish DeepFake data. For instance, Wang et al. [35] observe GAN signatures using discrete cosine transform for detecting CNN-based DeepFake samples. The last categories' approaches rely on DeepFake dataset to train the detectors [30, 7, 18, 29], regardless of the inconsistent or signals.

On the other hand, the recent liveness detector mainly defending the psychical level attacks called Replay Attacks. For instance, attackers could build a 3D mask of victims through a 3D printer or print out victims' face in the paper and wear it. To defense against such attacks, many detection methods have been proposed. Some approaches try to detect the differences between real faces and forgery faces. De et.al., [14] estimate the invariant of facial points for detection. Komulainen [24] believe detecting faces' dynamic muscle change can distinguish the spoofing. Wang et al., [36] detect the blood flow change under the skin to separate facial mask and real face. The recent payment or identity verification solutions with smartpohone usually combine pose verification with liveness detection, such as Alipay. They usually require users to do some specific action, such as blinking or yawing, to improve the performance. However, our experiments show that the facial animation data with specific actions could spoof the SOTA liveness detectors in the market.

### 2.3. DeepFake Datasets

A few DeepFake datasets have been proposed recently as a response to the increasing concerns on DeepFake techniques as they can generate realistic results to spoof people. However, most of the datasets are generated using identity swapping algorithms, but only a few works for facial animation. Table 1 shows the details of these DeepFake datasets.

**Celeb-DF [27]**: The Celeb-DF is a large face swap

dataset. It selects 590 real videos from Youtube, which all about the talking videos of celebrities. The dataset contains 5,639 synthesis videos using a high-resolution face generator.

**DFDC [16]** : The Facebook DeepFake detection challenge dataset is one of the largest face-swapping video dataset recently. It contains 104,500 face swap videos based on 48,190 source videos shot with 3,426 paid actors in different locations and light conditions.

**FF++ [30]**: The FaceForensics++ dataset is one of the popular DeepFake datasets. The autoher use multiple Deep-Fake methods to generate the dataset. The dataset to provide the data generated by expression manipulation, the related techniques with facial animation. Two manipulation methods, Face2Face [34] and NeuralTextures [33], are implemented to generated DeepFake videos, while the two face-swapping methods [6, 5] are also included. The dataset uses 1,000 real videos from Youtube to generate 1,000 DeepFake videos for each method.

the proposed dataset is a contemporaneous work with Forgerynet [21]. Our proposed dataset has several differences comparing with it: i) the proposed dataset includes animation data under ten specific categories (e.g., head movement and emotion changes), rather than applying the unknown action from a random video; ii) we boost the quality and challenge of the proposed dataset by filtering the generation with liveness detection; and iii) we present a comprehensive analysis about the proposed dataset, including action categories and video quality.

According to Table 1, the only dataset involving the face animation is FF++, which is a small dataset and is hard to cover the challenging deepfake data in the real world. We argue that the face animations in a deepfake dataset should be diverse enough and are better to cover the animation categories in the prospective downstream tasks, instead of the just causal talking in FF++. For example, the liveness detectors often require specific actions or expressions of the face as the input. For these reasons, it engages us to propose a large-scale and action-specific facial animation dataset.

| | #real | #fake | type(s) of generation | action specific |
|---|---|---|---|---|
| UADFV [38] | 49 | 49 | face swap | No |
| Celeb-DF [27] | 590 | 5,639 | face swap | No |
| DFDC [16] | 48,190 | 104,500 | face swap | No |
| FF++ [30] | 1000 | 4000 | animation & swap | No |
| ForgeryNet [21] | 91,630 | 121,617 | animation & swap | No |
| **DeepFake MNIST+** | 10,000 | 10,000 | image animation | Yes |

Table 1: Basic information of existing DeepFake datasets.



Figure 1: Facial animation video samples for different actions.

## 3. DeepFake MNIST+

The major contribution of this paper is our proposed human face animation video dataset, called DeepFake MNIST+. It includes 10,000 face animation videos performing ten different actions and 10,000 real human face videos selected from other datasets. Besides, all these animation videos can spoof the liveness detection solution in the market. Such that the videos are still challenging recent public detectors. It is the first large-scale dataset for face animation videos of variant actions to the best of our knowledge. We believe such a dataset allows us to train advanced detection models to distinguish the face animation videos for preventing spoofing.

### 3.1. Generation Model and Data Preparation

We select Siarohin's framework [31] to generate face animation videos. This SOTA animation framework taking as input a source identity face image and a driving video shot by another actor. The model performs local affine trans-

formations using first-order Taylor expansion to estimate the motion of the driving video. Then applying the motion features to the generator to provide high-quality face animation videos. As a result, it generates a video that animates the motion of driving video while keeping the identities of the source image. For more detail, please review the original paper. The major advantage of this framework is that the model is not identity-dependent. We can generate different videos for variant identities with arbitrary driving in the one trained model, while it only requires one image as the source image.

