

Visual Representations of Physiological Signals for Fake Video Detection

Kalin Stefanov
Monash University
Australia
kalin.stefanov@monash.edu

Bhawna Paliwal
IIT Ropar
India
2017chb1039@iitrpr.ac.in

Abhinav Dhall
IIT Ropar
India
abhinav@iitrpr.ac.in

Abstract—Realistic fake videos are a potential tool for spreading harmful misinformation given our increasing online presence and information intake. This paper presents a multimodal learning-based method for detection of real and fake videos. The method combines information from three modalities – audio, video, and physiology. We investigate two strategies for combining the video and physiology modalities, either by augmenting the video with information from the physiology or by novel learning the fusion of those two modalities with a proposed Graph Convolutional Network architecture. Both strategies for combining the two modalities rely on a novel method for generation of visual representations of physiological signals. The detection of real and fake videos is then based on the dissimilarity between the audio and modified video modalities. The proposed method is evaluated on two benchmark datasets and the results show significant increase in detection performance compared to previous methods.

I. INTRODUCTION

Advances in computer vision and deep learning have enabled the creation of very realistic fake versions of videos, known as *deepfakes*¹²³⁴. Highly realistic deepfakes are a potential tool for spreading harmful misinformation given our increasing online presence and information intake. Hence, it has become increasingly important to identify deepfakes with more accurate and reliable methods. The recent surge in synthesized fake video content on the Internet has also led to the release of several benchmark datasets (*e.g.*, [1]–[3]) and methods (*e.g.*, [3]–[5]) for fake content detection. The aim of these fake video detection methods is to correctly classify any given input video as either real or fake.

There is rich and growing literature on different methods for fake video detection. Previous methods can be loosely placed in two groups – methods that use a single modality (unimodal) and methods that use more than one modality (multimodal) for the detection of fake videos. Unimodal methods usually exploit artifacts in the visual modality, whereas, multimodal methods combine multiple modalities (*e.g.*, audio and video) that can provide complementary information and lead to better fake video detection rates.

The main idea of this work is to validate the hypothesis that *physiological signals* information from face videos can be used

to better classify fake and real videos, and to answer the related questions such as how can we *generate* suitable representations of these physiological signals, and how can we *incorporate* these signals in existing methods to improve the performance even further. The idea that fake videos exhibit considerably more inconsistencies in the associated physiological signals compared to real videos is thoroughly investigated in this work through evaluation of a video augmentation approach with novel visual representations of physiological signals and a proposed learning-based method. The main contributions of this work include:

- A method for generation of visual representations of physiological signals (*i.e.*, physiological maps) which are helpful in better determining if a video is real or fake.
- A data augmentation strategy that employs the physiological maps as facial “heat maps”.
- A Graph Convolutional Network for learning-based fusion of the physiological maps and face videos instead of direct augmentation.
- The developed multimodal methods are evaluated on two benchmark deepfake detection datasets and the results show that the proposed physiological maps increase the performance of fake video detection using both strategies; with Graph Convolutional Network based learning giving improvements over the augmentation strategy.

II. BACKGROUND

The topic of deepfake detection is well-researched and growing with increasing use cases [6], [7]. This section reviews some of the detection methods that specifically deal with deepfake content of human faces. Deepfake generation methods often generate artifacts in the resulting video that are too subtle for humans but can be easily detected using machine learning. According to [7] there are two groups of visual artifacts: 1) spatial artifacts and 2) temporal artifacts.

Blending is a spatial artifact related to the areas where the generated content is blended back into the original content. Researchers have proposed different strategies to emphasize these artifacts, for example, edge detectors (*e.g.*, [8]) and frequency analysis (*e.g.*, [9]). Environment is another spatial artifact related to the content of the fake face being anomalous in the context of the rest of the image. Residuals from the face

