

# Motion Matters: Neural Motion Transfer for Better Camera Physiological Measurement

Akshay Paruchuri<sup>1</sup>, Xin Liu<sup>2</sup>, Yulu Pan<sup>1</sup>, Shwetak Patel<sup>2</sup>, Daniel McDuff<sup>2,\*</sup>, Soumyadip Sengupta<sup>1,\*</sup>

<sup>1</sup>UNC Chapel Hill <sup>2</sup>University of Washington

{akshay, ronisen, yulupan}@cs.unc.edu, {xliu0, shwetak, dmcduff}@cs.washington.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 of interest. We adapt a neural video synthesis approach to augment videos for the task of remote photoplethysmography (rPPG) 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, we demonstrate a 47% improvement over existing inter-dataset results using various state-of-the-art methods on the PURE dataset. We also present inter-dataset results on five benchmark datasets to show improvements of up to 79% using TS-CAN, a neural rPPG estimation method. Our findings illustrate the usefulness of motion transfer as a data augmentation technique for improving the generalization of models for camera-based physiological sensing. We release our code for using motion transfer as a data augmentation technique on three publicly available datasets, UBFC-rPPG, PURE, and SCAMPS, and models pre-trained on motion-augmented data here: <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 [28]. Camera-based measurement of cardiac signals is one such application [20], in which cameras are used to measure the pulse via light re-

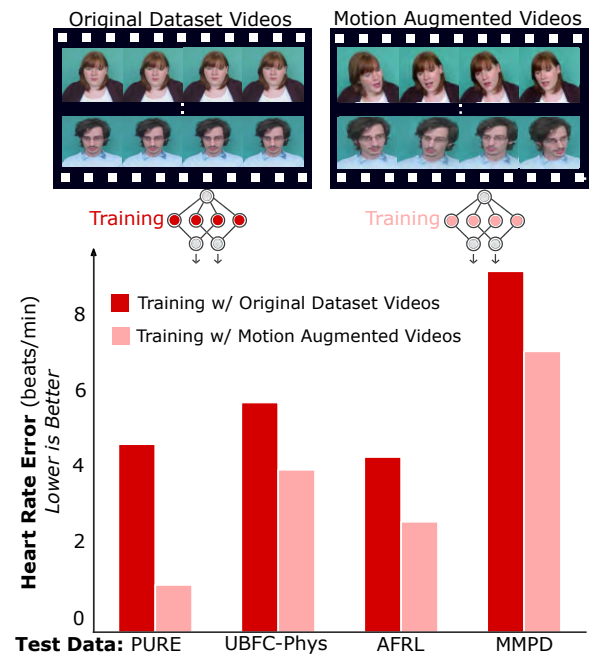


Figure 1. **Motion augmentation improves rPPG.** We present the first neural motion augmentation pipeline for the task of remote PPG estimation and empirically show it reduces error in heart rate estimation by up to 79% in inter-dataset results using TS-CAN and 47% over existing results using SOTA methods on PURE.

flected from the body, a principle known as photoplethysmography (PPG) [2, 42]. The PPG signals can be used to derive respiration [30], heart rate variability [30], arrhythmia [31], and blood pressure [12]. As a result this technology has the potential to turn webcams and smartphones into meaningful health sensors.

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, 15, 37, 50, 51]. 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 [38]) serve as an extremely valuable resource for the research

\*denotes equal advising

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 [21, 23, 47], body pose estimation [33] and eye tracking [39, 48].

However, creating synthetic data that preserve the subtle and nuanced peripheral pulse in a video is non-trivial. McDuff et al. [24] released a large dataset (2,800 videos) of avatars and cardiac signals; however, their computer graphics pipeline had an extremely high computational overhead. Wang et al. [46] used a learning based method to generate synthetic videos given a reference image and target PPG signal. Their creative approach successfully incorporated PPG signals to produce 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:

- A systematic investigation of the impact of motion augmentation on the physiological information within rPPG videos.
- Quantitative, empirical evidence that conveys the meaningfulness of training with motion-augmented data, including (1) the benefits of different kinds of motion-augmented data, (2) consistent motion augmentations across neural motion transfer algorithms, (3) the benefit of naturalistic head motion over other synthetic methods (e.g., SCAMPS), and (4) consistent improvements using motion-augmented data regardless of a chosen rPPG estimation model.
- Our motion augmentation pipeline, using which in Table 1 we demonstrate that our approach surpasses the SOTA when compared to other methods that train on UBFC-rPPG [3] and test on PURE [38]. We also provide comprehensive inter-dataset results (see Table 2) that highlight the usefulness of motion augmentation for improving the generalization of models for camera-based physiological sensing.