For the source identity images, we select the frames from video in VoxCeleb1 dataset [28]. VoxCeleb1 is a large-scale audio-visual dataset of human speech. It includes 1251 unique celebrities in 22,496 talking videos. All videos are face-cropped and have size 256x256 resolution. During our experiments, using the front face source images could achieve better generation quality. Therefore, we select the face frames mostly facing the front from the VoxCeleb1 dataset as our source images for generation.

The DeepFake MNIST+ dataset contains forgery videos in 10 actions. It includes: *Blink, Open mouth, Yaw, Nod, Right slope head, Left slope head, Look up and smile, surprise and embarrassment*. The driving videos of embarrassment are collected from ADFES dataset [19]. The dataset contains variant emotional expression videos (anger, disgust, fear, joy, sadness, surprise, contempt, pride and are shot by 22 actors. We select five actors' embarrass videos from the dataset as our driving video. The videos of the remaining actions are shot by one volunteer. The action videos are captured using a front camera of the iPhone 11 Pro, and each action has been executed by the volunteer for 5 times. All the driving videos are face cropped by using MTCNN modules [40] and resized into 256x256 resolution to align the format of the VoxCeleb1 dataset.

### 3.2. Generating High-Quality Facial Animations

We adopt two public liveness detection APIs to select the challenging samples that cannot be accurately recognized by the detector. The first one was provided by TianyanData [9]. TianyanData's API supports liveness detection with a specific action, including blinking, yawing, nodding, and opening mouth. In order to pass the detection, the input face video has to perform a particular action while passing the spoofing test. The second one comes from Baidu [2]. Their detector supports universal liveness detection regardless of actions. Both of these two companies claim their detector can achieve 99% accuracy for detecting spoofing.

We generate many animation videos for all actions and then pass the data to the liveness detection APIs to pick out the samples that are challenging for the liveness detector. As a result, the DeepFake MNIST+ dataset contains 10,000 face animation videos in 10 specific actions and 10,000

real face videos collected from VoxCeleb1. Each action includes 1,000 videos, and all of them can spoof the APIs. For blinking, yawing, nodding, and opening mouth, we use both TianyanData and Baidu APIs to filter videos. For the remaining actions, we use Baidu's API to collect spoofed data. The following graph presents each action's spoof rate by passing the animation videos to the two APIs.

The actions of blinking, yawing, nodding and opening mouth have lower average spoofing rates than others, such the situation could be caused by two APIs filter the videos of those actions. In addition, the TianyanData's API requires further action detection, which reduces the chance to attack. On the other hand, the actions that require large-angle head movement, e.g., yawing and nodding, the success spoofing rates are much lower than the other actions that don't need a significant motion change. One reasonable explanation could be that a single source image of the frontal face cannot provide sufficient detail of all head information, e.g., profile face, leading to lower head movement quality. The videos of simile have the highest spoofing rate, which has achieved 61%. It might because the smile action doesn't lead to significant head change, making it hard to detect the spoofing details. While the yaw videos are more likely to be detected, that has a 23% successful rate only.

|  | DeepFake Mnist+ | deepfake | NeuralTextures | Face2Face |
|---|---|---|---|---|
| original accuracy | 96.58% | 99.7% | 99.2% | 98.9% |
| accuracy on DM+ | - | 43.7% | 63.6% | 67.5% |
| fine-tuning | 95.3% | 98.46% | 98.1% | 98.49% |

Table 2: The performance of Resnet50 models trained with existed datasets from FF++ [30] to detect our proposed dataset. And the performance of DeepFake Mnist+ trained model, fine-tuning with FF++ data, to detect all these datsets.

In addition, we explore the transferability of the detector trained with existing datasets. We train Resnet50 [20] models with the data of deepfake [5], Face2Face [34] and NeuralTextures [33] provided by FF++ [30] and present the result for detecting DeepFake Mnist+. The first one is the face-swapping data. The last two are the expression manipulation dataset. The Table 2 presents the result. All these three models are trained with raw quality and achieve nearly 100% accuracy in their own dataset. The result shows that the face-swapping detector fails to distinguish our proposed dataset. It seems better in manipulation datasets but still has a huge gap compared with the performance in their own datasets. Furthermore, we fine-tune the DeepFake Mnist+ trained model with the three FF++ datasets. The result indicates that the model can gain the power to detect both animation data and other Deepfake data by using our proposed dataset with other data sets.

Based on the these results, we believe the current detectors still cannot defend against such attacks. These obser-

vations engage us to proposed a face animation dataset to improve the detectors and achieve better security.
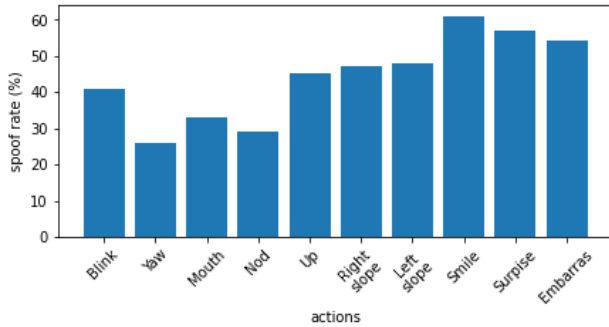


Figure 2: The spoofing successful rates for different actions. The first four actions' videos are detected by TianyanData [9] and Baidu [2] API, while the later six actions' videos are detected by Baidu API.