¹<https://github.com/deepfakes/faceswap>

²<https://github.com/dfaker/df>

³<https://github.com/iperov/DeepFaceLab>

⁴<https://github.com/shaoanlu/faceswap-GAN>

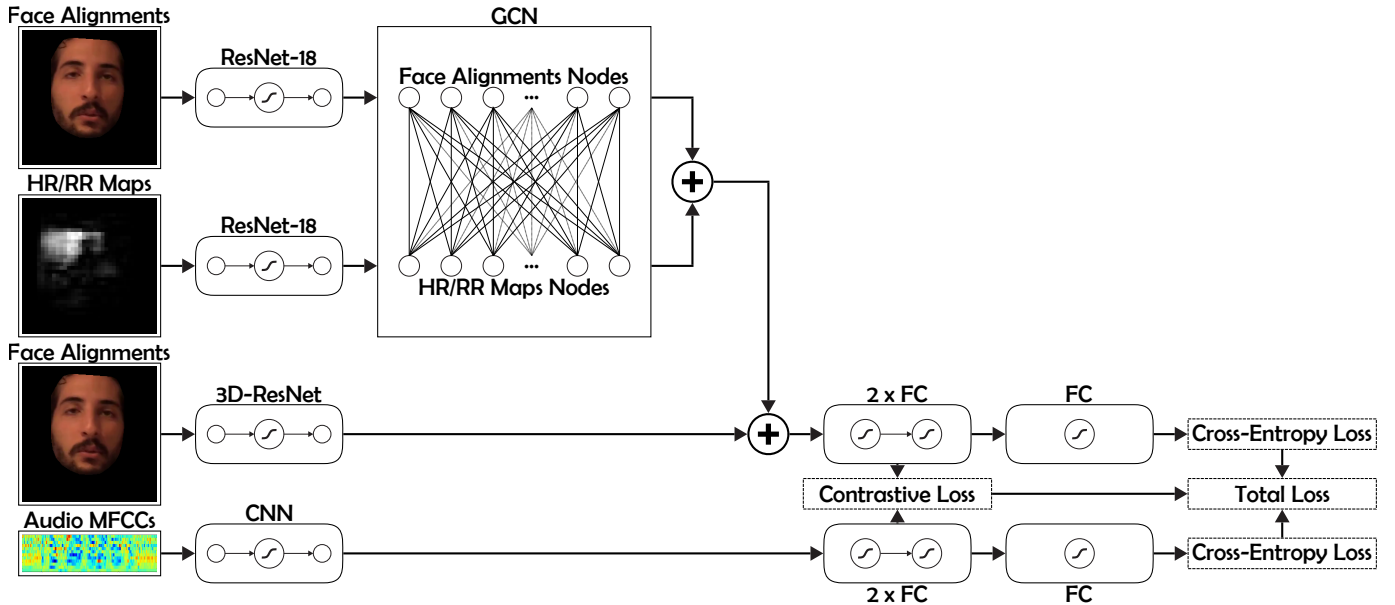


Fig. 1: Overview of the proposed method for physiology-based fake video detection. Graph Convolution Network is used to learn the fusion of the physiological maps and video streams. Then a contrastive loss based learning is performed incorporating information from the audio stream. Symbol \oplus indicates concatenation of vectors.

warping processes (e.g., [10]) and lighting (e.g., [11]) can indicate the presence of generated content in the image.

Behavior can exhibit temporal artifacts in fake videos. Given sufficient amount of data for a target person, different behaviors can be monitored for anomalies (e.g., [12]). Another temporal artifacts are related to physiology, for example, blood volume patterns under the skin [13]. Synchronization inconsistencies (between audio and video) are another temporal artifact usually found in fake videos (e.g., [14]). Realistic temporal coherence is challenging to generate, and some previous work focus on the resulting artifacts to detect the fake content (e.g., [15]).

Another approach to fake video detection is to train a generic classifier and let the model choose the features to analyze instead of focusing on the specific artifacts generated by deepfake methods. Deep neural networks have been shown to effectively detect deepfake videos when employed as classifiers (e.g., [16]–[18]). In contrast to classification, anomaly detection methods are trained on real data and then detect outliers during testing (e.g., [19], [20]).

The work described in this paper falls under the group of temporal artifacts. More specifically, the hypothesis as earlier mentioned is that fake videos exhibit considerably more temporal inconsistencies in the associated physiological signals compared to real videos. The literature on fake video detection using physiological signals is quite limited compared to other type of spatial and temporal artifacts. The closest related work is the one in [13]. The differences between this work and [13] include:

- This work does not consider specific face regions for extraction of the physiological signals, hence it is not

limited to only portrait videos as in [13].

