

DEEPAKE CARICATURES: AMPLIFYING ATTENTION TO ARTIFACTS INCREASES DEEPAKE DETECTION BY HUMANS AND MACHINES

Anonymous authors

Paper under double-blind review

1 LINK TO VIEW CARICATURES IN VIDEO FORM

The effect of our caricatures is best experienced in video form. We provide a compilation of examples in the supplementary mp4 file *caricature_applied_on_fake_vids.mp4*. In the clip, we show the impact on the Caricature generation process on deepfakes. The artifacts are amplified and the video is generally distorted. We additionally provide an anonymized web gallery, *gallery.html*, that allows for easy visualization of several additional examples. Once the *gallery.html* file is extracted from the provided zip, along with its accompanying videos, double click on the html file to view it on your browser.

2 VIDEO SELECTION FOR HUMAN ARTIFACT ANNOTATIONS

2.1 SELECTING CHALLENGING DEEPAKES ACCORDING TO HUMANS.

Selecting target videos for our human artifact annotations poses a challenge: we would like videos to be hard to detect as fakes by humans, to ensure that our artifacts are a) non-trivial, and b) capture the subtleties of well-functioning face alteration techniques. To facilitate this selection, we formulated a crowd sourced experiment to find fake videos that are relatively difficult to detect as fake, without scrutiny.

We randomly selected 1000 video pairs (1000 real, 1000 fake) from the Deepfake Detection Challenge Dataset Dolhansky et al. (2019) and used FaceNet Schroff et al. (2015) to crop the face regions, resulting in videos of 360×360 pixels. These 1000 were randomly divided into sets of 100, and each set received 10 responses. Participants were instructed to maintain their fixation at a center cross while a pair of videos were presented, one on the left and one on the right (Figure 1). One of those videos was fake, while the other one was real. Participants were asked to select which video was more likely to be fake by pressing a left or right keyboard key. At the start of the experiment, people had the possibility to adjust the distance between the videos (left and right videos were always at the same distance from the center). To help with the adjustment, the following instructions were given “The videos should not be too close to disturb you from fixating at the center, and not too far away to notably reduce your ability to distinguish them. Find the distance that feels more comfortable to you.” After every 10 videos shown to users, users took a quick break and then had 3 seconds to refocus their attention on the fixation cross before examining the next set of videos. Participants were asked to report whenever their gaze shifted towards either side of the videos, in which case the trial was discarded.

We computed video difficulty by averaging the responses over videos. Based on those responses, we selected 500 fake videos that were closest to having a 50% recognition rate, which indicated that human detection performance was close to chance. Those videos were considered “hard to detect”.

2.2 SELECTING CHALLENGING FAKES ACCORDING TO MODELS.

We would additionally like to represent the distribution of hard fakes for models, and collect human artifact maps on these fakes to ensure our distribution of maps covers the space of challenging artifacts. For this, we selected 500 fake videos from our DFDCp training split that minimized the “fake” classification confidence (output of the softmax) of an Xception network trained on FF++, while not being a duplicate of the 500 videos selected through our human experiment.

