

Motion Matters: Neural Motion Transfer for Better Camera Physiological Sensing

Akshay Paruchuri¹, Xin Liu², Yulu Pan¹, Shwetak Patel², Daniel McDuff², Soumyadip Sengupta¹

¹UNC Chapel Hill ²University of Washington

{akshay, ronisen}@cs.unc.edu, {xliu0, shwetak, dmcduff}@cs.washington.edu, ypan1@unc.edu

Abstract

Machine learning models for camera-based physiological measurement can have weak generalization due to a lack of representative training data. Body motion is one of the most significant sources of noise when attempting to recover the subtle cardiac pulse from a video. We explore motion transfer as a form of data augmentation to introduce motion variation while preserving physiological changes. We adapt a neural video synthesis approach to augment videos for the task of remote photoplethysmography (PPG) and study the effects of motion augmentation with respect to 1) the magnitude and 2) the type of motion. After training on motion-augmented versions of publicly available datasets, the presented inter-dataset results on five benchmark datasets show improvements of up to 75% over existing state-of-the-art results. Our findings illustrate the utility of motion transfer as a data augmentation technique for improving the generalization of models for camera-based physiological sensing. We release our code and pre-trained models for using motion transfer as a data augmentation technique on our project page: <https://motion-matters.github.io/>

1. Introduction

Scalable health sensors enable frequent, opportunistic, and more equitable access to vital information about the body’s internal state. Cameras are some of the most versatile and widely available sensors. Videos capture spatial, temporal, and ultimately frequency-specific information making them suitable for imaging dynamic processes, even below the surface of the skin [30]. Camera-based measurement of cardiac signals is one such application [23], in which cameras are used to measure the pulse via light reflected from the body, a principle known as photoplethysmography (PPG) [2, 43]. The PPG signals can be used to derive respiration [32], heart rate variability [32], arrhythmia [33], and blood pressure [14]. As a result this technology has the potential to turn webcams and smartphones into meaningful health sensors.

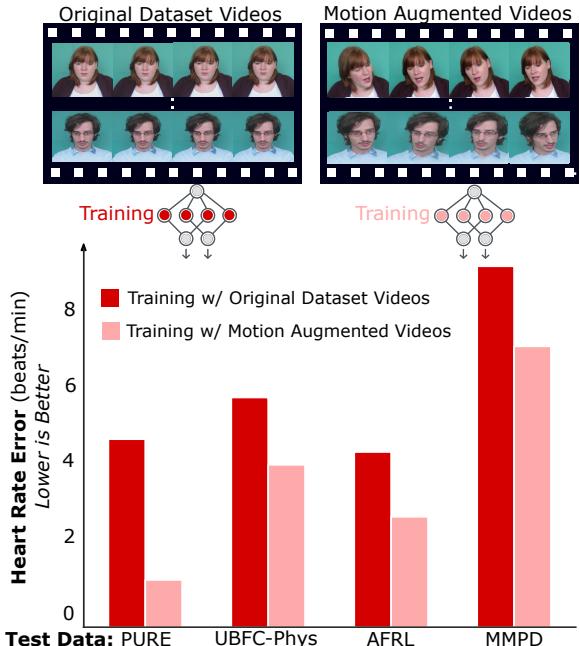


Figure 1: **Motion augmentation improves rPPG.** We present the first neural motion augmentation pipeline for the task of remote PPG estimation and empirically show it reduces error in heart rate estimation by up to 75%.

However, unlike traditional medical sensors, extracting physiological signals from a video requires more than filtering and simple signal processing. The state-of-the-art (SOTA) algorithms are supervised neural models [4, 38, 51, 18, 52]. Despite the prowess of these models, they are inherently limited by the diversity of the data used to train them. Public datasets (e.g., UBFC-rPPG [3], PURE [39]) serve as an extremely valuable resource for the research community, containing videos and synchronized physiological gold-standard measurements making them suitable for training and testing models. Building datasets such as these is challenging for two reasons: (1) collecting videos with gold-standard signals from a medical-grade sensor is time consuming and labor intensive, (2) it requires storing

and distributing privacy sensitive biometric data. Therefore, more data efficient methods for training rPPG sensing models would be desirable.

Synthetic data are a powerful resource in machine learning. The two main sources of synthetic data are (1) parametric computer graphics engines and (2) statistically-based generative machine learning models. Data created using these approaches have been used successfully for many computer vision tasks, including face detection, landmark localization, face parsing and face recognition [24, 48, 26], body pose estimation [35] and eye tracking [49, 40].

However, creating synthetic data that preserve the subtle and nuanced peripheral pulse in a video is non-trivial. McDuff et al. [27] released a large dataset (2,800 videos) of avatars and cardiac signals; however, their computer graphics pipeline had an extremely high overhead. Creating a pipeline for generating videos of avatars required years worth of investment, procuring assets (e.g., 3D facial scans and environments), building a parameterized rendering pipeline, and then generating each video frame-by-frame. Wang et al. [47] used a learning based method to generate synthetic videos given a reference image and target PPG signal. Their creative approach successfully incorporated PPG signals producing videos that benefited training. However, the videos created lacked the visual fidelity of other synthetics or real video datasets, and their pipeline involved several relatively complex components.

We question whether existing motion transfer algorithms can be used effectively for augmenting rPPG video data and explore what steps need to be taken to achieve optimal results. Our main contributions are as follows: (1) We perform a systematic investigation of the impact of motion augmentation on the physiological information within the video and in-turn the corresponding labels, (2) through experimentation, we provide quantitative, empirical evidence that training with certain kinds of motion-augmented data is effective for camera-based physiological measurement algorithms, and (3) we demonstrate, through inter-dataset results, the usefulness of motion augmentation for improving the generalization of models for camera-based physiological sensing. We achieve state-of-the-art results on multiple public benchmark datasets, including those with significant motion. We summarize the key findings of this paper about the effectiveness of motion transfer as a data augmentation tool in Sec. 5. We provide our code for augmenting datasets, training using these data, and pre-trained models trained on motion-augmented data (all assets are released with responsible use licenses [5]).

2. Background

Generative Synthetics for Training Models: Statistical generative models [9, 16, 15, 37, 12, 6] capture a probabilistic representation of a dataset from which samples can be drawn. These models are typically trained to mimic the

distribution of the training set and can be trained without the need for labels, allowing large sets of data to be used. Facial video generation using generative models has advanced rapidly over recent years [17, 34]. Numerous image-driven works have accomplished the ability to separate identity and pose in source and driving images used for high quality, robust video generation using generative adversarial networks (GANs) [53, 36, 45, 13]. Image-driven facial video generation methods attempt to preserve the identity of a given source image while manipulating the pose based on a driving video to generate a new video. The identity from the driving video is excluded with the help of a keypoint-based motion transfer approach, where keypoints are predicted for both a source image and a driving image in order to model local motion using shifts in the corresponding keypoints [36, 45, 13]. Face video generation that is achieved by using keypoints that take pose and expression into account can be successful for the task of head video generation, but can at times have a loss in source image identity and unwanted temporal artifacts [36, 53, 13]. FaceVid2Vid [45] utilizes canonical keypoints in addition to source and driving image keypoints in order to capture a target person’s geometry signature, which includes the shape of the target’s face, nose, and eyes. This allows for improved head video generation that minimizes source identity loss while effectively transferring motion from a driving video.

rPPG Models: The principle that photoplethysmography could be performed with a camera and without contact with the body was established by Blazek et al. [2] and replicated in a series of following experiments [41, 43]. The application of more advanced signal processing methods helped make measurement somewhat more robust under real-world conditions [32, 46], as did leveraging knowledge of physiological and physical properties [46]. Yet, these models were still very sensitive to body motions. Neural data-driven models currently achieve state-of-the-art results in most cases [4, 51, 18, 52, 19], but are a function of the data used to train them. While intra-dataset performance is generally strong, inter-dataset performance is often substantively worse. In order to alleviate the dependency on labeled data, several researchers have proposed unsupervised learning procedures [8, 44, 50]. However, most require fine-tuning on a labeled set and also reveal that supervised learning still holds some additional benefit. As an alternative or a complement, generative methods have been suggested to “create” data [25, 47].

rPPG Datasets: As with many health applications, those working in camera physiological measurement face challenges associated with collecting and managing data. Public datasets (such as UBFC-rPPG [3], PURE [39], VIPL-HR [29]) are valuable resources. However, given the challenging nature of the rPPG task researchers have col-

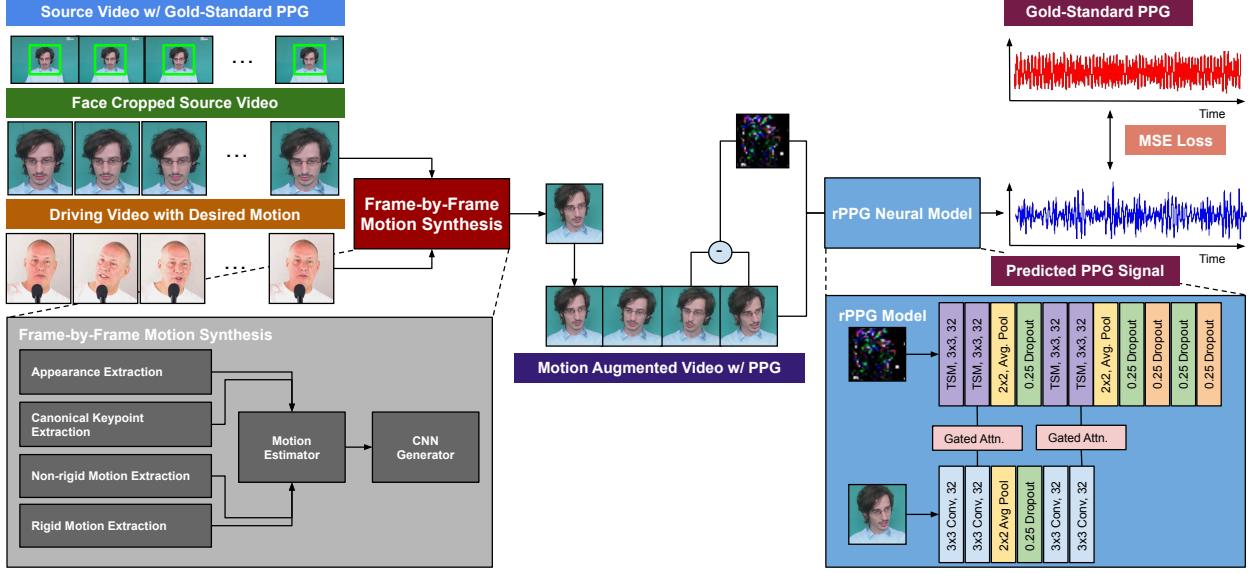


Figure 2: Motion augmentation and training pipeline. We augment each frames of a source video with corresponding frames of a randomly selected driving video to create an augmented video with the identity of the source video and motion of the target video. We then train a PPG estimation network on the augmented video with Mean Squared Error Loss.

lected and released data under heavily constrained conditions with very little physical motion. More recent datasets (such as UBFC-PHYS [28] and MMPD [42]) contain larger and more natural motions. However, the baseline results on these datasets are not very strong.