- This work proposes method for generation of physiological maps that can be directly applied as “heat maps” on the face videos, producing novel data augmentations.
- This work proposes a multimodal method whereas the method in [13] is vision-only.

III. PHYSIOLOGY-BASED FAKE VIDEO DETECTION

The proposed fake video detection method is based on the hypothesis that fake videos exhibit considerably more inconsistencies in the associated camera-based physiological measurements compared to real videos. This section offers a description of the proposed method for physiology-based fake video detection.

A. Method Overview

Given an input video, the goal is to classify it as real or fake, that is, given a training dataset $D^{train} = \{(v^1, y^1), (v^2, y^2), \dots, (v^N, y^N)\}$ consisting of N videos, where v^i denotes the input video and the label $y^i \in \{0, 1\}$ indicates whether the video is real ($y^i = 0$) or fake ($y^i = 1$), the goal is to create a computational model that can correctly classify a new video as either real or fake.

The process of data preparation and classifier training is similar to the work described in [22]. The audio signal is extracted from the input video v^i using the FFmpeg⁵ library, and then split into D -second long segments. Similarly, the video signal is divided into D -second long segments, and face tracking is performed on those video segments using the S3FD [23] face detector to extract face crops. The data

⁵<http://ffmpeg.org>

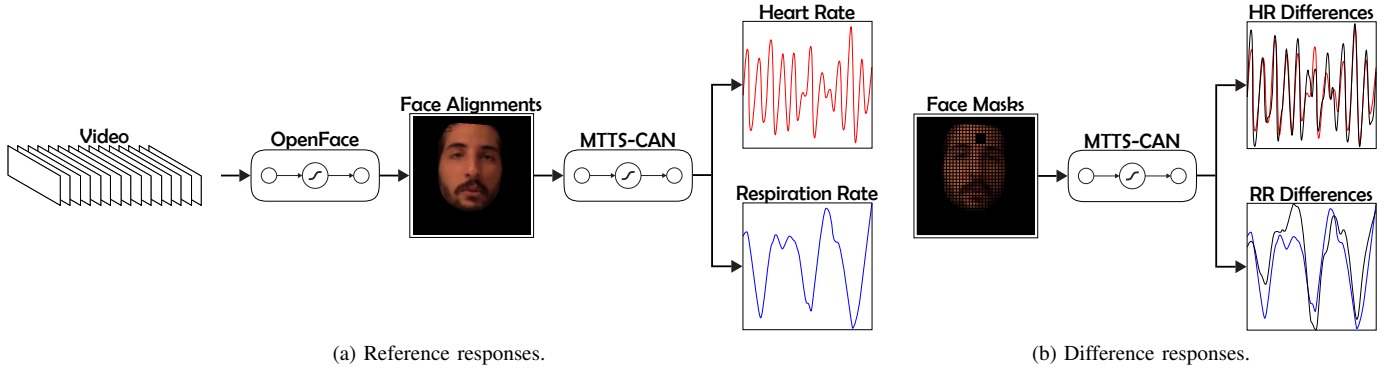


Fig. 2: Overview of the proposed method for generation of physiological maps. The aligned face crops are passed through the camera-based physiological measurement model (*i.e.*, MTTs-CAN [21]) to obtain the reference responses. Then an occlusion patch is slid over the aligned face crops, new responses are calculated and the differences with the reference responses are used as encoding for the relative importance of the image area under the occlusion patch.

preprocessing produces image segments $\{s_1^i, s_2^i, \dots, s_n^i\}$ and the corresponding audio segments $\{a_1^i, a_2^i, \dots, a_n^i\}$, where n denotes segment count for an input video v^i . The classifier consists of a bi-stream network, where each image segment s_t^i ($t = 1 \dots n$) is passed through a video stream S_v , and the corresponding audio segment a_t^i is passed through an audio stream S_a . The network is trained using a combination of contrastive loss (L_1) and cross-entropy loss (L_2 and L_3). The contrastive loss is meant to ensure that the video and audio streams are closer for real videos, and farther for fake videos. The cross-entropy loss is meant to ensure learning of robust video and audio feature representations. The overall loss L is a weighted sum of the three losses, L_1 , L_2 and L_3 [22]:

$$L_1 = \frac{1}{N} \sum_{i=1}^N (y^i)(d_t^i)^2 + (1 - y^i) \max(\text{margin} - d_t^i, 0)^2 \quad (1)$$