## 4. Benchmark Systems

We proposed a simple detection pipeline to detect the forgery videos in our dataset and distinguish them from those real videos. The forgery detection can be formulated as a binary classification task. We take the 10,000 real videos in the dataset as positive, while those forgery videos as negative.

In practice, the videos are usually suffering compression before uploading to the Internet, or the video might be shot by poor camera equipment, which suffering low video quality. Different compression rates are considered in our experiment to simulate the realistic detecting setting under different video qualities. We select two different compression rates - C23 and C40 under H.264 codec to compress the videos. To be mentioned that, a higher compression rate indicates worse video quality.

We exploited multiple models to accomplish the Deep-Fake detection task in our experiment:

**MesoInception-4**: MesoInception [10] is a CNN model consisting of two inception modules inspired by Inception-Net [10]. The model uses mean squared error between true and predicted labels rather than the ordinary cross-entropy loss. Following the training procedure in [30], we extracted the frames as the original size which is 256x256 resolution.

**XceptionNet**: The XceptionNet [13] is a traditional CNN model based on separable convolutions with residual connections. The model has shown high accuracy when detecting deepfake [5] videos and has been the baseline model introduced by Rossler et al. [30]. We used a pre-trained model on the ImageNet [15] dataset in this experiment. The CNN layers are frozen, and we only update the weights of newly inserted full connected layers for the prediction. Same with MesoInception, we also used the 256x256 resolution of frames as the input.

**Resnet**: The Residual Neural network (Resnet) [20] is one of the most popular neural networks. It utilizes skip connections or shortcuts to jump over some network layers, such that the networks are easier to be optimized even with the increasing depth. The Resnet networks are also pre-trained with ImageNet [15]. Three versions of Resnet - Resnet50, Resnet101, and Resent152 are included in the experiment to explore the performance change under different network depths. For all versions of Resnet, we capture and resize the image frames into 224x244 resolution.

All CNN models are trained with Adam optimizer with an initial learning rate of $0.0002$, and we set $\beta_1 = 0.999$ and $\beta_2 = 0.9$. The learning rate decreases under the poly-decay schedule. The total number of training epochs for each model is set as 50, and the batch size is 64. We pick 70% of videos from the DeepFake MNIST+ as the training data, that is, 700 videos from each of the ten action categories and 7000 real videos. 15% videos are the validation data. The remaining 15% will be the test data. We select the best versions of models based on the accuracy on the validation set.

## 5. Evaluation

This section presents different models' performance changes for detecting our proposed face animation video dataset under different situations. We present different models' accuracy for classifying test video data frames as real or fake ones to show the performance.

### 5.1. Overall Detection Performance

Table 3 compared the accuracy of different models under three different video compression levels (raw, light and heavy compression). The result indicates that the Resnet models have the best performance among all models. They achieve a 96.3% average accuracy on the raw video dataset, which is much higher than those of the other two CNN models. The XceptionNet, that also has a deep architecture, approaches a 92.38% accuracy for detecting forgery videos. In addition, the MesoNet shows the poorest performance with only a 60.39% accuracy for detecting the raw videos. The increasing of the network depth does not constantly boost the performance. The accuracy of three Resnet variants are very similar. It is hard to say, the deeper architectures, e.g., Resnet152, show improvement compared to other smaller models.

The video quality is also an important factor affecting the performance. Table 3 shows that worse video quality (higher compression rate) could lead to a performance downgrade. Compared to the raw dataset, the average accuracy of Resnet networks decreases around 2% under the light compression condition. The situation is worse when the videos are under heavy compression. 91.06% of compressed videos are classified correctly with Resnet networks on average, which has a 5% gap from that on the raw dataset. The XceptionNet

Table 3: The accuracy of different classifier models in testing set under different compression rate. The **light compression** corresponding to the c23 compression rate, while the **heavy compression** corresponding to the c40 compression size. We select the best version of models based on the validation accuracy during the training process.

| | Resnet50 | Resnet101 | Resnet152 | XceptionNet | MesoNet |
|---|---|---|---|---|---|
| raw | 96.58% | 96.64% | 96.18% | 92.38% | 60.39% |
| light compression | 94.32% | 94.87% | 94.90% | 85.52% | 58.58% |
| heavy compression | 91.49% | 90.44% | 91.27% | 83.143% | 57.90% |

| | Raw ->LC | Raw ->HC | LC ->Raw | LC ->HC | HC ->Raw | HC ->LC |
|---|---|---|---|---|---|---|
| Resnet50 | 85.52% | 71.3% | 95.16% | 82.98% | 59% | 61.02% |
| Resnet152 | 79.23% | 68.12% | 82.33% | 73.60% | 76.72% | 76.69% |
| Xception | 79.17% | 71.83% | 87.89% | 68.99% | 67.38% | 62.9% |
| MesoNet | 51.37% | 56.92% | 57.01% | 49.53% | 55.65% | 58.64% |

Table 4: The models' performances trained by one specific compression rate and detecting the videos from the datasets of the other two compression rates.

is more sensitive to compression rate than other networks, whose accuracy drops to 85.52% significantly when applying light compression and 83.143% for heavy compression.