We summarize the key findings of this paper about the effectiveness of motion transfer as a data augmentation tool in Section 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 [6, 9, 10, 13, 35] 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 [14, 32]. 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) [11, 34, 44, 52]. 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 [11, 34, 44]. 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 [11, 34, 52]. Face-Vid2Vid [44] 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 [40, 42]. The application of more advanced signal processing methods helped make measurement somewhat more robust under real-world conditions [30, 45], as did leveraging knowledge of physiological and physical properties [45]. Yet, these models were still very sensitive to body motions. Both task-specific and multi-task neural, data-driven models currently achieve SOTA results in most cases [4, 15, 17, 26, 50, 51], 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, 36, 43, 49].

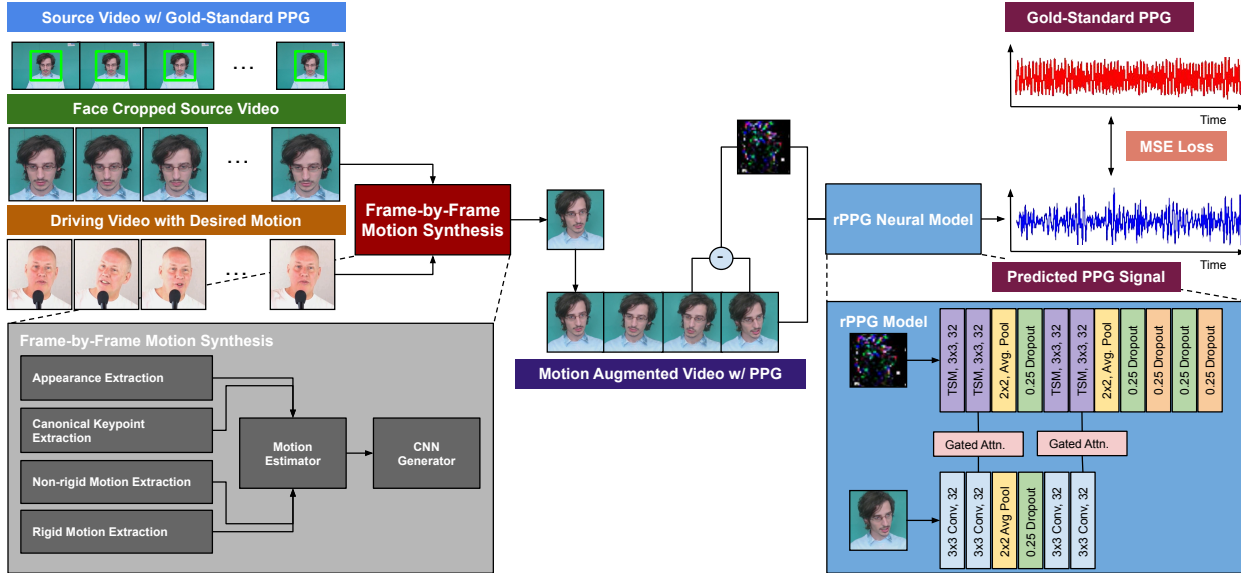


Figure 2. **Motion augmentation and training pipeline.** We augment frames of a source video with corresponding frames of a selected driving video to create an augmented video with the identity of the source video and motion of the target driving video. We then train a remote PPG estimation network on the augmented video with a mean squared error (MSE) loss.

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 [22, 46].

**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 [38], VIPL-HR [27]) are valuable resources. However, given the challenging nature of the rPPG task researchers have collected and released data under heavily constrained conditions with very little physical motion. More recent datasets (such as UBFC-PHYS [25] and MMPD [41]) 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. We provide additional quantitative evidence to highlight preservation of the underlying physiological signal through the signal-to-noise ratio (SNR) metric and after rPPG signal extraction using TS-CAN in Table 6.