3. Motion Augmented rPPG Video Pipeline

We propose neural motion transfer as a data augmentation technique to train machine learning models for predicting physiological measurements, specifically Photoplethysmography (PPG) signal, from facial videos. First, we describe our proposed pipeline to augment facial videos with naturalistic human head motion and expression in section 3.1. Neural motion transfer algorithms often use generative models to synthesize new videos of a person by transferring the rigid head motion and non-rigid facial expressions from a driving video of another person. Since these models generate image pixels from scratch, it is possible that images generated by neural motion transfer algorithms can destroy the underlying physiological signal. Thus, in section 3.2, we provide qualitative evidence to prove that neural motion transfer algorithms do not destroy the original PPG signal, and the original heart rate is preserved. This allows us to effectively use neural motion transfer as a data augmentation technique for training rPPG networks.

3.1. Motion Augmentation Pipeline

In a camera-based physiological sensing (e.g., rPPG) task, a machine learning model is trained on facial videos with time-aligned physiological labels. These may take the form of continuous waveforms (e.g., a gold-standard PPG

or a respiration wave) or vital statistics (e.g., heart or breathing rates). In this project, we consider video labels in the form of a PPG signal. The goal of designing a data augmentation strategy is to apply more naturalistic motion to the facial videos without changing the PPG labels.

To apply naturalistic motion to these facial videos, we consider neural talking-head video synthesis models that transfer more naturalistic motion from a *driving* video of a person to the *source* video with PPG signal labels. Our goal is to find a neural motion transfer algorithm that can: (a) inject a large variety of rigid and non-rigid head motions into the source video, (b) not introduce any artifacts that significantly degrade the generated video quality, and (c) maintain the key properties of the underlying PPG signal in terms of frequency information indicating physiological signals like heart rate.

Our pipeline takes in a source video with PPG signal labels from the training data, \mathbf{S} , and a driving video, \mathbf{D} , randomly selected from a curated driving video set as inputs for motion augmentation. Both \mathbf{S} and \mathbf{D} can be represented as a sequence of frames, respectively $\{s_1, s_2, \dots, s_n\}$ and $\{d_1, d_2, \dots, d_n\}$. Motion is transferred from driving video \mathbf{D} to source video \mathbf{S} on a frame-by-frame basis, such that an output video \mathbf{Y} represents the motion-augmented sequence of frames $\{y_1, y_2, \dots, y_n\}$. Thus we search for a motion transfer algorithm $M(\cdot; \theta)$, such that $y_t = M(s_t, d_t; \theta)$.

We choose Face-Vid2Vid [45], a neural talking-head synthesis model for transferring motion from a driving video to a source video. The original Face-Vid2Vid paper was intended for teleconferencing applications where a motion-augmented video is generated from a single source

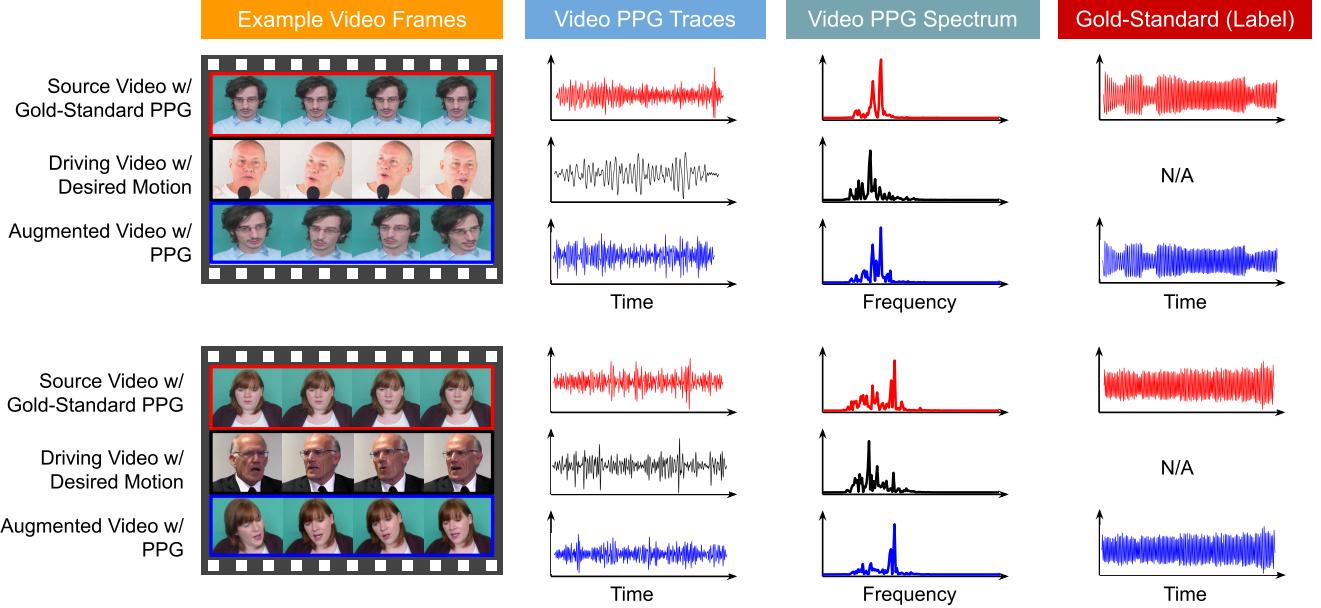


Figure 3: **Preserving physiological signals in motion augmented videos.** We show that applying neural motion transfer preserve the physiological signal corresponding to the heart-rate present in the peak of the frequency spectrum of the source and the augmented video.

image using a driving video. In contrast, we redesign and reimplement this algorithm such that each frame of the source video is augmented with motion from the corresponding frame of the driving video. The motion-augmented video \mathbf{Y} , along with the original PPG signal label, is ultimately used as training data for various deep learning-based camera physiological measurements. This pipeline is shown in Figure 2.

Source Video Datasets: We utilize the UBFC-rPPG [3] and PURE [39] rPPG video datasets as source videos. The UBFC-rPPG dataset contains videos with a very minimal amount of both rigid motion and non-rigid motion, making them ideal for motion augmentation. The PURE dataset contains videos of various tasks with a variety of constrained rigid and non-rigid motion.

Driving Video Datasets: The driving video datasets used include a self-captured, constrained driving video set (CDVS) and the TalkingHead-1KH [45] dataset. The CDVS contains 90 self-captured videos by 5 subjects with heavily constrained, unnatural motion used only for ablation studies to understand the impact of augmenting data with various degrees of rigid and non-rigid motion. The CDVS will be released in the future for research purposes. Talkinghead-1KH is a publicly available, large-scale talking-head video dataset used as a benchmark for Face-Vid2Vid [45] and entirely sourced from YouTube videos. It contains 180K unconstrained videos of people speaking in a variety of real-world contexts, leading to a rich diversity in both rigid and non-rigid motion.

Deep Networks for estimating PPG signal: For our experiments, we focus on using TS-CAN [18] to predict the 1st-

order derivative of the PPG signal after training on videos augmented with motion. We also use DeepPhys [4] and EfficientPhys [19] to highlight the consistent benefits of motion augmentation across different neural models.

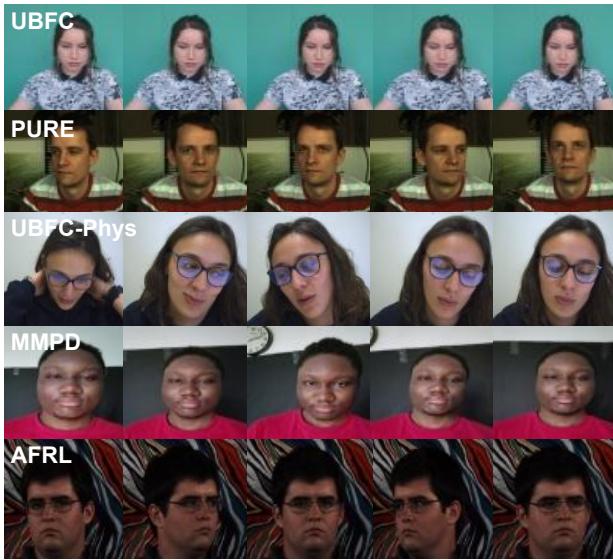
3.2. The Effect of Motion Transfer on PPG

Neural Motion Transfer algorithms are based on generative models where every pixel of the generated image is synthesized by a neural network. While these algorithms succeed in producing photorealistic facial images that are indistinguishable from real images, it is not obvious if the synthesized videos can preserve the underlying PPG signal.