The decreasing correction rate could be caused by the loss of detail in low video quality. The compression process leads to blur frames, hiding the forgery information so that the detector might not be able to capture such information for detecting the forgery areas.

## 5.2. Analysing the Impact of Video Quality

We further explore how the video quality could affect the performance. The Table 4 shows the models' generalization for the videos of different compression rates. It presents the accuracy of models trained with one quality (e.g., Raw) and predicts the data with the other two qualities(e.g., light and heavy compressed). The results indicate that the models cannot adapt well to datasets with other qualities, especially between raw and heavy compression datasets. The raw video models only achieve 70% accuracies on average for heavy compression videos and 67% conversely. It might show the heavy compression could change the data distribution leading to a dramatic accuracy drop. Also, light compression models have better generalizations than others, indicating the models could learn animation and compression information simultaneously.

One way to overcome the impact of different video quality is to train the network with the videos in all qualities. We train the Resnet50 and Resnet152 with the mixed video quality dataset for the experiment. More specifically, we select the video quality randomly for each sample to train the models and present their performances in testing sets under three qualities respectively. The Table 5 shows the result. With training in mixed quality videos, the models can adapt to different video qualities. However, it still has a slight

| | Resnet50 | Resnet152 | XceptionNet | MesoNet |
|---|---|---|---|---|
| Raw | 93.57% | 95.78% | 90.82% | 59.47% |
| LC | 90.69% | 92.11% | 83.28% | 56.94% |
| HC | 85.56% | 88.32% | 82.43% | 55.38% |

Table 5: The performances of models trained with mixed video quality dataset in different testing sets.

accuracy decrease, especially for heavy compression videos. This change supports our previous observation, which could have a large difference in distribution between raw and heavy compressed videos. The result also suggests that it might require a large model to learn the mixed video quality dataset. Resnet152, a deeper network, has smaller gaps with the models trained with a single-quality video set, which has 3% improvement for all qualities compare to the smaller Resnet50 network.

## 5.3. Evaluation of the Training Corpus Size

We evaluate how the training corpus size could affect the detection performance. We select 10000, 3000, 1000, 500, 100, and 10 animation videos and the corresponding number of real videos to train the Resnet50 model. The animation videos for each action are selected equally.

The chart on the left of Figure 3 presents the importance of training corpus size. It could only lead to a small downgrade in the raw dataset when keeping 10% of the data, but more impact when the videos are compressed. It could indicate that more low-quality videos are required to train the models. Besides, correction decreases dramatically if we use 1% of data for training in all quality. Simultaneously, the models tend to random guessing when we only have ten animation videos for training.

In addition, Figure 4 shows performances under different

proportions of animation and real data for training. More specifically, we reduce the number of one class's videos (either real or animation) to reach the expected proportions of videos for training. For example, "1:2 more real videos" means we remove half of the animation videos for training. Similar to decrease the total training corpus size, either reducing the real or animation videos will lead to an accuracy decrease. Also, preserving animation videos for training is more critical than introducing more real videos. When we use 10% of animation video for training, the performance only achieves 85%; a 10% decrease compares to the full dataset. On the contrary, keeping 10% of real videos with full animation videos leads to a smaller 5% downgrade.

The chart on the right of Figure 3 presents how the training corpus size could affect the performance of models trained with the mixed quality dataset. We trained the Resnet152 models with 3000, 1000, 500, 100, and 10 animation videos and the corresponding number of real videos under mixed video quality. Similar to the models trained with single quality datasets, a smaller corpus size also reduces corrections for the mixed quality situation. We also notice that the performance will slightly lower for large training corpus size than single-quality videos trained models. However, when the corpus size becomes smaller, the accuracy tends to be higher than the model trained with single-quality videos, especially for light compression videos, which suffering the most impact relate to decreased corpus size. One reason could be that the mixed video quality training strategy is similar to data augmentation, which increases the data diversity for small training corpus size to improve performance.
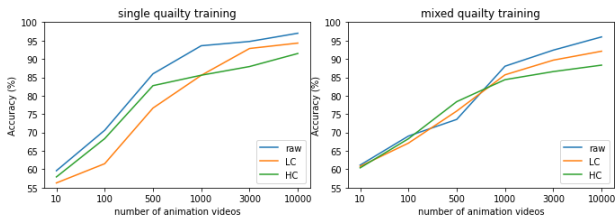


Figure 3: The performance changes under different training corpus sizes using single quality or mixed quality video dataset.