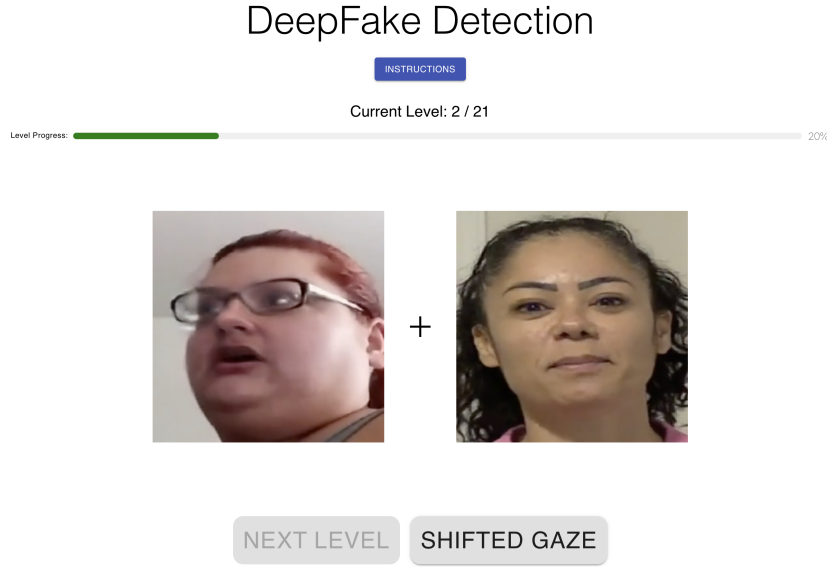


Figure 1: Screenshot of the interface used to collect human detection difficulty. Humans were tasked with indicating which of the two videos was fake without shifting their gaze from the cross.

3 VIDEO ARTIFACT ANNOTATIONS

3.1 METHODS

Videos were annotated in an online task that used a paintbrush interface. Participants were instructed to paint over the regions that were most informative about whether a video was fake. Videos were paused during painting, and the interface contained a separate visualization of where the participant has painted, so they could keep track of their annotations.

3.2 QUALITATIVE RESULTS

In total we collected over 11K annotations for 1000 videos. On average, each video had annotations from 22.6 participants, resulting in 4.1K sample points per video. For each video clip, we aggregated all annotation data and generated one 3D attention map: An anisotropic Gaussian kernel of size $(c \cdot s, c \cdot s, 6)$ in x,y and time dimensions was first applied. c is set to a constant 20px and s corresponds to the brush size. We normalize the attention map of each time frame to sum to one.

Humans tend to annotate areas of the face that are near the eyes or mouth, which are more salient facial features. Our model picks up on these features as well as other areas that are indicative of doctoring. For example, there is attention to inflection regions of the face, such as the jawline or near the hairline. Humans also gravitate towards inconsistencies around the seams generated by the face swapping pipeline, as well as flickering artifacts generated by tracking errors during the deepfake process.

4 MODEL LOSS

We utilize a dual loss to train our prediction and caricature generation network. The fake / real output of our classification module is supervised with a typical binary cross-entropy loss L_{BCE} :

$$L_{BCE} = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(\bar{y}_i) + (1 - y_i) \cdot \log(1 - \bar{y}_i)$$

Where N is the number of elements in our dataset.

Additionally, we impose another loss on the system: a correlation loss L_{human} ensures that the output of the artifact attention module reproduces the ground truth human maps. This loss is composed of a KL-divergence component, a Pearson Correlation Coefficient component, and an L1 loss component, all weighted by specific factors determined experimentally. The loss takes the following form:

$$L_{human} = \alpha L_{KL} + \beta L_{CC} + \gamma L^1$$

Let p and q be the predict and ground truth maps respectively, M the total number of pixels in the map, and $\epsilon = 1e - 7$ to avoid numerical issues. We have:

$$L_{KL}(p, q) = \sum_i q_i \log \left(\epsilon + \frac{q_i}{\epsilon + p_i} \right)$$

$$L_{CC}(p, q) = \frac{\sigma(p, q)}{\sigma(p) \cdot \sigma(q)}$$

$$L^1(p, q) = \frac{1}{M} \sum_i |p_i - q_i|$$

These losses are computed element-by-element and averaged over a minibatch during training. We perform a simple grid search to determine adequate loss parameters and converge to the following values: $\alpha = 1$, $\beta = -0.5$, $\gamma = 1$. Note that β is negative to ensure the correlation coefficient is maximized.

The total loss L of our model is $L = L_{BCE} + L_{human}$.

5 DETAILS OF BENCHMARKS AND BASELINES FOR MODEL EVALUATION

5.1 DATASETS

We evaluated models on four benchmarks:

FaceForensics++ (FF++) Rossler et al. (2019), which contains more than 1.8M images generated from 5000 YouTube videos. It features four different facial manipulation methods, which are Deepfake D (2020), Face2Face Thies et al. (2016), FaceSwap Kowalski (2020) and NeuralTextures Thies et al. (2019), and covers a large variation in video resolution ranging from 1080p to 480p. We use 8151 images for training, 280 images for validation, and 280 for testing. To ensure fair comparisons with Haliassos et al. (2021) and others, we use the first 270 frames of each training video and the first 110 frames for each validation video.