In an ideal world, a motion transfer algorithm is expected to perturb the PPG signal since head motion will induce certain changes in raw pixel intensities. However, the frequency domain analysis of the PPG signal should preserve the peaks related to the heart rate of the patient. It is highly unlikely that the peak frequency of head motion and heart rate will be exactly the same.

Thus, our goal is to first analyze if the motion transfer algorithm of Face-Vid2Vid [45] can preserve the peak heart rate indicated in the frequency domain analysis of the PPG signal extracted from the source video and the synthesized video. In Figure 3, we qualitatively analyze the time-domain and frequency domain PPG signals extracted from the source and the synthesized (augmented) video. We choose a simple unsupervised algorithm, POS [46], for extracting the PPG signal from all the facial videos to focus more on the original signal contents in the videos. We observe that the most prominent frequency peak, corresponding to the heart rate, is the same for the source video and

Table 1: A Summary of the rPPG Benchmark Datasets.



Dataset	Subjects / Videos	Motion Tasks
UBFC-rPPG	42 / 42	Stationary
PURE	10 / 59	Stationary, Talking, Rotation, Translation
UBFC-Phys	56 / 168	Stationary, Talking, Head Rotation
MMPD	33 / 660	Stationary, Talking, Walking, Head Rotation
AFRL	25 / 300	Stationary, Head Rotation

the augmented video. This appears to also hold true across different appearances and motion conditions, both in the source videos and the driving videos. We present additional results in the Appendix D.1. Thus, we can effectively claim that motion transfer algorithms like Face-Vid2Vid [45] do preserve the underlying physiological signal, like heart rate, and they can be a very effective tool for large scale augmentation of training videos for PPG estimation tasks. Our quantitative experimental results show that deep neural networks for camera physiological measurement can take advantage of this to significantly improve model performance by training on motion-augmented data.

4. Experiments

We consider five datasets for training and evaluation, **UBFC-rPPG** [3], **PURE** [39], **UBFC-PHYS** [28], **AFRL** [7], and **MMPD** [42] (see Table 1). They consist of facial videos and corresponding gold-standard PPG signal labels. We use some of these datasets for augmentation with neural motion transfer and training the rPPG models, and use the rest to evaluate different aspects of the effectiveness of neural motion augmentation. To our knowledge, we perform the most extensive inter-dataset evaluation of PPG estimation to date, testing on five independent test datasets.

Implementation Details: The predicted PPG signals were filtered using a band-pass filter with cut-offs 0.75 Hz and 2.5 Hz. The heart rate was calculated based on the predicted PPG signal using the Fast Fourier Transform (FFT), with a measurement window of the video length for UBFC-rPPG, PURE, UBFC-PHYS, and MMPD. To

evaluate the AFRL dataset, a measurement window of 30 seconds was utilized for heart rate calculations. All networks were trained using an NVIDIA RTX A4500 and PyTorch [31] implementations in the publicly available rPPG-Toolbox [20]. A cyclic learning rate scheduler was utilized with 30 epochs, a learning rate of 0.009, and a batch size of 4 for both training and inference.

4.1. Training with Motion Augmented Data

In Table 11, we compare the performance of a supervised PPG estimation network, TS-CAN [18], trained on existing video datasets and motion-augmented versions of those datasets. We also show the performance of unsupervised methods for comparison. For the sake of space and clarity, the following tables only show the mean absolute error (MAE) in heart rate estimation for each video (and the corresponding mean absolute percentage error (MAPE)). Equivalent tables with root mean squared error (RMSE) and Pearson correlation metrics can be found in Appendix C. The driving videos used for augmentation in Table 11 contain significant amounts of unconstrained motion – both rigid and non-rigid.

We observe that training TS-CAN on augmented videos produces state-of-the-art (SOTA) performance in most cases. Additionally, we observe that in most cases, the augmented versions outperform the non-augmented versions, with a gain in performance up to 75% and an average gain of 26%. However, when comparing the performance of MAPURE versus PURE when tested on UBFC-PHYS, we note a minor drop in performance rather than an improvement due to the difficulty in effectively augmenting the PURE dataset. This is because the PURE dataset already contains significant amounts of rigid motion, and when augmented, it may provide training data with artifacts that make the learned rPPG task less useful in the face of a highly unconstrained dataset with natural rigid and non-rigid motion.

Details: We utilize all downloadable videos from the TalkingHead-1KH [45] dataset as our driving videos for augmenting various PPG estimation datasets with motion. We analyze the videos using OpenFace [1] to obtain the intensity (0 to 5) of 17 Facial Action Units (AUs) and the head pose rotations R_x , R_y , and R_z in radians (rad). To generate MAUBFC-rPPG, we choose driving videos from a pool of 60 driving videos with a range of mean standard deviation in head pose rotations from 0.10 to 0.14 rad to augment as much rigid motion as possible into a source video dataset that has very little of both rigid and non-rigid motion. We do not constrain for non-rigid motion in this case, so we observe a wide range of mean standard deviation in facial AUs from 0.15 to 0.5 intensity. To generate MAPURE, we choose driving videos with a range of mean standard deviation in facial AUs from 0.45 to 0.55 intensity to augment as much non-rigid motion as possible into a source

video dataset that has very little non-rigid motion. We do not constrain for rigid motion in this case, so we observe a wide range of mean standard deviation in head pose rotations from 0.03 to 0.14 rad.

4.2. Effect of Motion Types

A key question in designing a motion augmentation strategy is deciding what type of motion should be applied to obtain the best performance on a certain evaluation dataset. To answer this question, we separately analyze two types of motion: rigid and non-rigid, by augmenting training data with different magnitudes of motion. Rigid motion refers to head pose rotation, while having minimal change in facial action units or expressions. Non-rigid motion refers to changes in facial expression, i.e. motion in facial action units for various tasks like talking, while having minimal head pose rotation.

Rigid Motion: For rigid motion, we consider UBFC-rPPG as training data, which has very little head motion and AFRL as test data which has large variations in rigid head motion. We classify videos in the AFRL dataset into different rigid head motion categories: ‘very small motion’, ‘small motion’ (10 deg rotation per sec), and ‘large motion’ (30 deg rotation per sec). Based on this categorization, we also select driving videos from our captured CDVS to have ‘small motion’ and ‘large motion’ using the mean standard deviation in estimated head pose rotations across all the frames of a video. Specifically, for ‘small motion’ we used mean standard deviation between 0.03 to 0.07 rad and for ‘large motion’ between 0.10 to 0.14 rad. These parameters are chosen to roughly match the distribution of head pose rotation in ‘small motion’ and ‘large motion’ categories of AFRL. We then use these videos from the CDVS dataset to augment the source videos of UBFC-rPPG to create 3 separate categories of augmented videos for ‘very small motion’ (which is the original UBFC-rPPG dataset), ‘small motion’, and ‘large motion’ respectively. We then train TS-CAN on augmented data in each category and test on the same categories of the AFRL dataset. We present these results in Table 10.

We observe that when the test data of AFRL has ‘very small motion’ or ‘small motion’, augmenting UBFC-rPPG with small motion performs the best. In fact, augmenting with large motion worsens the result by 19% in this case. However, when testing on the ‘large motion’ split of AFRL, UBFC-rPPG augmented with ‘large motion’ outperforms ‘small motion’ by 13.5% and ‘very small motion’ by 52%.

Non-rigid Motion: For non-rigid motion, we also consider UBFC-rPPG as training data since it has very little motion, and the speech task of the PURE dataset [39] as the test data which has significant non-rigid head motion. We also augment the UBFC-rPPG dataset with non-rigid head motion from our captured CDVS with ‘small’ and ‘large’

non-rigid motions and minimal rigid motion. For this experiment, we define small non-rigid motion to have a range of mean standard deviation in facial action units from 0.15 to 0.25 intensity and large non-rigid motion to have a range of mean standard deviation in facial action units from 0.45 to 0.55 intensity. We train TS-CAN on ‘small’ and ‘large’ motion augmented versions of UBFC-rPPG and test it on the speech task of PURE, in which recorded participants are asked to talk while avoiding head movements as much as possible. We observe that augmenting UBFC-rPPG with ‘large’ non-rigid motion improves over ‘very small motion’ (original UBFC-rPPG) by 89.2% and over ‘small’ non-rigid motion by 37%.

4.3. Effect of Multiple Augmentations

We consider whether it is plausible to augment the same source video with multiple driving videos using neural motion transfer. Thus, the newly augmented dataset has the same number of identities as the original dataset but a significantly larger variation in motions. Our goal is to analyze how many times one can augment a single source video before the performance starts to saturate or drop. We consider UBFC-rPPG as training data that we augment with randomly sampled driving videos from the TalkingHead-1K dataset to produce MAUBFC-rPPG. We augment the same source video from 1 to 4 times with different driving videos and evaluate on the PURE [39] dataset and the UBFC-PHYS [28] dataset and report the results in Table 12. We notice that the results saturate pretty quickly and can start to decline after augmenting more than 2 times.