## 5.4. Evaluating the Impact of Type of Actions

The type of actions could affect the detection performance. We train the models with videos of one single action and evaluate whether the models could adapt to other unseen animation videos. For each action, we select all animation videos of that action and the corresponding number of real videos (1000 videos for each label) to train the Resnet50 models and test the performance with the whole animation testing video set. We select the raw quality videos for the experiment, and Figure 5 shows the result. With the training of single-action, it is not doubted that models cannot
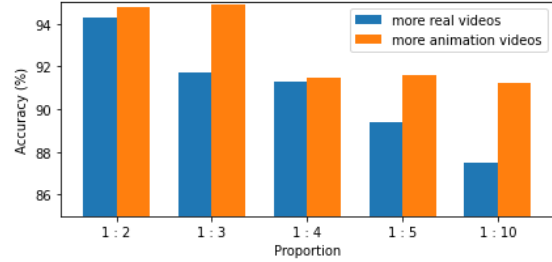


Figure 4: The detection accuracy of Resnet50 models trained with raw videos under different proportions of real and animation videos.

keep similar performance to the one trained with full actions raw quality videos with 1000 corpus size. The accuracy drops to 74.3% on average from 93.4%. The videos of nodding and surprising actions provide a better generalization to adapt unseen actions, which achieve 80.46% and 79.98% respectively. Right and left slope videos also introduce relatively high performance than remaining actions, which reach 77.02% and 75.89% respectively. On the other hand, the models trained with smile and blink videos have poor correction rates, which drop to 67.36% and 69.08%. In summary, using large movement action videos to train classifiers could lead to better performance for detecting new actions' videos. The actions that only include small changes, e.g., smiling and blinking, cannot provide sufficient information for the networks to adapt to unseen videos.

We also compare the full dataset trained models' performances for detecting different actions under three video qualities, and Figure 6 shows the result. We can notice that some actions' videos are relatively hard to detect in all video qualities. The left slope videos are the most difficult ones, which the accuracy is 92.3% under raw videos and drop to 88.24% under heavy compression videos. In addition, the detection of embarrassment videos might require higher video quality. Its performance decreases to 87.5%, the lowest one under heavy compression and a huge 7% gap compare to the one in raw quality. On the other hand, some actions, e.g., smiling, blinking, are much easier to be detected with the networks. They have achieved 100% correctness for raw videos, decrease to 97% on average if the videos are heavily compressed, which still keep in a higher level than other actions.

Comparing with Figure 5, we observe that some hard-to-detect actions, e.g., right and left slope, could provide more generalization if we only use those actions' videos for training. On the contrary, the models trained with the videos of easy-to-detect actions, e.g., smiling and blinking, showing poorer performances for adapting unseen actions.

## 5.5. Visualizing the Attention Parts of Models

We analyze what the classifier learned for distinguishing the animation and real videos. More specifically, we try
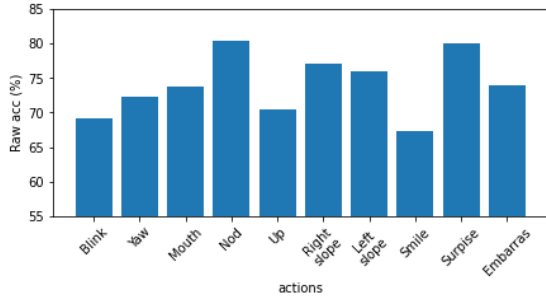
Figure 5: The detection accuracy for the Resnet50 models trained with one specific actions videos. Each label indicates the model of the whole testing data performance, trained with the videos of that specific action and 1,000 real videos.
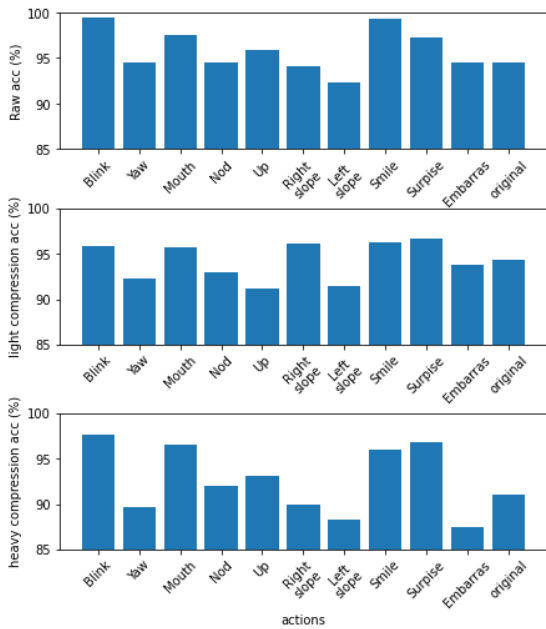


Figure 6: The detection accuracy for each action under different video quality. The models are trained with full training dataset. The original class indicate the selected real videos.

to visualize the network pay attention to which part of the video frames for detection. We use Class Activation Map (CAM) [41] to achieve visualization. The CAM is a technique to visualize what the classifier is looking at. CAM relies on the output from the models' global average pooling(GAP) layer, which right after the last convolution layer. The GAP layer could keep the spatial information of the convolution layer. By multiplying the weight of the Softmax layer of one specific class with GAP output, we can visualize the models' attention regions for classifying the given input images as that specific class. Our baseline models are binary classifiers to separate frames from either real and animation videos. In this case, the CAM results present semantic information about which parts of input image are critical regions for the models to decide whether it is from the animation videos.