The Deepfake Detection Challenge Dataset (DFDCp) Dolhansky et al. (2019). We use the preview version of this dataset, which contains 3754 fakes and 1460 real videos. Two facial modification methods were used to generate the fake videos, however, no specific labels were given regarding the generation method. In this work, we report performance on the publicly available validation set of 400 labeled videos.

Celeb-DF v2 Li et al. (2019c) dataset provides 5639 synthesized fake videos of celebrities generated from 590 real videos. Based on a tandem auto-encoder architecture, a video synthesis method that puts special emphasis on reducing visual artifacts was used to generate high quality fake videos. In our experiment, we split the dataset into a training set of 6011 videos, and a validation set of 518 videos.

DeeperForensics (DFo) Li et al. (2019c), which contains 10000 fake videos obtained by taking 1000 real videos from FF++, and swapping in faces of 100 new paid actors.

FaceShifter (FShifter) Li et al. (2019a), which also consists of 10000 fake videos generated by doctoring the real videos from FF++.

Video Standardization Videos from the different datasets under study were preprocessed to bring them into a standard configuration. All datasets were downloaded at full resolution, except FaceForensics++, which was downloaded in the provided *c23* compressed format to avoid large storage costs. We performed multi-person face extraction on all videos using a pretrained PyTorch implementation of FaceNetSchroff et al. (2015), with a minimum face size of 50 pixels and a face margin size of 100 pixels. For frames without a detected face, bounding box coordinates were produced via linear interpolation of the next available neighboring bounding boxes. To avoid abrupt motion artifacts (e.g. jittering) from face extraction, the 4 corner coordinates of the detected bounding boxes were temporally smoothed using hanning smoothing with a window size equal to 20% the video length. These face frames were resized to 360×360 pixels and saved as an MP4 file.

5.2 BASELINE MODELS

Clip-based ResNet: Building on features learned from ImageNet, the inflation of pretrained 2D kernels into 3D has shown significant improvement on video-based tasks Carreira & Zisserman (2017); Hara et al. (2018); Tran et al. (2018). As a baseline of 3D CNN, We use a 3D ResNet model where the 3D kernels are inflated from ResNet18 pretrained on the Kinetics dataset Carreira et al. (2019).

XceptionNet, which was trained on the FaceForensics++ dataset.

Face X-ray, which attends to the seams in face-swapping to make its predictions. We use existing performance values from Haliassos et al. (2021) for the FF++ trained results, and evaluate an HRNet-W48 trained on synthetic blended images in conjunction with FF++ videos for our DFDCp results.

The *CNN-GRU* from Sabir et al., where we train a DenseNet-161 followed by a GRU on FF++ to evaluate on DFDCp, and report existing performance values for FF++.

The *Multi-task* model from Nguyen et al. (2019). We report existing numbers from Haliassos et al. (2021) and retrain using the publicly available implementation.

DSP-FWA, a model that tries to detect face warping artifacts. We evaluate the publicly available pretrained model for FF++, and train a ResNet50 with the pyramidal setup introduced in Li & Lyu (2018) for DFDCp.

The *Two-branch* model from Masi et al. (2020). We report existing numbers for CelebDFv2 evaluation, and retrain on the rest.

The *Multi-attention* model from Zhao et al. (2021). We report the values from the existing paper, and retrain following the setup in the paper on FF++ for DFDCp results.

The Lip Forensics approach from Haliassos et al. (2021). We use their publicly available values for our main tables.

6 USER STUDY: EXTENDED METHODS

6.1 STUDY 1: HOW DO CARICATURES COMPARE TO TEXT-BASED INDICATORS

Stimuli. Stimuli consisted of 50 real and 50 fake videos selected from DFDC. The 50 fake videos were high-quality deepfakes, which were difficult for human participants to detect. These were selected by running a separate pilot deepfake detection task on a larger subset of the DFDC, and selecting 50 videos with detection rates between 30% and 10%, to have a range of difficulty. We focused on difficult deepfakes because this provides the best test-case of how convincing a visual indicator is; users would not need to rely on a visual indicator if they can see the artifacts for themselves.