#### 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 associated with a PPG signal label. 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 a PPG signal

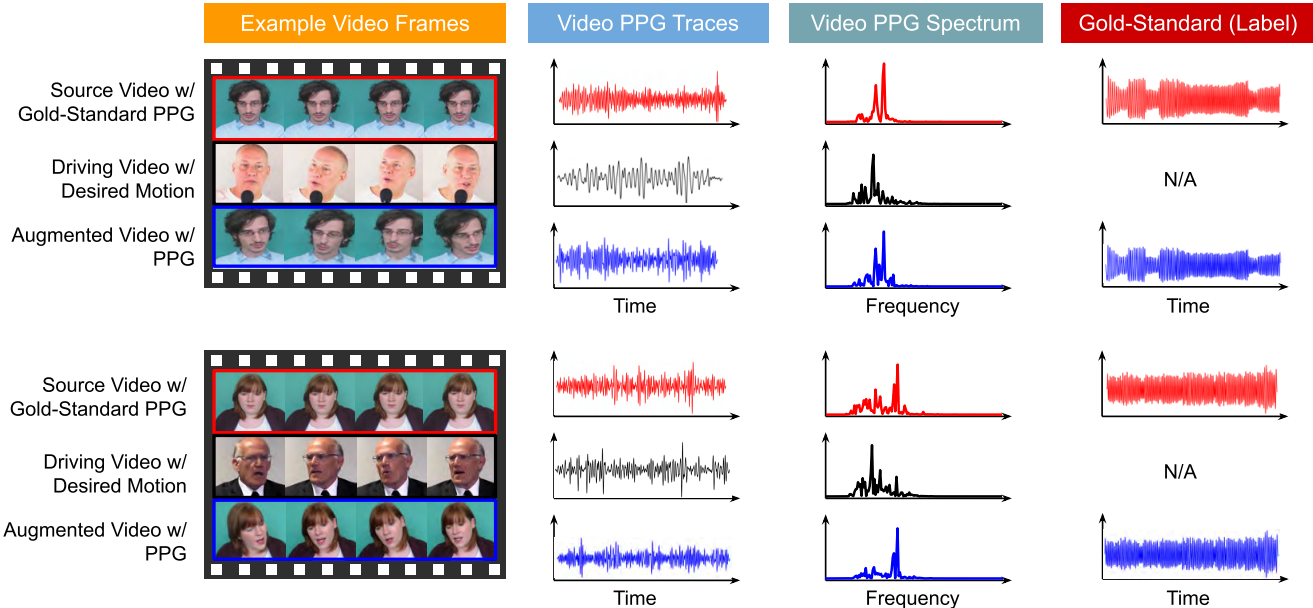


Figure 3. **Preserving physiological signals in motion augmented videos.** We show that applying neural motion transfer preserve the physiological signal corresponding to the heart-rate present in the peak of the frequency spectrum of the source and the augmented video.

label 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  $y_t = \mathbf{M}(s_t, d_t; \theta)$ .

We choose Face-Vid2Vid [44], 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 image using a driving video. In contrast, we re-purpose the same core 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 and is further described in our supplementary materials, alongside provided code that will be open-sourced.

**Source Video Datasets:** We utilize the UBFC-rPPG [3] and PURE [38] 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 [44] 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 [44] 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 [15] to predict the 1st-order derivative of the PPG signal after training on videos augmented with motion. We also use DeepPhys [4] and PhysNet [50] 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 [44] 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 [45], 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 the augmented video. This appears to also hold true across different appearances and motion conditions, both in the source videos and the driving videos. Again, we also provide quantitative evidence to support our observation in Table 6. We also present additional qualitative results and information about prior works analyzing the presence of physiological signals in deep fake videos in the supplementary material. Thus, we can effectively claim that keypoint-based motion transfer algorithms like Face-Vid2Vid [44] 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 rPPG 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** [38], **UBFC-PHYS** [25], **AFRL** [7], and **MMPD** [41] (see supplementary materials for more details). 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 rPPG estimation to date, testing on five independent test datasets.

tensive inter-dataset evaluation of rPPG 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. All networks were trained using an NVIDIA RTX A4500 and PyTorch [29] implementations in a publicly available toolbox for the rPPG task [18]. 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 1, we compare our approach with TS-CAN and motion-augmented source data to other SOTA methods using the same source data, the UBFC-rPPG dataset. The PhysNet [50] result was obtained from [36], and differs from our reproduced PhysNet result in Table 7 due to pre-processing and implementation differences. We achieve a 47% improvement over SOTA results on the PURE dataset with our data augmentation strategy using neural motion transfer. In Table 2, we comprehensively compare the performance of a supervised PPG estimation network, TS-CAN [15], 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 aforementioned tables only show limited metrics such as the MAE or MAPE in heart rate estimation. Equivalent tables with additional metrics, including root mean squared error (RMSE) and Pearson correlation metrics can be found where applicable in the supplementary material. The driving videos used for augmentation in Table 2 contain significant amounts of unconstrained motion – both rigid and non-rigid.