In our experiment, we visualize the CAM results of the Resnet50 model trained with raw quality videos. To be noticed that, the GAP layer has been added to the model in the original design, so we don't require to do the further modification. The Figure 7 demonstrates some CAM results of selected video frames of different actions. We adapt the results to the original images to highlight the attention parts. The red color region indicates the important region, while the model pays less attention to the blue color areas. The results indicate that the model could learn semantic information to detect the animation videos. The model relies on the forgery regions to make the decision. For the opening mouth video frames, the model focuses on the detailed information of mouth. Similarly, the head and neck regions of the frames could be significant for detecting up videos. And the network pays attention to the profile face for head movement action's video frame, like yawing and sloping.



Figure 7: The Class Activation Map (CAM) results of the Resnet50 model trained with raw videos.

## 6. Discussion and Conclusion

We present a new large-scale, action-specified facial animation video dataset - Deepfake Mnist+, and evaluate the dataset's detecting performance with the proposed baseline detection method in different situations. We mainly explore the impact of the compression and actions on classification accuracy. It indicates the low-quality videos could significantly affect the performance and large movement actions could provide further generalization for unseen data. We expect that our proposed dataset could improve the detection performance of facial animation videos and increase the robustness and security of the recent liveness detectors as a response to the concern about pervasive animation videos online. As future work, we will explore other advanced facial animation methods and enlarge our datasets with more actions shot by more actors.

# References

[1] Avatarify. https://avatarify.ai/. (Accessed on 03/17/2021). 2

[2] Baidu ai. https://ai.baidu.com/tech/face/faceliveness. (Accessed on 03/17/2021). 4, 5

[3] Fakeapp 2.2.0 - download for pc free. https://www.malavida.com/en/soft/fakeapp/. (Accessed on 03/17/2021). 1

[4] Github - deepfakes/faceswap: Deepfakes software for all. https://github.com/deepfakes/faceswap. (Accessed on 03/17/2021). 1

[5] Github - dfaker/df: Larger resolution face masked, weirdly warped, deepfake,. https://github.com/dfaker/df. (Accessed on 03/17/2021). 1, 2, 3, 4, 5

[6] Github - iperov/deepfacelab: Deepfacelab is the leading software for creating deepfakes. https://github.com/iperov/DeepFaceLab. (Accessed on 03/17/2021). 1, 3

[7] Github - selimsef/dfdc_deepfake_challenge: A prize winning solution for dfdc challenge. https://github.com/selimsef/dfdc_deepfake_challenge. (Accessed on 03/17/2021). 1, 2

[8] Github - shaoanlu/faceswap-gan: A denoising autoencoder + adversarial losses and attention mechanisms for face swapping. https://github.com/shaoanlu/faceswap-GAN. (Accessed on 03/17/2021). 1

[9] Tianyandata. https://www.tianyandata.cn/. (Accessed on 03/17/2021). 4, 5

[10] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Mesonet: a compact facial video forgery detection network. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7. IEEE, 2018. 5

[11] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. Recycle-gan: Unsupervised video retargeting. In *Proceedings of the European conference on computer vision (ECCV)*, pages 119–135, 2018. 1

[12] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. Neural head reenactment with latent pose descriptors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13786–13795, 2020. 1

[13] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. 5

[14] Maria De Marsico, Michele Nappi, Daniel Riccio, and Jean-Luc Dugelay. Moving face spoofing detection via 3d projective invariants. In *2012 5th IAPR International Conference on Biometrics (ICB)*, pages 73–78. IEEE, 2012. 2

[15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5

[16] Brian Dolhansky, Russ Howes, Ben Pflaum, Nicole Baram, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) dataset. *arXiv preprint arXiv:2006.07397*, 1(2), 2020. 1, 3, 12

[17] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. 1

[18] David Güera and Edward J Delp. Deepfake video detection using recurrent neural networks. In *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6. IEEE, 2018. 2

[19] ST Hawk, J Van der Schalk, and AH Fischer. Moving faces, looking places: The amsterdam dynamic facial expressions set (adfes). In *12th european conference on facial expressions, geneva, switzerland*, volume 4, 2008. 4, 11

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 5

[21] Yinan He, Bei Gan, Siyu Chen, Yichun Zhou, Guojun Yin, Luchuan Song, Lu Sheng, Jing Shao, and Ziwei Liu. Forgerynet: A versatile benchmark for comprehensive forgery analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4360–4369, 2021. 3

[22] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 2

[23] Liming Jiang, Ren Li, Wayne Wu, Chen Qian, and Chen Change Loy. Deeperforensics-1.0: A large-scale dataset for real-world face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2889–2898, 2020. 12

[24] Jukka Komulainen, Abdenour Hadid, and Matti Pietikäinen. Face spoofing detection using dynamic texture. In *Asian Conference on Computer Vision*, pages 146–157. Springer, 2012. 2