$$d_t^i = \|f_v - f_a\|_2 \quad (2)$$

$$L_2 = -\frac{1}{N} \sum_{i=1}^N y^i \log \hat{y}_v^i + (1 - y^i) \log(1 - \hat{y}_v^i) \quad (3)$$

$$L_3 = -\frac{1}{N} \sum_{i=1}^N y^i \log \hat{y}_a^i + (1 - y^i) \log(1 - \hat{y}_a^i) \quad (4)$$

$$L = \lambda_1 L_1 + \lambda_2 L_2 + \lambda_3 L_3 \quad (5)$$

where λ_1 , λ_2 and λ_3 all equal to 1, margin is a hyper-parameter and f_v and f_a are feature representations of the video and audio streams respectively.

This work extends [22] with the introduction of novel visual representations of physiological signals into the network. The proposed visual representations are either used to augment the original face crops or the relationship between the face crops and the proposed visual representations is learned from data

through a novel Graph Convolutional Network based architecture. The next subsections detail this further. An overview of the proposed method is illustrated in Figure 1.

B. Generation of Physiological Maps

For the generation of visual representations of physiological signals we utilize recent advancements in camera-based physiological measurement. In particular, we leverage a novel Multi-Task Temporal Shift Convolutional Attention Network (MTTS-CAN) [21] that enables real-time cardiovascular (*i.e.*, heart rate) and respiratory measurements (*i.e.*, respiration rate), that is, given an RGB video of a frontal face (*i.e.*, face crop), the network estimates the waveform of the two physiological signals.

One of the contributions of this work is a method for generation of visual representations (*i.e.*, physiological maps) of physiological signals based on the estimated signal waveform. Visual representations of the signal waveform is important since a single float number estimate for the physiological signal for each face crop is not useful data representation. Therefore, we propose a novel method for generation of visual representations that can be used to learn relationships between the video and physiological signals and provide a significant improvement in the downstream task of fake video detection. The process of physiological maps generation is divided into several steps: 1) the face is detected and aligned with OpenFace [24] and most of the background is removed, 2) the aligned face crops are passed through MTTs-CAN and the two physiological signals (*i.e.*, heart rate and respiration rate) are estimated and kept as reference, 3) a square occlusion patch is defined and used to mask-out part of all aligned face crops and the occluded face crops are passed through MTTs-CAN to estimate the two physiological signals, 4) the difference between the estimated reference signals without occlusion and the estimated signals with occlusion is considered as the relative contribution (importance) of the masked-out region for the accurate estimation of the two physiological signals, and 5)

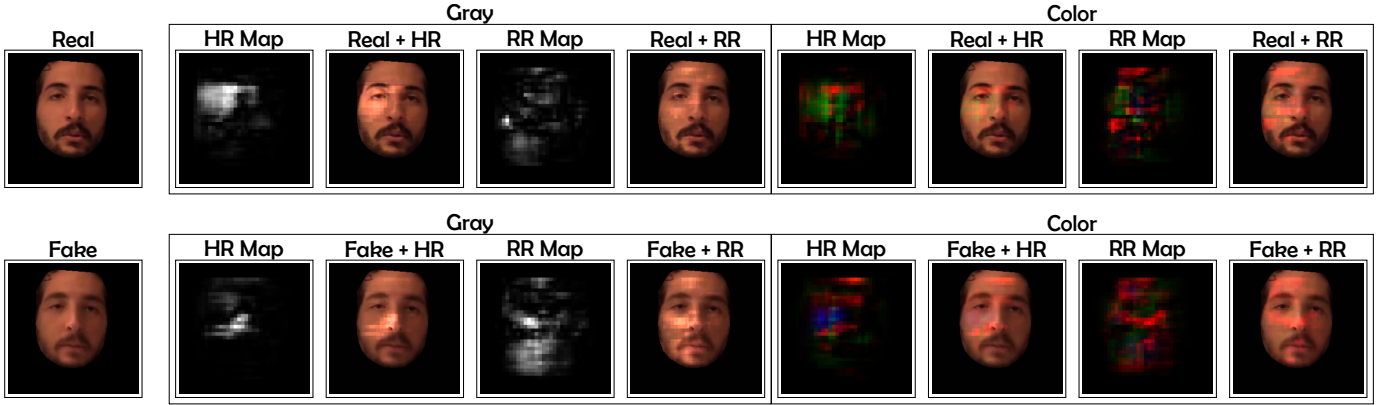


Fig. 3: Example of data augmentation. The real (top row) and the corresponding fake (bottom row) aligned face crops are augmented with the generated gray and color HR (*i.e.*, heart rate) and RR (*i.e.*, respiration rate) physiological maps.