Behavioral Task. The behavioral task was a simple 2-alternative forced choice task. Users were shown one video at a time, and were asked to assess whether the video was real or fake. There were three condition, presented in a between-subjects manner. In the unaided detection condition, real and fake videos were presented as-is, with no further modifications. In the text-based indicator condition, a badge was added to the top-left corner of all fake videos which said "video modified by AI" in bold white letters on a semi-opaque bright turquoise background. In the caricature condition, the Deepfake Caricature transformation was applied to all fake videos prior to presentation.

Participants. Participants consisted of 90 workers on Prolific.com. Participation was restricted to individuals over the age of eighteen, located in the US, with approval ratings over 90%. Participants were recruited and compensated according to procedure approved by the [University's] Committee on the Use of Humans as Experimental Subjects. Participants were paid according to an hourly rate of \$11.25 per hour.

Analysis. Hit rates were calculated for each subject, and averaged across subjects for each condition. Hit rates were compared across conditions using planned t-tests.

6.2 STUDY 2: WHAT RANGE OF CONDITIONS ARE CARICATURES EFFECTIVE IN

Stimuli. Stimuli consisted of 400 videos selected from DFDCp Dolhansky et al. (2019). Half of the videos contained real, unaltered faces, while the other half contained deepfake faces. Videos were randomly selected, then deepfakes with obvious artifacts were discarded. All videos were subjected to the standardization procedure described above, and were presented at a resolution of 300x300 pixels. These 200 fake videos made up the stimuli for the "standard deepfake" condition. These videos were then passed through CariNet to receive modulation from the Caricature Creation Module, yielding stimuli for the "Caricature" condition.

Behavioral Task. Participants performed a two-alternative forced choice detection task. They were shown one video at a time and were asked to indicate if the video was real or fake by using the mouse to click on buttons labeled "Real" and "Fake". Prior to the tasks, participants were shown instructions which included 3 example deepfakes, to inform them about what a fake video looks like.

Each task session contained 100 videos, divided into 5 blocks of 20 videos with breaks in between. Videos could be presented at 6 different presentation durations: 300ms, 500ms, 1000ms, 3000ms and 5000ms. Each trial started with a 3 second countdown. In addition, each block contained one compliance check, and one engagement probe. Compliance checks consisted of real or fake videos (not included in the experimental set) with the text "THIS IS AN ATTENTION CHECK, PLEASE CLICK REAL" written across the video in light grey font. These trials were included to provide an objective basis by which to exclude low-performing participants (see Analysis section below). Engagement probes consisted of deepfakes (not included in the experimental set) containing extremely obvious artifacts, such as large patches of noise covering the face. These trials were included to provide an on-line measure of participants attentiveness. We reasoned that while highly-engaged participants would succeed on all of these trials, medium to low-engagement participants would miss some proportion of them.

Trials were divided into experimental sessions (or "HITS") of 100 trials. Trials were assigned to HITS pseudo randomly, with the exception that no fake or caricature could be shown twice in the same HIT. Data were collected for 19 HITS of randomly selected videos, leading to an average of

72 videos per presentation time for both standard fakes and caricatures. We collected 10 participants per each HIT.

Participants. Participants consisted of 170 workers on Amazon Mechanical Turk. Participation was restricted to individuals over the age of eighteen, located in the US, with approval ratings over 99%. Participants were recruited and compensated according to procedure approved by the [University’s] Committee on the Use of Humans as Experimental Subjects. Participants were paid according to an hourly rate of \$11.25 per hour.

Analysis. We collected 10 judgments per video per condition. This number was selected after a power analysis (performed on a pilot set of 100 HITs), which suggested that 6 responses was sufficient to obtain stable estimates of the mean detectability of a given video.