[25] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. Advancing high fidelity identity swapping for forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5074–5083, 2020. 1, 2

[26] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5001–5010, 2020. 1, 2

[27] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3207–3216, 2020. 1, 2, 3, 12

[28] Arsha Nagrani, Joon Son Chung, Weidi Xie, and Andrew Zisserman. Voxceleb: Large-scale speaker verification in the wild. *Computer Speech & Language*, 60:101027, 2020. 4, 11

[29] Huy H Nguyen, Junichi Yamagishi, and Isao Echizen. Capsule-forensics: Using capsule networks to detect forged images and videos. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2307–2311. IEEE, 2019. 2

[30] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1–11, 2019. 1, 2, 3, 4, 5

[31] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *arXiv preprint arXiv:2003.00196*, 2020. 1, 2, 3, 11

[32] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 1, 2

[33] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 1, 2, 3, 4

[34] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. 1, 3, 4

[35] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8695–8704, 2020. 2

[36] Shun-Yi Wang, Shih-Hung Yang, Yon-Ping Chen, and Jyun-We Huang. Face liveness detection based on skin blood flow analysis. *symmetry*, 9(12):305, 2017. 2

[37] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*, 2018. 1

[38] Xin Yang, Yuezun Li, and Siwei Lyu. Exposing deep fakes using inconsistent head poses. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8261–8265. IEEE, 2019. 1, 2, 3

[39] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9459–9468, 2019. 1

[40] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503, 2016. 4, 11

[41] B. Zhou, A. Khosla, Lapedriza. A., A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. *CVPR*, 2016. 8

[42] Michael Zollhöfer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Pérez, Marc Stamminger, Matthias Nießner, and Christian Theobalt. State of the art on monocular 3d face reconstruction, tracking, and applications. In *Computer Graphics Forum*, volume 37, pages 523–550. Wiley Online Library, 2018. 2

## A. Data Generation And Collection Pipeline

The Figure 8 shows the pipeline to generate and collect our proposed DeepFake MNIST+ dataset. First, we collect the real videos from the VoxCeleb1 [28], extract the frames from these real videos as the source identity images. Then we shot driving videos with ten actions through the volunteers. In order to match the format of VoxCeleb1 videos, which are face-cropped and have the size of 256x256 resolution. We made a face-cropped version of driving videos with MTCNN modules [40] and also resize them into 256x256 resolution. We use Siarohin's framework [31] for animation video generation to produce face animation videos. It is a SOTA animation framework such that it even could capture the detail of eyeball moving. A single source image and a driving video were passed to the generator each time to produce single animation videos with a specific action. The generated videos were filtered with Liveness detector APIs to collect the challenging videos. Finally, 10,000 passed animation videos with ten actions (1000 videos for each action) and selected 10,000 real videos from VoxCeleb1 from our proposed DeepFake MNIST+ dataset.

## B. More Examples of Data

### B.1. Compressed Images

We demonstrate some raw and compressed (in both C23 and C40 compression rate under H.264 codec) video frames in Figure 9. The higher rate means heavier compression. As we can see, the c23 compression rate only leads to a minor impact on visual video quality. However, the frames suffer significant detail loss and blur effect under the c40 compression rate.

### B.2. Driving Video samples

In this section, we present some driving video samples of different actions for animating the source images in Figure 10. The driving videos of embarrassment are picked from ADFES dataset [19].
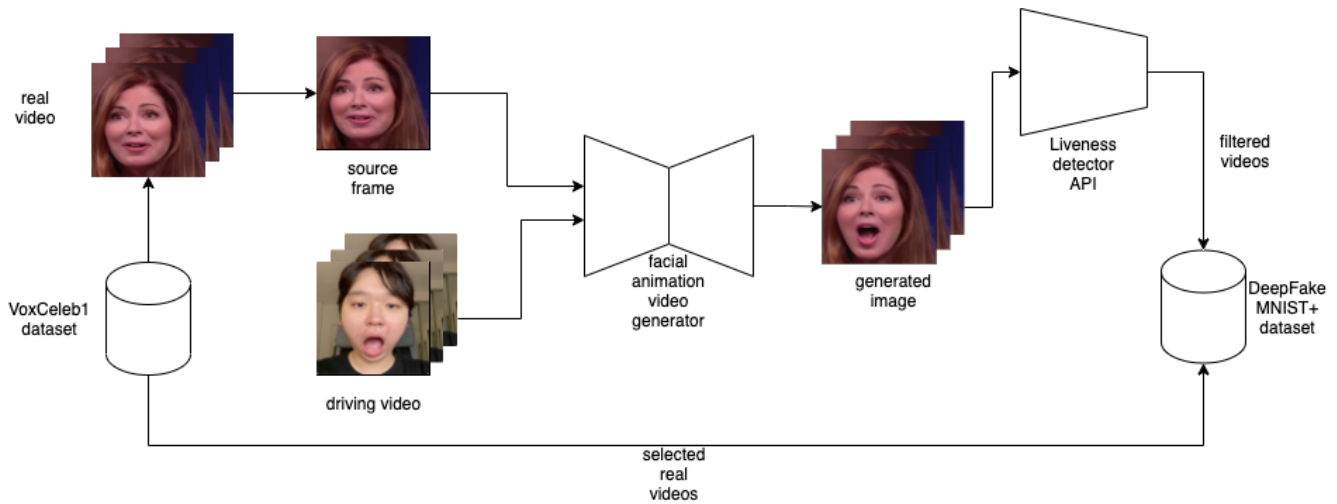


Figure 8: The pipeline to generate and collect our proposed DeepFake MNIST+ dataset.