the occlusion patch is moved and steps 3 and 4 are repeated. Additionally, the occlusion patch in step 3 is either applied to all color channels or separately for each color channel. This results in two sets of physiological maps – *gray* (a single difference value per patch) and *color* (three difference values per patch for each color channel). The overall result of this process is the generation of visual representations that encode the relative contribution of all masked-out regions for the accurate estimation of the two physiological signals. Figure 2 offers a visualization of the steps involved in the generation of the proposed physiological maps.

C. Physiological Maps Augmentation

Given the visual representations of the physiological signals, one approach to utilize this additional information is to augment the aligned face crops. In particular, this work employs simple multiplication of the aligned face crops and the corresponding physiological maps as a data augmentation strategy. This process results in an augmented dataset that incorporates the information for the two physiological signals. An example of data augmentation is provided in Figure 3.

D. Cross-Modal Learning

Given the physiological maps, another approach to utilize this additional information is to learn the relationship between the aligned face crops and physiological maps from data. This work proposes a multimodal fusion of the aligned face crops and physiological maps with a graph-based model. For a given video segment, each aligned face crop and the corresponding physiological map are used as nodes. The node features for both the crops and maps are obtained using ResNet-18 [25] (with the pre-trained weights on ImageNet [26]). Inside the graph, each of the maps is connected to all crops with an edge with weight given by the cosine similarity between the node features. This representation enables the model to learn the weights of each map/crop. Given this graph representation, this work uses a Graph Convolutional Network (GCN) [27] to learn the node features. Then the learned node features are

concatenated with the features generated by the video stream of the bi-stream classifier network as shown in Figure 1.

The graph-based multimodal feature fusion consists of N nodes, $X = \{x_1, x_2, \dots, x_N\}$, corresponding to each physiological map and aligned face crop across T time steps. Each of the map nodes is connected to all aligned face crop nodes. With this graph representation, a Graph Convolutional Network performs message passing based updates in the graph to update the features for each node. Formally, given a node feature matrix of dimension $N \times d$, there exists an edge between all the pairs of nodes of the form $\langle \text{physiological map node} : \text{aligned face crop node} \rangle$; hence there are T^2 edges in the graph representation weighted by the similarity between the nodes they connect. The similarity criterion (*i.e.*, edge weight) between two nodes x_1 and x_2 is defined by the cosine similarity between the node features:

$$\text{EdgeWeight}(x_1, x_2) = \frac{\sum_{i=1}^d x_{1i} x_{2i}}{\sqrt{\sum_{i=1}^d (x_{1i})^2} \sqrt{\sum_{i=1}^d (x_{2i})^2}} \quad (6)$$

IV. EXPERIMENTS

This section describes the datasets used and the general setup of the experiments.

A. Datasets

The deepfake detection challenge (DFDC) dataset was released gradually with the preview dataset [28], comprising of 5214 videos and the complete dataset [1] with 119146 videos. The details of the video manipulations were not disclosed in order to represent the real adversarial space of facial manipulation. The manipulations can be present in either the audio or video or both of the channels. In order to bring out a fair comparison, this work uses 18000 videos⁶ of DFDC in all experiments. The videos are of ≈ 10 seconds duration each with a frame rate of 30 frames per second.

The Deepfake-TIMIT [2] dataset contains videos of 16 similar looking pairs of people, which are manually selected from

⁶Identical to those used in [22] and [29]

the publicly available VIDTIMIT⁷ dataset and manipulated using an open-source Generative Adversarial Network based approach. There are two different models for generating fake videos, one low quality (LQ), with 64×64 input/output size, and the other high quality (HQ), with 128×128 input/output size. Each of the 32 subjects has 10 videos, resulting in a total of 640 fake videos in the dataset. The videos are of ≈ 4 seconds duration each with a frame rate of 25 frames per second. In this dataset the audio channel is not manipulated.