We observe that training TS-CAN on augmented videos produces 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 79% 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 [44] dataset as our driving videos for augmenting various rPPG video datasets with motion. We analyze the videos using OpenFace [1] to obtain the inten-

Method	MAE↓
EfficientPhys-C [16]	5.47
SiNC [36]	4.02
PhysNet [50]	3.81
Physformer [51]	1.99
Dual-GAN [19]	1.81
Ours (Motion Augmented)	<b>0.96</b>
Ours vs. Best Baseline	<b>+46.96%</b>

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

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 UBFC-PHYS		AFRL		MMPD	
		MAE↓	MAPE↓	MAE↓	MAPE↓	MAE↓	MAPE↓	MAE↓	MAPE↓	MAE↓	MAPE↓
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	-	-	<b>0.96</b>	<b>1.13</b>	<b>3.93</b>	<b>5.24</b>	2.67	3.65	<b>6.80</b>	<b>7.97</b>
PURE	TS-CAN	1.34	1.55	-	-	4.43	5.89	2.63	3.51	8.96	10.33
MAPURE	TS-CAN	<b>1.03</b>	<b>1.17</b>	-	-	4.39	5.90	<b>2.37</b>	<b>3.26</b>	8.08	9.54
MAUBFC-rPPG vs. UBFC-rPPG		-	-	<b>+78.90%</b>	<b>+75.08%</b>	<b>+29.32%</b>	<b>+27.72%</b>	<b>+37.03%</b>	<b>+37.50%</b>	<b>+22.20%</b>	<b>+24.17%</b>
MAPURE vs. PURE		<b>+23.13%</b>	<b>+24.52%</b>	-	-	<b>+0.90%</b>	<b>-0.17%</b>	<b>+9.89%</b>	<b>+7.12%</b>	<b>+9.82%</b>	<b>+7.65%</b>

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

sity (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 3.

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 [38] 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 present these results in Table 4. 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%.

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
		MAE↓	MAE↓	MAE↓	MAE↓
UBFC-rPPG	Very Small	1.00	2.28	7.59	4.72
MAUBFC-rPPG	Small	<b>0.84</b>	<b>1.44</b>	4.21	<b>3.19</b>
MAUBFC-rPPG	Large	1.00	1.78	<b>3.64</b>	3.39
OURS VS. BASELINE		<b>+16.0%</b>	<b>+36.8%</b>	<b>+52.0%</b>	<b>+32.4%</b>

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 [38]. The best results are shown in bold.

Training Set	Non-Rigid Motion	Testing Set	
		MAE↓	MAPE↓
UBFC-rPPG	Very Small	10.84	11.40
MAUBFC-rPPG	Small	1.86	2.94
MAUBFC-rPPG	Large	<b>1.17</b>	<b>1.55</b>
OURS VS. BASELINE		<b>+89.2%</b>	<b>+86.4%</b>

### 4.3. Naturalistic versus Synthetic 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 synthetic data with motion generated using parametric motion animation, as used in the SCAMPS dataset [24]. The SCAMPS dataset consists of synthetic human heads that can be rigged to induce parametric motion. We consider 200 such samples from Table 5. **Naturalistic versus Synthetic Head Motion.** We evaluate the effect of adding head motions to SCAMPS and UBFC-rPPG. The best results are shown in bold.