Our primary analysis was to examine the difference in correct detection rates between the standard deepfake condition and the caricature condition, both pooling over presentation times and for each presentation time individually. As such, we pre-registered an ANOVA analysis to test for a main effect of condition on accuracy, with the hypothesis that correct detection would be higher in the caricature than standard deepfake condition. We additionally hypothesized that overall detection (pooling across conditions) would be higher for longer presentation times (main effect of presentation time), and that the difference between standard and caricature deepfakes would increase with higher presentation time (significant interaction), although these analyses were orthogonal to the main question of this experiment. We also planned further comparisons between the standard and caricature conditions for each time point individually. We pre-registered individual unpaired t-tests at a conservative alpha level of $p = 0.0083$ (Bonferroni correction for performing 6 individual t-tests, one per time point). To ensure quality data for this analysis, we excluded any sessions where participants used only one key, or failed three or more compliance checks. These trimming and analysis procedures were developed in an initial pilot dataset of 100 HITs, then applied in the experimental dataset (i.e. they were pre-registered). We elected to use the conservative Bonferroni correction even though these were a-priori analyses because this experiment had high enough power to detect effects even at conservative thresholds, and applying this threshold would reduce the chance of Type I error.

We also performed a second analysis to examine the effects of condition across levels of participant attentiveness. We divided HITs into participant engagement bins based on the number of engagement probes that were successfully detected throughout the HIT. As the dependent measure, we subtracted the correct detection rate in the standard condition from the caricature condition, to get a measure of the relative benefit of the caricature condition, and tested whether this difference was greater than 0 using a one-sided one-sample t-test at an alpha level of $p = 0.05$. We used a traditional rather than conservative p values because some conditions were likely to contain a low number of trials, and thus have lower power. Statistics were only carried out for conditions containing 10 or more datapoints. This analysis used all experimental sessions.

6.3 COMPLETE STATISTICAL REPORTING

Analysis 1: ANOVA testing for an effect of condition

	SS	df	F	p
Condition	2.17	4.0	16.46	3.13-13
Presentation time	2.43	5.0	14.75	3.50e-14
Interaction	6.15	20.0	9.32	3.98e-20
Residual	54.79	1661.0	-	-

Analysis 2: planned comparisons of detectability difference between standard fake and caricatures for each timepoint.

	300ms	500ms	1000ms	3000ms	5000ms
df	155	160	159	129	172
t	4.58	7.97	9.11	11.23	12.86
p	9.2e-6	2.9e-13	3.1e-17	7.7e-21	5.9e-27

	0	0.2	0.4	0.6	0.8	1
df	-	9	18	20	58	76
t	-	2.63	1.73	6.27	9.82	15.64
p	-	0.015	0.051	2.5e-06	3.6e-14	1.3e-25

Analysis 3: Significance of the difference between standard fake and caricatures for each level of participant engagement.

7 RESULTS INCLUDING THE DFDCP-TRAINED ARTIFACT ATTENTION MODULE

In the paper, we reported the results from a model with a FF++-trained Artifact Attention Module to ensure fair comparisons with prior work. Our DFDCp-trained model had similar performance, but did not perform with the same consistency as the FF++ trained model. We report the values here for completeness.

Model	CelebDFv2	DFDCp	FShifter	DFo	Overall
Xception Rossler et al. (2019)	73.7	65.7	72.0	84.5	75.3
Face X-ray Li et al. (2019b)	79.5	62.1	92.8	86.8	81.2
CNN-GRU Sabir et al.	69.8	63.7	80.8	74.1	73.4
Multi-task Nguyen et al. (2019)	75.7	63.9	66.0	77.7	71.9
DSP-FWA Li & Lyu (2018)	69.5	64.5	65.5	50.2	63.1
Two-branch Masi et al. (2020)	73.4	64.0	-	-	-
Multi-attention Zhao et al. (2021)	67.4	67.1	-	-	-
LipForensics Zhao et al. (2021)	82.4	70.0	97.1	97.6	86.8
FTCN Zheng et al. (2021)	86.9	74.0	98.8	98.8	-
DCL Sun et al. (2022)	82.3	<u>76.7</u>	92.4	97.1	-
RF Haliassos et al. (2022)	86.9	75.9	99.7	<u>99.3</u>	-
S-B Shiohara & Yamasaki (2022)	93.2	72.4	-	-	-
X+PCC Hua et al. (2023)	54.9	62.7	-	-	-
CariNet DFDCp (ours)	<u>87.6</u>	77.1	<u>98.9</u>	99.3	-