## C. More experiments

|  | DFDC[16] | DF-1.0[23] | Celeb-DF[27] |
|---|---|---|---|
| without finetune | 61.2% | 60.7% | 57.2% |
| finetune | 95.6% | 96.1% | 95.1% |

Table 6: Accuracy of detecting DeepFake Mnist+ using the models trained with previous datasets. Fine-tuning means fine-tuned with proposed dataset.

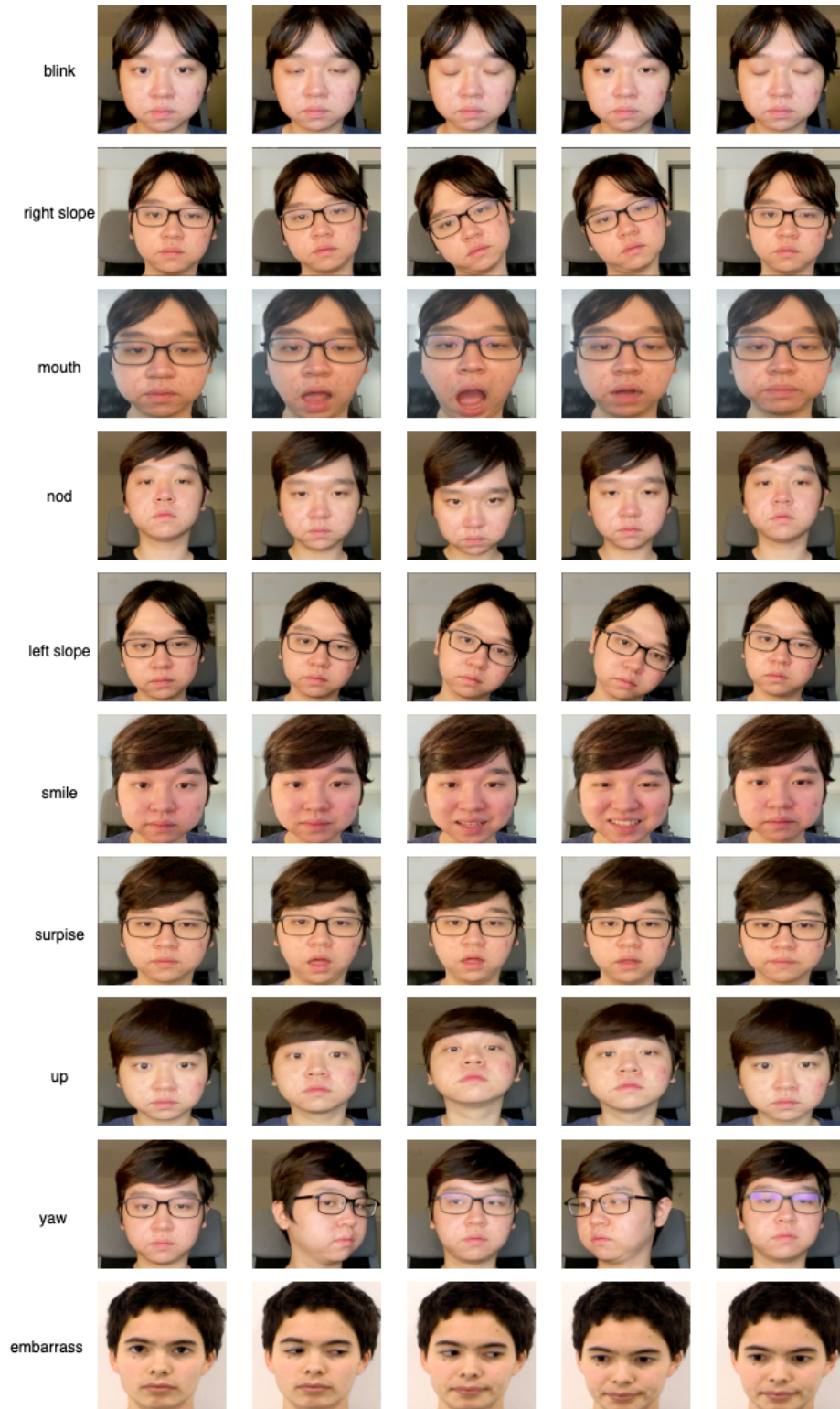Figure 9: The samples of raw and corresponding compressed (both c23 and c40) video frames.

Figure 10: The driving video samples for animating the source images.