4.4. Synthetic vs Naturalistic Head Motion

In order to further evaluate the impact of motion transfer as a data augmentation technique, we explore whether data augmented with natural head motion using a neural motion transfer algorithm is better than augmenting data with synthetically generated motion using parametric motion animation, as used in the SCAMPS dataset [27]. The SCAMPS dataset consists of synthetic human heads that can be rigged to induce parametric motion. We consider 200 such samples from the SCAMPS dataset that consist of significant synthetically generated rigid and non-rigid head motion (ID 1801 to 2000) as SCAMPS-200 (Motion). We then take instances from the SCAMPS dataset with no head motion (ID 1 to 200) and augment them with naturalistic head motion using our motion synthesis pipeline and a subset of driving videos from the TalkingHead-1KH dataset to produce MASCAMPS-200. We choose driving videos with a range of mean standard deviation in AUs from 0.35 to 0.40 intensity and a range of mean standard deviation in head pose rotations from 0.05 to 0.125 rad. Note that both SCAMPS-200 (Motion) and MASCAMPS-200 consist of synthetics with the same number of identities, with the only difference

Table 2: **Evaluation across all datasets.** We motion-augment two training datasets, UBFC-rPPG and PURE, to create MAUBFC-rPPG and MAPURE, respectively. We observe that the motion-augmented versions produce significant improvements (shown in bold).

Training Set	Method	UBFC-rPPG		PURE		Testing Set		AFRL		MMPD	
		MAE \downarrow	MAPE \downarrow								
Unsupervised	Green	19.82	18.78	10.09	10.28	13.45	16.00	7.01	9.24	16.27	20.09
	ICA	14.70	14.34	4.77	4.47	8.00	9.48	6.77	8.96	13.10	16.33
	CHROM	3.98	3.78	5.77	11.52	4.68	6.20	5.41	7.95	8.85	11.93
	POS	4.00	3.86	3.67	7.25	4.62	6.29	6.93	10.00	8.18	11.12
UBFC-rPPG	TS-CAN	-	-	4.55	4.67	5.56	7.25	4.24	5.84	8.74	10.51
MAUBFC-rPPG	TS-CAN	-	-	1.14	1.30	3.93	5.24	2.67	3.65	6.80	7.97
PURE	TS-CAN	1.34	1.55	-	-	4.43	5.89	2.63	3.51	8.96	10.33
MAPURE	TS-CAN	1.03	1.17	-	-	4.39	5.90	2.37	3.26	8.08	9.54
MAUBFC-rPPG vs. UBFC-rPPG	-	-	+74.95%	+72.16%	+29.32%	+27.72%	+37.03%	+37.50%	+22.20%	+24.17%	
MAPURE vs. PURE	-	+23.13%	+24.52%	-	-	+0.90%	-0.17%	+9.89%	+7.12%	+9.82%	+7.65%

MAE = Mean Absolute Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation

Table 3: **Effect of Motion Types – Rigid.** We augment UBFC-rPPG with various types of rigid head motions and test on AFRL [7]. The best results are shown in bold.

Training Set	Rigid Motion	Testing Set			
		No Motion	Small Motion	Large Motion	All Motion
UBFC-rPPG	Very Small	1.00	2.28	7.59	4.72
MAUBFC-rPPG	Small	0.84	1.44	4.21	3.19
MAUBFC-rPPG	Large	1.00	1.78	3.64	3.39
OURS vs. BASELINE	-	+16.0%	+36.8%	+52.0%	+32.4%

Table 4: **Effect of Motion Types – Non-rigid.** We augment UBFC-rPPG with various types of non-rigid motions (expressions) and test on the speech task, in PURE [39]. The best results are shown in bold.

Training Set	Non-Rigid Motion	Testing Set	
		MAE \downarrow	MAPE \downarrow
UBFC-rPPG	Very Small	10.84	11.40
MAUBFC-rPPG	Small	1.86	2.94
MAUBFC-rPPG	Large	1.17	1.55
OURS vs. BASELINE	-	+89.2%	+86.4%

being synthetic and naturalistic head motion, respectively.

We train TS-CAN on both SCAMPS-200 (Motion) and MASCAMPS-200, and evaluated its performance on PURE and AFRL, as shown in Table 13. We observed that adding naturalistic motion improved performance by 13.2% on PURE and 31.1% on AFRL compared to synthetically generated motion. It is worth noting that the average time taken to add synthetic motion to each frame of a sequence is 37 seconds, compared to only 1.2 seconds for adding naturalistic motion using the neural motion transfer algorithm. For comparison, we also included real-world training data,

Table 5: **Effect of Multiple Augmentations.** Augmenting each source video of UBFC-rPPG 1x, 2x, 3x, and 4x, we test on PURE and UBFC-PHYS datasets. The best results are shown in bold.

Training Set	Size	Subjects	Testing Set			
			PURE		UBFC-PHYS	
UBFC-rPPG	42	42	4.55	4.67	5.56	7.25
MAUBFC-rPPG	42	42	1.14	1.30	3.93	5.24
MAUBFC-rPPG 2x	84	42	1.12	1.29	3.90	5.22
MAUBFC-rPPG 3x	126	42	1.11	1.27	3.97	5.31
MAUBFC-rPPG 4x	168	42	1.19	1.33	4.10	5.40
OURS vs. BASELINE	-	-	+2.63%	+2.31%	+0.76%	+0.38%

Table 6: **Synthetic vs Naturalistic Head Motion.** We compare the effect of adding head motions to SCAMPS and UBFC-rPPG and contrast this with using motion data in SCAMPS. Average time for augmenting each frame of a sequence is presented. The best results are shown in bold.

Training Set	Testing Set			
	PURE		AFRL	
SCAMPS-200 (No motion)	MAE \downarrow	MAPE \downarrow	MAE \downarrow	MAPE \downarrow
SCAMPS-200 (Motion)	5.38	5.42	7.25	10.20
UBFC-rPPG	4.55	4.67	4.72	6.59
MASCAMPS-200	4.67	4.22	5.00	6.69
MAUBFC-rPPG	1.14	1.30	3.24	4.37
MASCAMPS vs. SCAMPS	+13.2%	+22.1%	+31.1%	+34.4%
MAUBFC vs. UBFC-rPPG	+74.4%	+72.2%	+31.4%	+33.7%
Avg. Synth. Time	= time (in seconds) to synthesize a frame			

UBFC-rPPG, which showed that having real images significantly improved performance over synthetic images. Furthermore, the only way to augment real images is to use the neural motion transfer algorithm, as parametric rigged head motion cannot be applied to real data.

4.5. Effect of PPG Estimation Models

It is important to decouple any data augmentation technique from additional factors that affect its usefulness for a given set of training data. One such factor is the neural network model used for training and evaluation. Thus,

Table 7: **Effect of PPG Estimation Models.** We train different PPG estimation networks on UBFC-rPPG and MAUBFC-rPPG and evaluate on PURE. The best results are shown in bold.

Training Set	Method	Testing Set	
		PURE	
		MAE \downarrow	MAPE \downarrow
UBFC-rPPG	DeepPhys [4]	5.14	4.90
MAUBFC-rPPG	DeepPhys	1.24	1.56
UBFC-rPPG	EfficientPhys [19]	4.95	4.56
MAUBFC-rPPG	EfficientPhys	1.45	1.76
UBFC-rPPG	TS-CAN [18]	4.55	4.67
MAUBFC-rPPG	TS-CAN	1.14	1.30

in addition to TS-CAN, we evaluate two more rPPG models - DeepPhys and EfficientPhys - in Table 14. We utilize MAUBFC-rPPG as training data and evaluate on PURE. We observe that the results are reasonably consistent across neural rPPG models.

5. Discussion

Can motion augmented videos achieve SOTA results?

In Section 3.2, we demonstrated that motion transfer algorithms can be used to create motion-augmented videos while still preserving the variations in skin appearance from the cardiac pulse. This means we can augment the training data to create novel samples with greater variance than those in the original set. We conducted a set of systematic empirical validation studies that show that these videos can be used to effectively train rPPG models that generalize to independent benchmark datasets (see Table 11). Cross-dataset experiments show a 23.1% reduction in HR MAE on UBFC-rPPG when using the motion-augmented PURE datasets for training and a 74.95% reduction in HR MAE on PURE when using the motion-augmented UBFC-rPPG dataset for training. Other than PURE, the largest gains were observed training on MAUBFC-rPPG and testing on videos with large rigid and/or non-rigid head motions (UBFC-PHYS: 29.32%, AFRL: 37.03% and MMPD: 22.20% reduction in HR MAE).

What type of motion is best to augment? In learning tasks, designing training data that matches the distribution of the testing data is advantageous. Does augmenting motion in the training set that is similar to that in a testing set lead to optimal results? Our experiments show that this is the case for both rigid (see Table 10) and non-rigid (see Table 9) head motions. Beyond the type of motion, if the motions have a larger magnitude, then including larger magnitude motions in the training set seems to have a benefit.

Does the effect of motion augmentation saturate?

While source videos with gold-standard PPG data might be limited, it is possible to augment each source video with many different motions by leveraging large video datasets like TalkingHead-1KH [45]. We looked at whether aug-

menting the same source videos multiple times with different motions improved results. We observed that most of the improvements were obtained by augmenting the data once. Incremental improvements were obtained by augmenting a second or third set of videos, but the results quickly saturated (see Table 12). This is presumably due to the fact that we were not augmenting other aspects of the subjects' appearance (e.g., skin tone, identity, etc.).

Is natural motion augmentation best? Finally, there are different methods for synthesizing motion in video data. State-of-the-art synthetic datasets are generated using parametric computer graphics, but they require a large amount of computational resources. As a result, if the motions present in those datasets are sub-optimal, it is costly to remedy. Can motion augmentation add motions to these datasets "cheaply" and still obtain the performance benefits of graphics approaches? Our results in Table 13 suggest that the motion in the SCAMPS dataset is sub-optimal when tested on PURE and AFRL. We were able to obtain a performance gain by using our simple motion augmentation.