Table 1: **Detection performance results on unseen datasets.**

Method	Train on remaining				Avg
	DF	FS	F2F	NT	
Xception Rossler et al. (2019)	93.9	51.2	86.8	79.7	77.9
CNN-GRU Sabir et al.	97.6	47.6	85.8	86.6	79.4
Face X-ray Li et al. (2019b)	99.5	93.2	94.5	92.5	94.9
LipForensics Haliassos et al. (2021)	99.7	90.1	99.7	99.1	99.5
CariNet DFDCp (ours)	99.9	99.9	99.8	99.4	99.8

Table 2: **Generalization to unseen manipulations.**

Method	Clean	Cont.	Noise	Blur	Pixel
Xception	99.8	98.6	53.8	60.2	74.2
CNN-GRU	99.9	98.8	47.9	71.5	86.5
Face X-ray	99.8	88.5	49.8	63.8	88.6
LipForensics	99.9	99.6	73.8	96.1	95.6
CariNet DFDCp (ours)	99.9	99.9	78.3	93.3	97.5

Table 3: **Generalization performance over unseen perturbations.**

REFERENCES

- Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6299–6308, 2017.
- Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-700 human action dataset. *arXiv preprint arXiv:1907.06987*, 2019.
- D. Deepfake. <https://github.com/deepfakes/faceswap>, 2020.
- Brian Dolhansky, Russ Howes, Ben Pfau, Nicole Baram, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) preview dataset. *arXiv preprint arXiv:1910.08854*, 2019.
- Alexandros Haliassos, Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Lips don’t lie: A generalisable and robust approach to face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5039–5049, 2021.
- Alexandros Haliassos, Rodrigo Mira, Stavros Petridis, and Maja Pantic. Leveraging real talking faces via self-supervision for robust forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14950–14962, 2022.
- Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- Yingying Hua, Ruixin Shi, Pengju Wang, and Shiming Ge. Learning patch-channel correspondence for interpretable face forgery detection. *IEEE Transactions on Image Processing*, 32:1668–1680, 2023.
- Marek Kowalski. Faceswap. <https://github.com/MarekKowalski/FaceSwap/>, 2020.
- Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457*, 2019a.
- Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. *arXiv preprint arXiv:1912.13458*, 2019b.
- Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts, 2018.
- Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A new dataset for deepfake forensics. *arXiv preprint arXiv:1909.12962*, 2019c.
- Iacopo Masi, Aditya Killekar, Royston Marian Mascarenhas, Shenoy Pratik Gurudatt, and Wael AbdAlmageed. Two-branch recurrent network for isolating deepfakes in videos. In *European Conference on Computer Vision*, pp. 667–684. Springer, 2020.
- Huy H Nguyen, Fuming Fang, Junichi Yamagishi, and Isao Echizen. Multi-task learning for detecting and segmenting manipulated facial images and videos. *arXiv preprint arXiv:1906.06876*, 2019.
- 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 International Conference on Computer Vision*, pp. 1–11, 2019.

- Ekraam Sabir, Jiabin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. Recurrent convolutional strategies for face manipulation detection in videos. *Interfaces (GUI)*, 3: 1.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- Kaede Shiohara and Toshihiko Yamasaki. Detecting deepfakes with self-blended images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18720–18729, 2022.
- Ke Sun, Taiping Yao, Shen Chen, Shouhong Ding, Jilin Li, and Rongrong Ji. Dual contrastive learning for general face forgery detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 2316–2324, 2022.
- 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*, pp. 2387–2395, 2016.
- 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.
- Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 6450–6459, 2018.
- Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2185–2194, 2021.
- Yinglin Zheng, Jianmin Bao, Dong Chen, Ming Zeng, and Fang Wen. Exploring temporal coherence for more general video face forgery detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15044–15054, 2021.