B. Implementation

In all experiments the hyper-parameters are kept the same as in [22] that were arrived at by a series of ablation studies. The segment duration D is 1 second and the margin hyper-parameter for the contrastive loss is set to 0.99. The graph is composed of $30 * 30$ edges (with $T = 30$ time steps for DFDC) and $25 * 25$ edges (with $T = 25$ time steps for Deepfake-TIMIT). The node features size is $d = 512$ and the Graph Convolutional Network layer size is $(512, 16)$. The physiological map generation is performed on 36×36 down-sampled versions of the aligned face crops since this is the input size expected by the MTTs-CAN model. The occlusion patch size is 9×9 resulting in $36 \times 36 \times 3$ (color) or $36 \times 36 \times 1$ (gray) dimensional physiological maps. The models are trained with a batch size of 8 using the Adam [30] optimizer with default parameters ($\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$) and with a learning rate of 0.001. All models are implemented in PyTorch [31].

C. Evaluation

Similarly to [22], during model evaluation, the video segments $\{s_1^i, s_2^i, \dots, s_n^i\}$ and corresponding audio segments $\{a_1^i, a_2^i, \dots, a_n^i\}$ of a test video are passed through S_v and S_a . Then for each segment, a dissimilarity score $d_t^i = \|f_v - f_a\|_2$ is accumulated, where f_v and f_a are feature representations of the video and audio streams, respectively. Finally, the Modality Dissonance Score (MDS) is computed as the average of the accumulated dissimilarity scores:

$$MDS_i = \frac{1}{n} \sum_{t=1}^n d_t^i \quad (7)$$

In order to classify a test video as either real or fake, the MDS_i is compared with a threshold τ using $I\{MDS_i < \tau\}$, where $I\{\cdot\}$ denotes the logical indicator function and τ is determined using the train set. Then the Area Under the Curve (AUC) metric is used for evaluation, computed using video-wise real/fake classification. This evaluation strategy is consistent with the one in [22] and ensures the fair comparison between the performance of the different methods.

V. RESULTS

This section reports that the proposed physiology-based fake video detection method achieves competitive performance on both the DFDC and Deepfake-TIMIT datasets. Table I,

starting from second row, reports the classification results on the test set of the DFDC and DF-TIMIT datasets using the proposed data augmentation strategy. The results from MDS models using three of the proposed physiological maps for data augmentation are in par with the MDS method (without augmentation). The second and third column present the detection results on the test set of the Deepfake-TIMIT HQ and LQ dataset, respectively. Here the results from all MDS models using physiological maps for data augmentation are better than the original MDS method (*i.e.*, 96.8% for HQ and 97.9% for LQ).

TABLE I: Comparison of the video-level classification performance using the proposed physiological maps for data augmentation. The MDS models are evaluated on the DFDC and Deepfake-TIMIT HQ/LQ datasets with %AUC metric.

Maps/Datasets	DFDC	DF-TIMIT (HQ)	DF-TIMIT (LQ)
Not Augmented [22]	91.5	96.8	97.9
HR Gray Augmented	91.4	98.1	99.7
HR Color Augmented	91.6	99.2	1.0
RR Gray Augmented	91.5	1.0	99.2
RR Color Augmented	89.8	99.4	1.0

Comparison of the proposed learning-based fusion method and previous methods in terms of classification performance is provided in Table II. The table reports the performance of the best performing fusion method with HR color physiological maps. Using the HR (*i.e.*, heart rate) color maps (the best performing maps during data augmentation) for learning the fusion between physiological maps and aligned face crops with the Graph Convolutional Network yields performance of 93.1% which is better than the MDS method with 1.6%. By extension, this suggests that the proposed method outperforms the other visual-only and audio-visual approaches discussed in [22].

We note that the results from Graph Convolutional Network based model on Deepfake-TIMIT dataset are slightly worse than but comparable to the proposed approach using physiological maps for augmentation as given in Table I. Since the Deepfake-TIMIT dataset is a small dataset (only 1281 test video samples), the results are highly oscillating and a few misclassifications result in change of evaluation metric scores. Augmentation of the video frames with the generated physiological maps performs better on these datasets. Overall, methods involving fusion with HR Color maps either directly augmented or the relationships learned through Graph Convolutional Network, outperform previous methods as seen in Table II.