Training Set	Testing Set				Synth. Time
	PURE		AFRL		
	MAE↓	MAPE↓	MAE↓	MAPE↓	
SCAMPS-200 (No motion)	10.29	11.09	7.75	10.54	37.00s
SCAMPS-200 (Motion)	5.38	5.42	7.25	10.20	37.00s
UBFC-rPPG	4.55	4.67	4.72	6.59	-
MASCAMPS-200	4.67	4.22	5.00	6.69	1.20s
MAUBFC-rPPG	<b>0.96</b>	<b>1.13</b>	<b>3.24</b>	<b>4.37</b>	2.39s
MASCAMPS vs. SCAMPS	<b>+13.2%</b>	<b>+22.1%</b>	<b>+31.1%</b>	<b>+34.4%</b>	<b>+96.8%</b>
MAUBFC vs. UBFC-rPPG	<b>+78.9%</b>	<b>+75.8%</b>	<b>+31.4%</b>	<b>+33.7%</b>	-

Avg. Synth. Time = time (in seconds) to synthesize a frame

the SCAMPS dataset that consist of significant synthetically generated rigid and non-rigid head motion. We then take instances from the SCAMPS dataset with no head motion and augment them with naturalistic head motion to produce MASCAMPS-200. Further details on this experiment, including comparisons to additional results using Wang et al.’s [46] synthetic rPPG video data, are included in the supplementary materials.

We train TS-CAN on both SCAMPS-200 (Motion) and MASCAMPS-200, and evaluated its performance on PURE and AFRL, as shown in Table 5. 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, 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.

### 4.4. Effect of Neural Motion Transfer Algorithms

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 motion transfer algorithm used for motion augmentation. In Table 6, we evaluate two additional neural motion transfer methods in addition to face-vid2vid [44] - FOMM [34] and DaGAN [11]. MAE and SNR are calculated using predictions from TS-CAN and the ground truth label. These results show that most motion transfer algorithms can serve as an effective data augmentation tool as long as they utilize a keypoint-based approach for transferring motion toward applications such as neural talking head synthesis.

Table 6. **Effect of Motion Transfer Methods.** We compare numerous SOTA neural motion transfer methods on UBFC-rPPG [3].

Method	MAE↓	SNR↑
Baseline (No Augmentation)	3.93	4.72
FOMM [34]	0.92	8.64
DaGAN [11]	1.23	8.37
face-vid2vid [44]	0.96	8.70

### 4.5. Effect of rPPG 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, in addition to TS-CAN, we evaluate two more rPPG models - DeepPhys and PhysNet - in Table 7. We utilize MAUBFC-rPPG as training data and evaluate on PURE. We observe

that the results are reasonably consistent across neural rPPG models.

Table 7. **Generalization to Different rPPG 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↓	MAPE↓
UBFC-rPPG	DeepPhys [4]	5.14	4.90
MAUBFC-rPPG	DeepPhys	1.24	1.56
UBFC-rPPG	PhysNet [50]	8.06	13.67
MAUBFC-rPPG	PhysNet	2.38	2.44
UBFC-rPPG	TS-CAN [15]	4.55	4.67
MAUBFC-rPPG	TS-CAN	<b>0.96</b>	<b>1.13</b>

## 5. Discussion

### Can motion augmented videos achieve SOTA results?

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 2). 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 79% 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). In Table 1, we also demonstrate, using UBFC-rPPG as a source dataset, the effectiveness of our method in contrast to other SOTA methods using the same source dataset to test on PURE.

**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 3) and non-rigid (see Table 4) head motions. Furthermore, if the motions have a larger magnitude, then including larger magnitude motions in the training set empirically seems to have a benefit.

### Does the type of motion transfer algorithm or the type of PPG estimation model matter?

It’s important to decouple other factors that may significantly affect performance such as the neural motion transfer algorithm or the type of rPPG estimation model. As per Table 6 and Table 7, it is clear that our motion augmentation strategy provides significant improvements 1) regardless of the core neural motion transfer algorithm utilized and 2) despite differences (e.g., 2DCNN versus 3DCNN) in neural rPPG estimation models. This in turn shows great promise in motion augmentation as a general data augmen-

tation strategy for rPPG videos.

**Is natural motion augmentation best?** Finally, there are different methods for synthesizing motion in video data. SOTA 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 5 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 [46] 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 79% reduction in error in cross-dataset experiments using TS-CAN and a 47% reduction in error when compared to other state-of-the-art methods using the same source dataset.