What are the limitations of our method? There are several limitations that we would like to highlight. First, detecting artifacts in augmented videos is not always trivial, and we used motion driving videos without extreme motions to mitigate the chance of augmented videos with unnatural artifacts. We did not conduct an extensive investigation to determine if other physiological changes (e.g., respiration) that might be correlated with the PPG signal are preserved in the augmented videos. However, empirically we have shown that these data can be used to effectively train *heart rate* estimation models. We did not thoroughly test whether the waveform dynamics, beyond the dominant frequency, were faithfully preserved in the augmented videos. For tasks such as blood pressure estimation from PPG waveforms, morphological information is important. Our method does not address diversity across other dimensions, particularly identity diversity. The augmented datasets we produced, while contributing to significant improvements over the baselines, only contain examples from the same number of subjects as the original dataset. Other synthetic generation techniques [47] could help in these regards alongside more generic neural rendering approaches such as ours.

6. Conclusion

Motion artifacts are a significant challenge in camera physiological measurement. The PPG signal presents only very subtle changes in diffuse light reflections from the skin, whereas motion of the head causes large changes in specular reflections. We have shown that neural motion augmentation can be used to create training data with more motion, while still preserving the pulse signal. Motion augmented data leads to up to 75% reduction in error in cross-dataset experiments compared to training with unaugmented data.

A. Overview of Appendices

Our appendices contain the following additional details and results:

- Section B, and the corresponding Table 8, describe and show intra-dataset results using the PURE [39] dataset.
- Tables 9, 10, 11, 12, 13, and 14 in Section C include additional metrics, RMSE and the Pearson correlation coefficient, for experimental results included in the main paper. We also provide Scatter and Bland-Altman Plots in Figures 4 and 5 that correspond to the overall results shown in Table 11. Section C.1 contains additional details regarding our experimental process.
- Section D briefly describes additional materials that we provide for research purposes, including our motion augmentation pipeline code, pre-trained models, and motion analysis scripts. Additionally, Section D.1 shows more qualitative examples of the effects of motion augmentation on the underlying PPG signal.
- Sections E.1, E.2, and E.3, provide further details on source, driving, and evaluation datasets used in the main paper.
- Section F is our broader impact statement.

B. Intra-dataset Results

We include intra-dataset results not included in the main paper here for reference. We utilize all of the tasks from the PURE dataset. We train on subjects 1, 2, 3, 4, and 5 and then test on subjects 6, 7, 8, 9, and 10. We then train on subjects 6, 7, 8, 9, and 10 and then test on subjects 1, 2, 3, 4, and 5. We average the results from these two experiments and repeat the aforementioned process for the motion-augmented version of PURE. We find that motion augmentation helps as an intra-dataset augmentation technique.

Table 8: PURE Intra-dataset Results. We use motion augmentation to augment half of the subjects in the PURE dataset at a time, while testing on the corresponding other half. The averaged results are shown below, with the best result in each column bolded.

Training Set	Testing Set PURE			
	MAE \downarrow	RMSE \downarrow	MAPE \downarrow	$\rho \uparrow$
PURE	2.52	8.92	2.55	0.92
MAPURE	1.61	5.50	1.77	0.97
OURS VS. BASELINE	+36.1%	+38.34%	+30.59%	+5.43%

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation

C. Experimental Results

The following section contains tables that include additional metrics, RMSE and the Pearson correlation coefficient, for experimental results already included in the main paper. We also provide scatter and Bland-Altman plots in Figures 4 and 5 that correspond to results shown in Table 11.

Table 9: Effect of Motion Types – Non-rigid. We augment UBFC-rPPG with various types of non-rigid motions (expressions) and test on the speech task, in PURE [39]. The best results are shown in bold.

Training Set	Non-Rigid Motion	Testing Set Non-rigid Motion Task			
		MAE \downarrow	RMSE \downarrow	MAPE \downarrow	$\rho \uparrow$
UBFC-rPPG	Very Small	10.84	24.64	11.40	0.46
MAUBFC-rPPG	Small	1.86	2.79	2.94	0.99
MAUBFC-rPPG	Large	1.17	1.90	1.55	0.99
OURS VS. BASELINE		+89.21%	+92.29%	+86.40%	+0.00%

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation

C.1. Experimental Details

The predicted PPG signals were filtered using a band-pass filter with cut-offs 0.75 Hz and 2.5 Hz. The heart rate was calculated based on the predicted PPG signal using the Fast Fourier Transform (FFT), with a measurement window of the video length for UBFC-rPPG, PURE, UBFC-PHYS, and MMPD. To evaluate the AFRL dataset, a measurement window of 30 seconds was utilized for heart rate calculations. All networks were trained using an NVIDIA RTX A4500 and PyTorch [31] implementations in the publicly available rPPG-Toolbox [20]. All pre-processing steps and evaluation was also done in a reproducible fashion using the toolbox. The AdamW [21] optimizer, a mean squared error (MSE) loss, and a cyclic learning rate scheduler was utilized with 30 epochs, a learning rate of 0.009, and a batch size of 4 for both training and inference.

For both the UBFC-rPPG dataset and PURE datasets, all subjects were augmented with motion. For our experiments, we elect to use all of the subjects in our training and train to the very last epoch. A variety of appropriately titled pre-trained models corresponding to results in the main paper and the appendices are included alongside our code in the *pretrained_models* folder.

D. Motion Augmented rPPG Videos

We provide code for augmenting various camera physiology datasets and pre-trained models trained on motion-augmented data. Additionally, we provide various files to easily train on baselines and motion augmented data using the publicly available rPPG-Toolbox [20]. Pre-trained

Table 10: **Effect of Motion Types – Rigid.** We augment UBFC-rPPG with various types of rigid head motions and test on AFRL [7]. The best results are shown in bold.

Training Set	Rigid Motion	No Motion				Small Motion				Large Motion				All			
		MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$
UBFC-rPPG	Very Small	1.00	3.86	1.48	0.95	2.28	6.36	3.44	0.85	7.59	12.91	10.99	0.49	4.72	10.01	6.59	0.67
MAUBFC-rPPG	Small	0.84	3.25	1.18	0.96	1.44	4.44	2.03	0.93	4.21	9.11	5.96	0.74	3.19	7.96	4.36	0.79
MAUBFC-rPPG	Large	1.00	3.61	1.37	0.96	1.78	5.23	2.49	0.90	3.64	8.14	5.12	0.78	3.39	8.26	4.58	0.77
OURS vs. BASELINE		+16.00%	+15.80%	+20.27%	+1.05%	+36.84%	+30.19%	+40.99%	+9.41%	+52.04%	+36.95%	+53.41%	+59.18%	+32.42%	+20.48%	+33.84%	+17.91%

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation

Table 11: **Evaluation across all datasets.** We motion-augment two training datasets, UBFC-rPPG and PURE, to create MAUBFC-rPPG and MAPURE, respectively. We observe that the motion-augmented versions produce significant improvements (shown in bold).

Training Set	Method	UBFC-rPPG				PURE				Testing Set UBFC-PHYS				AFRL				MMPD			
		MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$
Unsupervised	Green	19.82	31.49	18.78	0.37	10.09	23.85	10.28	0.34	13.45	19.11	16.00	0.31	7.01	12.52	9.24	0.52	16.27	21.74	20.09	-0.04
	ICA	14.70	23.71	14.34	0.53	4.77	16.70	4.47	0.72	8.00	13.51	9.48	0.48	6.77	12.25	8.96	0.51	13.10	17.84	16.33	0.03
	CHROM	3.98	8.72	3.78	0.89	5.77	14.93	11.52	0.81	4.68	8.09	6.20	0.77	5.41	10.71	7.95	0.60	8.85	12.77	11.93	0.29
	POS	4.00	7.58	3.86	0.92	3.67	11.82	7.25	0.88	4.62	8.02	6.29	0.78	6.93	11.89	10.00	0.49	8.18	13.04	11.12	0.31
UBFC-rPPG	TS-CAN	-	-	-	-	4.55	14.47	4.67	0.80	5.56	9.88	7.25	0.68	4.24	8.72	5.84	0.75	8.74	15.55	10.51	0.25
MAUBFC-rPPG	TS-CAN	-	-	-	-	1.14	5.95	1.30	0.97	3.93	7.50	5.24	0.81	2.67	6.55	3.65	0.85	6.80	14.20	7.97	0.29
PURE	TS-CAN	1.34	3.01	1.55	0.99	-	-	-	-	4.43	8.12	5.89	0.78	2.63	7.35	3.51	0.82	8.96	16.59	10.33	0.15
MAPURE	TS-CAN	1.03	2.70	1.17	0.99	-	-	-	-	4.39	8.10	5.90	0.78	2.37	6.28	3.26	0.87	8.08	15.38	9.54	0.18
MAUBFC-rPPG vs. UBFC-rPPG		-	-	-	-	+74.95%	+58.88%	+72.16%	+21.25%	+29.32%	+24.09%	+22.72%	+19.12%	+37.03%	+24.89%	+37.50%	+13.33%	+22.20%	+8.68%	+24.17%	+16.00%
MAPURE vs. PURE		+23.13%	+10.30%	+24.52%	+0.00%	-	-	-	-	+0.90%	+0.25%	-0.17%	+0.00%	+0.89%	+14.56%	+7.12%	+6.10%	+9.82%	+7.29%	+7.65%	+20.00%

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation

Table 12: **Effect of Multiple Augmentations.** Augmenting each source video of UBFC-rPPG 1x, 2x, 3x, and 4x, we test on PURE and UBFC-PHYS datasets. The best results are shown in bold.