VI. DISCUSSION

Along with leading to an improved performance on AUC metric, our graph-based fusion framework is an interpretable model, that is, the edges in the graph (and hence the frames) contributing most to the given model decision can be identified based on the edge weights. The framework can be helpful in making insightful progress on the misclassified samples

⁷<https://conradsanderson.id.au/vidtimit>

TABLE II: Comparison of the video-level classification performance using the proposed physiological maps for GCN-based learning with other methods on the DFDC and Deepfake-TIMIT HQ/LQ datasets using %AUC metric.

Methods/Datasets	DFDC	DF-TIMIT (HQ)	DF-TIMIT (LQ)
Capsule [32]	53.3	74.4	78.4
Multi-task [33]	53.6	55.3	62.2
HeadPose [5]	55.9	53.2	55.1
Two-stream [34]	61.4	73.5	83.5
VA-LogReg [35]	66.2	77.3	77.0
Meso4 [16]	75.3	68.4	87.8
Xception-c23 [3]	72.2	94.4	95.9
DSP-FWA [10]	75.5	99.7	99.9
Siamese [29]	84.4	94.9	96.3
MDS [22]	91.5	96.8	97.9
Ours	93.1	96.2	97.8

by providing frame-based information from the graph edge weights joining two modalities. This interpretability aspect can be very useful in real-world setting where it becomes important to work out the causes of misclassifications.

The current work investigates the fusion/augmentation of only one physiological measurement with the video stream. Given that both physiological signals might provide useful and complementary information for the task of fake video detection, a direction for future work includes fusion of both types of physiological masks and video data at the same time.

MTTS-CAN used for estimation of the heart and respiration rate has been shown to perform well for frontal faces with limited movement. However, the data in the DFDC dataset is unconstrained and includes lots of head movement, lighting conditions, and backgrounds. We partially address this problem with the face alignment step to provide higher quality data for MTTS-CAN. We expect that the fake video classification performance will be even higher in case more advanced and accurate methods are used for alignment and estimation of physiological signals.

Finally we note that in [22] one of the arguments for the increased performance of the MDS model is due to the relaxed crop around the face region. This suggests that the contribution of the proposed physiological maps is bigger than the 1.6% because the visual inputs in this work are aligned faces (with removed background). Indeed, re-training the original MDS model with aligned faces yields a performance of 88.6% on the DFDC dataset which implies that the relative contribution of the HR color maps is an increase of 4.5%.

VII. CONCLUSION

We propose a novel method for generation of visual representations of physiological signals extracted from facial videos. Additionally, we propose two strategies for utilizing those representations for the problem fake video detection, either by augmenting the video modality with information from the physiology or by novelly learning the fusion of those two modalities with a Graph Convolutional Network. Experiments show that this novel physiological maps help to achieve competitive performance. Future work will focus

on the simultaneous fusion of both types of physiological maps and video stream and investigation on different similarity measurements used in the Graph Convolutional Network.