## 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. 5
- [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. 1, 2
- [3] Serge Bobbia, Richard Macwan, Yannick Benezeth, Alamin Mansouri, and Julien Dubois. Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters*, 124:82–90, 2019. 1, 2, 3, 4, 5, 7
- [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. 1, 2, 4, 8
- [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. 2
- [6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 2
- [7] Justin R Estep, 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. 5, 7
- [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. 2
- [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. 2
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2
- [11] 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. 2, 7
- [12] 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. 1
- [13] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018. 2
- [14] 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. 2
- [15] Xin Liu, Josh Fromm, Shwetak Patel, and Daniel McDuff. Multi-task temporal shift attention networks for on-device contactless vitals measurement. *NeurIPS*, 2020. 1, 2, 4, 5, 8
- [16] Xin Liu, Brian Hill, Ziheng Jiang, Shwetak Patel, and Daniel McDuff. Efficientphys: Enabling simple, fast and accurate camera-based cardiac measurement. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5008–5017, January 2023. 5
- [17] 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. 2
- [18] Xin Liu, Girish Narayanswamy, Akshay Paruchuri, Xiaoyu Zhang, Jiankai Tang, Yuzhe Zhang, Yuntao Wang, Soumyadip Sengupta, Shwetak Patel, and Daniel McDuff. rppg-toolbox: Deep remote ppg toolbox. *arXiv preprint arXiv:2210.00716*, 2022. 5
- [19] Hao Lu, Hu Han, and S Kevin Zhou. Dual-gan: Joint bvp and noise modeling for remote physiological measurement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12404–12413, 2021. 5
- [20] Daniel McDuff. Camera measurement of physiological vital signs. *ACM Computing Surveys (CSUR)*, 2021. 1
- [21] Daniel McDuff, Roger Cheng, and Ashish Kapoor. Identifying bias in ai using simulation. 2018. 2
- [22] Daniel McDuff, Xin Liu, Javier Hernandez, Erroll Wood, and Tadas Baltrušaitis. 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. 3
- [23] 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. 2
- [24] Daniel McDuff, Miah Wander, Xin Liu, Brian L Hill, Javier Hernandez, Jonathan Lester, and Tadas Baltrušaitis. Scamps: Synthetics for camera measurement of physiological signals. *arXiv preprint arXiv:2206.04197*, 2022. 2, 7
- [25] Rita Meziasabour, Yannick Benezeth, 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. 3, 5
- [26] Girish Narayanswamy, Yujia Liu, Yuzhe Yang, Chengqian Ma, Xin Liu, Daniel McDuff, and Shwetak Patel. Bigsmall: Efficient multi-task learning for disparate spatial and temporal physiological measurements. *arXiv preprint arXiv:2303.11573*, 2023. 2
- [27] Xuesong Niu, Hu Han, Shiguang Shan, and Xilin Chen. Vipl-hr: A multi-modal database for pulse esti-

- mation from less-constrained face video. *arXiv preprint arXiv:1810.04927*, 2018. 3
- [28] 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. 1
- [29] 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. 5
- [30] 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. 1, 2
- [31] 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. 1
- [32] 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. 2
- [33] 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. 2
- [34] 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. 2, 7
- [35] 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. 2
- [36] Jeremy Speth, Nathan Vance, Patrick Flynn, and Adam Czajka. Non-contrastive unsupervised learning of physiological signals from video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14464–14474, June 2023. 2, 5
- [37] 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. 1
- [38] 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. 1, 2, 3, 4, 5, 6, 7
- [39] 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. 2
- [40] Chihiro Takano and Yuji Ohta. Heart rate measurement based on a time-lapse image. *Medical engineering & physics*, 29(8):853–857, 2007. 2
- [41] Jiankai Tang, Kequan Chen, Yuntao Wang, Yuanchun Shi, Shwetak Patel, Daniel McDuff, and Xin Liu. Mmpd: Multi-domain mobile video physiology dataset, 2023. 3, 5
- [42] Wim Verkrusse, Lars O Svaasand, and J Stuart Nelson. Remote plethysmographic imaging using ambient light. *Optics express*, 16(26):21434–21445, 2008. 1, 2
- [43] 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. 2
- [44] 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. 2, 4, 5, 7
- [45] 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. 2, 5
- [46] 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. 2, 3, 7, 8
- [47] 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. 2
- [48] 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. 2
- [49] 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. 2
- [50] 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. 1, 2, 4, 5, 8
- [51] 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. 1, 2, 5
- [52] 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. 2