Training Set	Size	Subjects	PURE				Testing Set UBFC-PHYS				AFRL				MMPD			
			MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$
UBFC-rPPG	42	42	4.55	14.47	4.67	0.80	5.56	9.88	7.25	0.68	4.24	8.72	5.84	0.75	8.74	15.55	10.51	0.25
MAUBFC-rPPG	42	42	1.14	5.95	1.30	0.97	3.93	7.50	5.24	0.81	2.67	6.55	3.65	0.85	6.80	14.20	7.97	0.29
MAUBFC-rPPG 2x	84	42	1.12	5.95	1.29	0.97	3.90	7.42	5.22	0.81	-	-	-	-	-	-	-	-
MAUBFC-rPPG 3x	126	42	1.11	5.95	1.27	0.97	3.97	7.80	5.31	0.81	-	-	-	-	-	-	-	-
MAUBFC-rPPG 4x	168	42	1.19	6.05	1.33	0.97	4.10	7.95	5.40	0.80	-	-	-	-	-	-	-	-
OURS vs. BASELINE			+2.63%	+0.00%	+2.31%	+0.00%	+0.76%	+1.07%	+0.38%	+0.00%	-	-	-	-	-	-	-	-

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation

Table 13: **Synthetic vs Naturalistic Head Motion.** We compare the effect of adding head motions to SCAMPS and UBFC-rPPG and contrast this with using motion data in SCAMPS. Average time for augmenting each frame of a sequence is presented. The best results are shown in bold.

Training Set	Testing Set				AFRL				Per Frame			
	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$	Synthesis Time			
SCAMPS-200 (No motion)	10.29	23.81	11.09	0.35	7.75	13.08	10.54	0.48	37s			
SCAMPS-200 (Motion)	5.38	16.98	5.42	0.72	7.25	12.85	10.20	0.48	37s			
UBFC-rPPG	4.55	14.47	4.67	0.80	4.72	10.01	6.59	0.67	-			
MASCAMPS-200	4.67	16.35	4.22	0.75	5.00	10.10	6.69	0.67	1.20s			
MAUBFC-rPPG	1.14	5.95	1.30	0.97	3.24	7.89	4.37	0.79	2.39s			
MASCAMPS vs. SCAMPS BASELINE	+13.20%	+3.71%	+22.14%	+4.17%	+31.03%	+21.40%	+34.41%	+39.58%	+96.76%			
MAUBFC vs. UBFC-rPPG BASELINE	+74.4%	+58.88%	+72.16%	+21.25%	+31.36%	+21.18%	+33.69%	+17.91%	-			

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation, Synthesis Time = the amount of time (in seconds) it takes to synthesize a single frame, when relevant

Table 14: **Effect of PPG Estimation Models.** We train different PPG estimation networks on UBFC-rPPG and MAUBFC-rPPG and evaluate on PURE. The best results are shown in bold.

Training Set	Method	Testing Set PURE			
		MAE↓	RMSE↓	MAPE↓	$\rho \uparrow$
UBFC-rPPG	DeepPhys	5.14	17.20	4.90	0.72
MAUBFC-rPPG	DeepPhys	1.24	6.01	1.56	0.97
UBFC-rPPG	EfficientPhys	4.95	16.28	4.56	0.74
MAUBFC-rPPG	EfficientPhys	1.45	6.07	1.76	0.97
UBFC-rPPG	TS-CAN	4.55	14.47	4.67	0.80
MAUBFC-rPPG	TS-CAN	1.14	5.95	1.30	0.97

MAE = Mean Absolute Error in HR estimation (Beats/Min), RMSE = Root Mean Square Error in HR estimation (Beats/Min), MAPE = Mean Absolute Percentage Error in HR estimation, ρ = Pearson Correlation in HR estimation

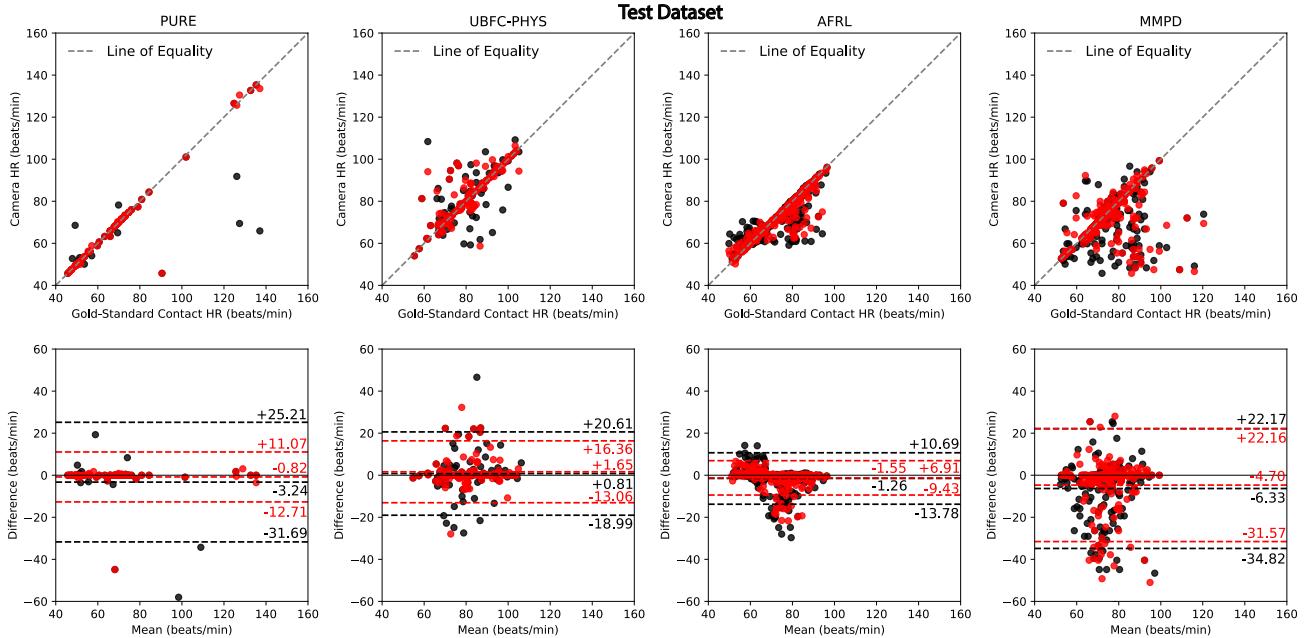


Figure 4: **Scatter and Bland-Altman Plots.** Scatter (top row) and Bland-Altman (bottom row) plots for models trained on UBFC-rPPG (black) and MAUBFC-rPPG (red) and tested on (from left to right), PURE, UBFC-PHYS, AFRL, and MMPD.

models using the baseline UBFC-rPPG or PURE datasets can be found in the rPPG-Toolbox. We also include motion analysis scripts that utilize OpenFace [1] to analyze both rigid and non-rigid motion in videos and generate plots. All of these additional materials be found through our project page: <https://motion-matters.github.io/>.

D.1. The Effect of Motion Transfer on PPG

In Figure 6, we provide additional qualitative examples of the effect of motion transfer on the underlying PPG signal in rPPG videos. All examples include a plot of the gold-standard PPG signal, the predicted PPG signal from the unaugmented source rPPG video by a TS-CAN model,

and the predicted PPG signal from the motion-augmented source rPPG video by a TS-CAN model. The TS-CAN models utilized are trained on a larger superset of the shown examples (e.g., UBFC, MAUBFC) with the same experimental settings mentioned in Section C.1. As shown by these qualitative examples, a neural method such as TS-CAN is capable of recovering the underlying PPG signal despite the application of motion augmentation.

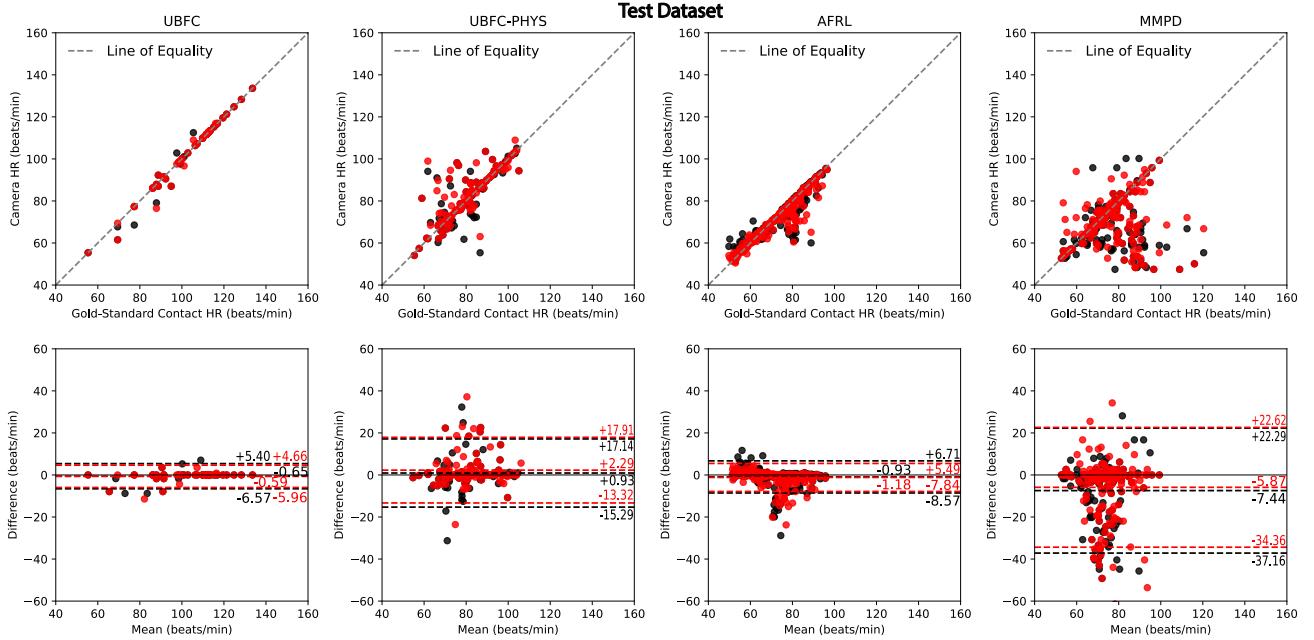


Figure 5: **Scatter and Bland-Altman Plots.** Scatter (top row) and Bland-Altman (bottom row) plots for models trained on PURE (black) and MAPURE (red) and tested on (from left to right), UBFC-rPPG, UBFC-PHYS, AFRL, and MMPD.