REFERENCES

- [1] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, and C. C. Ferrer, "The deepfake detection challenge dataset," 2020.
- [2] P. Korshunov and S. Marcel, "Deepfakes: a new threat to face recognition? assessment and detection," 2018.
- [3] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "FaceForensics++: Learning to detect manipulated facial images," in *Proceedings of the International Conference on Computer Vision*, 2019.
- [4] L. Verdoliva and P. Bestagini, "Multimedia forensics," in *Proceedings of the ACM International Conference on Multimedia*, 2019, pp. 2701–2702.
- [5] X. Yang, Y. Li, and S. Lyu, "Exposing deep fakes using inconsistent head poses," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 8261–8265.
- [6] L. Zheng, Y. Zhang, and V. L. Thing, "A survey on image tampering and its detection in real-world photos," *Journal of Visual Communication and Image Representation*, vol. 58, pp. 380–399, 2019.
- [7] Y. Mirsky and W. Lee, "The creation and detection of deepfakes: A survey," *ACM Computing Surveys*, 2021.
- [8] H. Mo, B. Chen, and W. Luo, "Fake faces identification via convolutional neural network," in *Proceedings of the ACM Workshop on Information Hiding and Multimedia Security*, 2018, pp. 43–47.
- [9] R. Durall, M. Keuper, F.-J. Pfrendt, and J. Keuper, "Unmasking deepfakes with simple features," 2020.
- [10] Y. Li and S. Lyu, "Exposing deepfake videos by detecting face warping artifacts," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [11] J. Straub, "Using subject face brightness assessment to detect 'deep fakes'," in *Proceedings of the Real-Time Image Processing and Deep Learning*, 2019.
- [12] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, and H. Li, "Protecting World Leaders Against Deep Fakes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [13] U. A. Ciftci, I. Demir, and L. Yin, "Fakecatcher: Detection of synthetic portrait videos using biological signals," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [14] P. Korshunov and S. Marcel, "Speaker inconsistency detection in tampered video," in *Proceedings of the European Signal Processing Conference*, 2018, pp. 2375–2379.
- [15] D. Guera and E. J. Delp, "Deepfake video detection using recurrent neural networks," in *Proceedings of the IEEE International Conference on Advanced Video and Signal Based Surveillance*, 2018, pp. 1–6.
- [16] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, "Mesonet: a compact facial video forgery detection network," in *Proceedings of the IEEE Workshop on Information Forensics and Security*, 2018, pp. 1–7.
- [17] X. Ding, Z. Raziei, E. C. Larson, E. V. Olinick, P. Krueger, and M. Hahsler, "Swapped face detection using deep learning and subjective assessment," 2019.
- [18] S. Tariq, S. Lee, H. Kim, Y. Shin, and S. S. Woo, "Detecting both machine and human created fake face images in the wild," in *Proceedings of the International Workshop on Multimedia Privacy and Security*, 2018, pp. 81–87.
- [19] R. Wang, F. Juefei-Xu, L. Ma, X. Xie, Y. Huang, J. Wang, and Y. Liu, "Fakespotter: A simple yet robust baseline for spotting ai-synthesized fake faces," 2020.
- [20] H. Khalid and S. S. Woo, "Oc-fakedect: Classifying deepfakes using one-class variational autoencoder," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 2794–2803.
- [21] X. Liu, J. Fromm, S. N. Patel, and D. McDuff, "Multi-task temporal shift attention networks for on-device contactless vitals measurement," 2020.
- [22] K. Chugh, P. Gupta, A. Dhall, and R. Subramanian, "Not made for each other- audio-visual dissonance-based deepfake detection and localization," in *Proceedings of the ACM International Conference on Multimedia*, 2020, pp. 439–447.

- [23] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Li, "S3fd: Single shot scale-invariant face detector," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 192–201.
- [24] T. Baltrusaitis, A. Zadeh, Y. C. Lim, and L.-P. Morency, "Openface 2.0: Facial behavior analysis toolkit," in *Proceedings of the IEEE International Conference on Automatic Face and Gesture Recognition*, 2018, pp. 59–66.
- [25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [26] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [27] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," 2017.
- [28] B. Dolhansky, R. Howes, B. Pflaum, N. Baram, and C. C. Ferrer, "The deepfake detection challenge (dfdc) preview dataset," 2019.
- [29] T. Mittal, U. Bhattacharya, R. Chandra, A. Bera, and D. Manocha, "Emotions don't lie: An audio-visual deepfake detection method using affective cues," in *Proceedings of the ACM International Conference on Multimedia*, 2020, pp. 2823–2832.
- [30] D. P. Kingma and J. Ba, "Adam: a method for stochastic optimization," 2014.
- [31] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in *Proceedings of the NeurIPS Autodiff Workshop*, 2017.
- [32] H. H. Nguyen, J. Yamagishi, and I. Echizen, "Capsule-forensics: Using capsule networks to detect forged images and videos," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 2307–2311.
- [33] H. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen, "Multi-task learning for detecting and segmenting manipulated facial images and videos," in *Proceedings of the IEEE International Conference on Biometrics Theory, Applications and Systems*, 2019, pp. 1–8.
- [34] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis, "Two-stream neural networks for tampered face detection," 2018.
- [35] F. Matern, C. Riess, and M. Stamminger, "Exploiting visual artifacts to expose deepfakes and face manipulations," in *Proceedings of the IEEE Winter Applications of Computer Vision Workshops*, 2019, pp. 83–92.