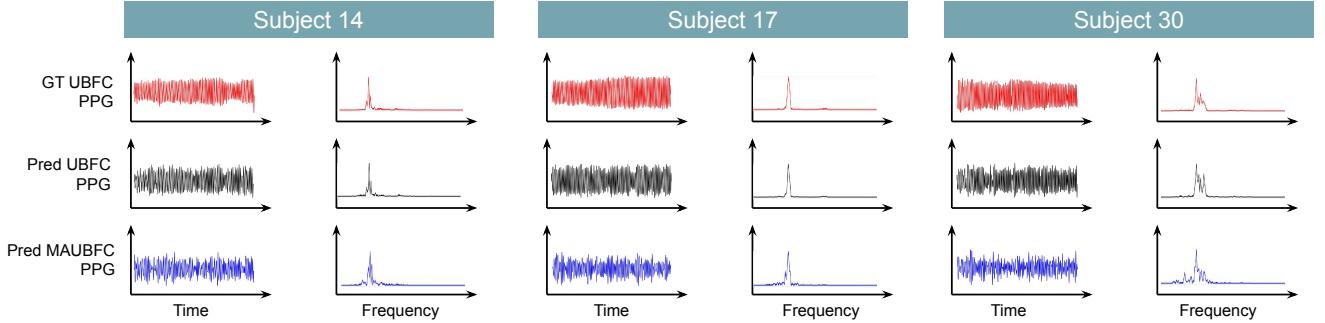


Figure 6: **Signal Prediction on UBFC-rPPG and MAUBFC-rPPG.** We provide three, subject-wise examples of signal prediction using a neural method, TS-CAN, on unaugmented and motion-augmented videos from the UBFC-rPPG dataset.

E. Datasets

E.1. Source Videos for Motion Synthesis

We use the following state-of-the-art datasets for source videos used in our motion synthesis pipeline:

UBFC-rPPG [3]: The UBFC-rPPG video dataset contains 42 RGB videos, one per subject, at 30 Hz. The videos were collected with a Logitech C920 HD Pro with a resolution of 640x480 and a CMS50E transmissive pulse oximeter was utilized in order to record gold-standard PPG signals. The UBFC-rPPG dataset contains minimal motion, with subjects being asked to simply sit one meter away from the camera in an environment with both artificial and natural lighting. When utilized as source videos, we utilized

all videos from the UBFC-rPPG dataset. When utilized for evaluation, we also utilized all videos from the UBFC-rPPG dataset.

PURE [39]: The PURE dataset contains 59 videos, each corresponding to a unique task, per a subject. The six tasks involve staying steady, talking, slow head translation, fast head translation, small head rotation, and medium head rotation. There are 10 subjects total with subject 6’s talking task video being excluded. All of the videos were captured with an RGB eco274CVGE camera (SVS-Vistek GmbH) at a resolution of 640x480 and 60 Hz. During all tasks, the subject was asked to be seated in front of the camera at an average distance of 1.1 meters and lit from the front with ambient natural light through a window. A

gold-standard measure of PPG was collected with a pulse oximeter, CMS50E, attached to the finger. When utilized as source videos, we utilized all videos from the PURE dataset. When utilized for evaluation, we also utilized all videos from the PURE dataset.

SCAMPS [27]: The SCAMPS dataset contains 2,800 synthetic videos that were generated using a blendshape-based rig with 7,667 vertices and 7,414 polygons, with distinct identities being learned from a set of high-quality facial scans. Blood flow, and subsequently the underlying physiological signals, are simulated using the modification of physically-based shading materials. The SCAMPS dataset contains a variety of rigid and non-rigid head motions, with varying intensities. The dataset also contains a variety of lighting conditions and background conditions. Each SCAMPS video is 20 seconds in length, with 600 frames at a sampling rate of 30 Hz. We only utilize portions of the SCAMPS dataset as source videos in our ablation study regarding synthetic versus naturalistic head motion.

E.2. Driving Videos for Motion Synthesis

We use the following datasets for driving videos used in our motion synthesis pipeline:

TalkingHead-1KH [45]: The TalkingHead-1KH dataset is a publicly available, large-scale talking-head video dataset used as a benchmark for Face-Vid2Vid [45] and entirely sourced from YouTube videos. It contains 180K unconstrained videos of people speaking in a variety of real-world contexts, leading to a rich diversity in both rigid and non-rigid motion. The videos are of varied resolutions, but there is an emphasis on collecting high quality, high resolution videos which compose a significant portion of the dataset (with a resolution of at least 512x512). We elect to filter the dataset for head pose such that any videos where the head pose, on average, is outside +/- 20 degrees are removed. This prevents damaging motion augmentation artifacts due to impractical differences in the head pose in the source video and the head pose in the driving videos, but comes at the cost of reduced head pose variations. We also filter by facial action units (AUs) (0 to 5, in units of intensity) such that any videos below a mean standard deviation in facial AUs of 0.15 is filtered out. This prevents driving videos that are not suitable for our application from being used - for example, a driving video that is effectively a slide show and doesn't have any naturalistic motion upon qualitative inspection.

CDVS: The CDVS contains 90 self-captured videos by 5 subjects with heavily constrained, unnatural motion used only for ablation studies to understand the impact of augmenting data with various degrees of rigid and non-rigid motion. Subjects self-capture the videos in a variety of settings with artificial lighting of the face in an indoors setting. When capturing a video to show one of the two types of mo-

tion we study, subjects are asked to constrain the other motion type as much as possible. The CDVS will be released in the future for research purposes.

E.3. Additional Datasets for Evaluation

In addition to using UBFC-rPPG [3] and PURE [39] as both source video datasets in the motion synthesis pipeline and datasets for evaluation, we use three additional state-of-the-art datasets for evaluation:

UBFC-PHYS [28]: The UBFC-PHYS dataset is a multimodal dataset with 168 RGB videos, with 56 subjects (46 women and 10 men) per a task. There are three tasks with significant amounts of both rigid and non-rigid motion - a rest task, a speech task, and an arithmetic task. Gold-standard BVP and electrodermal activity (EDA) measurements were collected via the Empatica E4 wristband. The videos were recorded at a resolution of 1024x1024 and 35Hz with a EO-23121C RGB digital camera. We utilized all of the tasks and the same subject sub-selection list provided by the authors of the dataset in the second supplementary material of Sabour et al. [28] for evaluation. This means we eliminated 14 subjects (s3, s8, s9, s26, s28, s30, s31, s32, s33, s40, s52, s53, s54, s56) for the rest task, 30 subjects (s1, s4, s6, s8, s9, s11, s12, s13, s14, s19, s21, s22, s25, s26, s27, s28, s31, s32, s33, s35, s38, s39, s41, s42, s45, s47, s48, s52, s53, s55) for the speech task, and 23 subjects (s5, s8, s9, s10, s13, s14, s17, s22, s25, s26, s28, s30, s32, s33, s35, s37, s40, s47, s48, s49, s50, s52, s53) for the arithmetic task.

AFRL [7]: The AFRL dataset contains 300 videos of 25 participants (17 males, 8 females) recorded at 658x492 resolution and 120 FPS. Gold-standard physiological signals were measured using the fingertip reflectance PPG method. Participants were asked to perform a series of tasks, resulting in 12 tasks total. With a black background behind the participant, the tasks entailed sitting still with a chin-rest, sitting still without a chin rest, rotating the head with an angular velocity of 10 degrees/second, 20 degrees/second, and 30 degrees/second, and finally randomly orienting their head once per a second to a predefined location. This resulted in six recordings, which were repeated once with a colorful background, resulting in 12 videos per a participant. As a part of our pre-processing steps for AFRL, we down-sampled the videos to 30 FPS. We utilized all of the videos for evaluation.

MMPD [42]: The Multi-domain Mobile Video Physiology Dataset (MMPD) dataset contains 11 hours of recordings from mobile phones of 33 subjects. Gold-standard PPG signals were simultaneously recorded using an HKG-07C+ oximeter. The dataset was designed to capture variations in skin tone, body motion, and lighting conditions in videos useful for the rPPG task. Videos were collected under three artificial light sources: i) low LED light (100 lumens on the

face region), ii) mid-level incandescent light (200 lumens on the face region), and iii) high LED light (300 lumens on the face region). Videos were also collected under natural light, which varied from 300-800 lumens intensity on the face region. Videos were recorded following an experimental procedure in which participants performed a variety of tasks in different lighting conditions - a stationary task, a head rotation task, a talking task, and a walking task. We evaluated on videos with artificial lighting, Fitzpatrick scale skin tone type 3, and any of the four tasks (stationary, head rotation, talking, and walking) that correspond to varying degrees of rigid and non-rigid motion.

F. Broader Impact Statement

While generating synthetic videos that are indistinguishable from those of real people has concerning use cases, there are positive applications of this technology can enabled. In the medical domain simulators are increasingly being tested within specific applications [11, 10]. It is important that the limitations of generative models are understood as these may impact the performance of the resulting models trained using simulated data. It is possible for generative approaches to compound harmful biases [22] and motion augmentation algorithms can be used for troubling negative applications. To mitigate negative outcomes, we license our source code using responsible behavioral use licenses used across a large number of publicly released machine learned models [5].

References

- [1] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. Openface: an open source facial behavior analysis toolkit. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, pages 1–10. IEEE, 2016.
- [2] Vladimir Blazek, Ting Wu, and Dominik Hoelscher. Near-infrared ccd imaging: Possibilities for noninvasive and contactless 2d mapping of dermal venous hemodynamics. In *Optical Diagnostics of Biological Fluids V*, volume 3923, pages 2–9. International Society for Optics and Photonics, 2000.
- [3] Serge Bobbia, Richard Macwan, Yannick Benzeth, Alamin Mansouri, and Julien Dubois. Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters*, 124:82–90, 2019.
- [4] Weixuan Chen and Daniel McDuff. Deepphys: Video-based physiological measurement using convolutional attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 349–365, 2018.
- [5] Danish Contractor, Daniel McDuff, Julia Katherine Haines, Jenny Lee, Christopher Hines, Brent Hecht, Nicholas Vincent, and Hanlin Li. Behavioral use licensing for responsible ai. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 778–788, 2022.
- [6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- [7] Justin R Estepp, Ethan B Blackford, and Christopher M Meier. Recovering pulse rate during motion artifact with a multi-imager array for non-contact imaging photoplethysmography. In *Systems, Man and Cybernetics (SMC), 2014 IEEE International Conference on*, pages 1462–1469. IEEE, 2014.
- [8] John Gideon and Simon Stent. The way to my heart is through contrastive learning: Remote photoplethysmography from unlabelled video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3995–4004, 2021.
- [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [10] Waldemar Hahn, Katharina Schütte, Kristian Schultz, Olaf Wolkenhauer, Martin Sedlmayr, Ulrich Schuler, Martin Eichler, Saptarshi Bej, and Markus Wolfien. Contribution of synthetic data generation towards an improved patient stratification in palliative care. *Journal of Personalized Medicine*, 12(8):1278, 2022.
- [11] Mikel Hernandez, Gorka Epelde, Ane Alberdi, Rodrigo Cilla, and Debbie Rankin. Synthetic data generation for tabular health records: A systematic review. *Neurocomputing*, 2022.
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [13] Fa-Ting Hong, Longhao Zhang, Li Shen, and Dan Xu. Depth-aware generative adversarial network for talking head video generation. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3387–3396, 2022.
- [14] In Cheol Jeong and Joseph Finkelstein. Introducing contactless blood pressure assessment using a high speed video camera. *Journal of medical systems*, 40(4):77, 2016.
- [15] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.
- [16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [17] Ming-Yu Liu, Xun Huang, Jiahui Yu, Ting-Chun Wang, and Arun Mallya. Generative adversarial networks for image and video synthesis: Algorithms and applications. *Proceedings of the IEEE*, 109:839–862, 2020.
- [18] Xin Liu, Josh Fromm, Shwetak Patel, and Daniel McDuff. Multi-task temporal shift attention networks for on-device contactless vitals measurement. *NeurIPS*, 2020.
- [19] Xin Liu, Brian L Hill, Ziheng Jiang, Shwetak Patel, and Daniel McDuff. Efficientphys: Enabling simple, fast and accurate camera-based vitals measurement. *arXiv preprint arXiv:2110.04447*, 2021.
- [20] Xin Liu, Xiaoyu Zhang, Girish Narayanswamy, Yuzhe Zhang, Yuntao Wang, Shwetak Patel, and Daniel McDuff. Deep physiological sensing toolbox. *arXiv preprint arXiv:2210.00716*, 2022.

- [21] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101, 2017.
- [22] Vongani H. Maluleke, Neerja Thakkar, Tim Brooks, Ethan Weber, Trevor Darrell, Alexei A. Efros, Angjoo Kanazawa, and Devin Guillory. Studying bias in gans through the lens of race, 2022.
- [23] Daniel McDuff. Camera measurement of physiological vital signs. *ACM Computing Surveys (CSUR)*, 2021.
- [24] Daniel McDuff, Roger Cheng, and Ashish Kapoor. Identifying bias in ai using simulation. 2018.
- [25] Daniel McDuff, Xin Liu, Javier Hernandez, Erroll Wood, and Tadas Baltrusaitis. Synthetic data for multi-parameter camera-based physiological sensing. In *2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2021.
- [26] Daniel McDuff, Shuang Ma, Yale Song, and Ashish Kapoor. Characterizing bias in classifiers using generative models. *Advances in Neural Information Processing Systems*, 32:5403–5414, 2019.
- [27] Daniel McDuff, Miah Wander, Xin Liu, Brian L Hill, Javier Hernandez, Jonathan Lester, and Tadas Baltrusaitis. Scamps: Synthetics for camera measurement of physiological signals. *arXiv preprint arXiv:2206.04197*, 2022.
- [28] Rita Mezatabisabour, Yannick Benzezeth, Pierre De Oliveira, Julien Chappe, and Fan Yang. Ubfc-phys: A multimodal database for psychophysiological studies of social stress. *IEEE Transactions on Affective Computing*, 2021.
- [29] Xuesong Niu, Hu Han, Shiguang Shan, and Xilin Chen. Vipl-hr: A multi-modal database for pulse estimation from less-constrained face video. *arXiv preprint arXiv:1810.04927*, 2018.
- [30] Ewa Nowara, Daniel McDuff, Ashutosh Sabharwal, and Ashok Veeraraghavan. Seeing beneath the skin with computational photography. *Communications of the ACM*, 65(12):90–100, 2022.
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [32] Ming-Zher Poh, Daniel McDuff, and Rosalind W Picard. Advancements in noncontact, multiparameter physiological measurements using a webcam. *IEEE transactions on biomedical engineering*, 58(1):7–11, 2010.
- [33] Ming-Zher Poh, Yukkee Cheung Poh, Pak-Hei Chan, Chun-Ka Wong, Louise Pun, Wangie Wan-Chiu Leung, Yu-Fai Wong, Michelle Man-Ying Wong, Daniel Wai-Sing Chu, and Chung-Wah Siu. Diagnostic assessment of a deep learning system for detecting atrial fibrillation in pulse waveforms. *Heart*, 104(23):1921–1928, 2018.
- [34] Tong Sha, Wei Zhang, Tong Shen, Zhoujun Li, and Tao Mei. Deep person generation: A survey from the perspective of face, pose and cloth synthesis. *ACM Computing Surveys*, 2021.
- [35] Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard Moore, Alex Kipman, and Andrew Blake. Real-time human pose recognition in parts from single depth images. In *CVPR 2011*, pages 1297–1304. Ieee, 2011.
- [36] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *Conference on Neural Information Processing Systems (NeurIPS)*, December 2019.
- [37] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [38] Radim Špetlík, Vojtech Franc, and Jirí Matas. Visual heart rate estimation with convolutional neural network. In *Proceedings of the british machine vision conference, Newcastle, UK*, pages 3–6, 2018.
- [39] Ronny Stricker, Steffen Müller, and Horst-Michael Gross. Non-contact video-based pulse rate measurement on a mobile service robot. In *The 23rd IEEE International Symposium on Robot and Human Interactive Communication*, pages 1056–1062. IEEE, 2014.
- [40] Lech Świrski and Neil Dodgson. Rendering synthetic ground truth images for eye tracker evaluation. In *Proceedings of the Symposium on Eye Tracking Research and Applications*, pages 219–222, 2014.
- [41] Chihiro Takano and Yuji Ohta. Heart rate measurement based on a time-lapse image. *Medical engineering & physics*, 29(8):853–857, 2007.
- [42] Jiankai Tang, Kequan Chen, Yuntao Wang, Yuanchun Shi, Shwetak Patel, Daniel McDuff, and Xin Liu. Mmpd: Multi-domain mobile video physiology dataset, 2023.
- [43] Wim Verkruyse, Lars O Svaasand, and J Stuart Nelson. Remote plethysmographic imaging using ambient light. *Optics express*, 16(26):21434–21445, 2008.
- [44] Hao Wang, Euijoon Ahn, and Jinman Kim. Self-supervised representation learning framework for remote physiological measurement using spatiotemporal augmentation loss. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2431–2439, 2022.
- [45] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10034–10044, 2020.
- [46] Wenjin Wang, Albertus C den Brinker, Sander Stuijk, and Gerard de Haan. Algorithmic principles of remote ppg. *IEEE Transactions on Biomedical Engineering*, 64(7):1479–1491, 2017.
- [47] Zhen Wang, Yunhao Ba, Pradyumna Chari, Oyku Deniz Bozkurt, Gianna Brown, Parth Patwa, Niranjan Vaddi, Laleh Jalilian, and Achuta Kadambi. Synthetic generation of face videos with plethysmograph physiology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20587–20596, 2022.
- [48] Erroll Wood, Tadas Baltrusaitis, Charlie Hewitt, Sebastian Dziadzio, Thomas J Cashman, and Jamie Shotton. Fake it till you make it: Face analysis in the wild using synthetic data alone. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3681–3691, 2021.

- [49] Erroll Wood, Tadas Baltrusaitis, Xucong Zhang, Yusuke Sugano, Peter Robinson, and Andreas Bulling. Rendering of eyes for eye-shape registration and gaze estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3756–3764, 2015.
- [50] Yuzhe Yang, Xin Liu, Jiang Wu, Silviu Borac, Dina Katabi, Ming-Zher Poh, and Daniel McDuff. Simper: Simple self-supervised learning of periodic targets. *arXiv preprint arXiv:2210.03115*, 2022.
- [51] Zitong Yu, Wei Peng, Xiaobai Li, Xiaopeng Hong, and Guoying Zhao. Remote heart rate measurement from highly compressed facial videos: an end-to-end deep learning solution with video enhancement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 151–160, 2019.
- [52] Zitong Yu, Yuming Shen, Jingang Shi, Hengshuang Zhao, Philip Torr, and Guoying Zhao. Physformer: Facial video-based physiological measurement with temporal difference transformer. *arXiv preprint arXiv:2111.12082*, 2021.
- [53] 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.