

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Motion Matters: Neural Motion Transfer for Better Camera Physiological Sensing

Anonymous ICCV submission

Paper ID 1080

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.

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 [26]. Camera-based measurement of cardiac signals is one such application [19], in which cameras are used to measure the pulse via light reflected from the body, a principle known as photoplethysmography (PPG) [2, 39]. The PPG signals can be used to derive respiration [28], heart rate variability [28], arrhythmia [29], 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

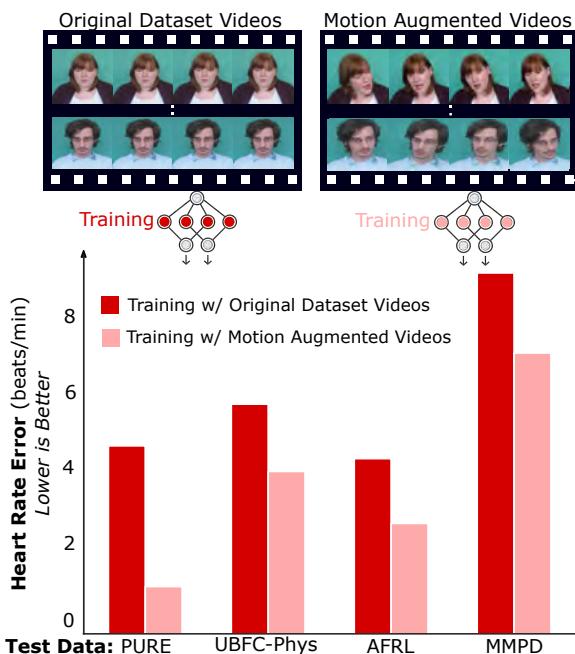


Figure 1: **Motion augmentation improves rPPG.** We present the first neural motion augmentation pipeline for the task of remote PPG estimation and empirically show it reduces error in heart rate estimation by up to 75%.

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, 34, 47, 16, 48]. 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 [35]) 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,

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

108 more data efficient methods for training rPPG sensing mod-
109 els would be desirable.
110

111 Synthetic data are a powerful resource in machine learning.
112 The two main sources of synthetic data are (1) parametric
113 computer graphics engines and (2) statistically-based
114 generative machine learning models. Data created using
115 these approaches have been used successfully for many
116 computer vision tasks, including face detection, landmark
117 localization, face parsing and face recognition [20, 44, 22],
118 body pose estimation [31] and eye tracking [45, 36].
119

120 However, creating synthetic data that preserve the subtle
121 and nuanced peripheral pulse in a video is non-trivial.
122 McDuff et al. [23] released a large dataset (2,800 videos)
123 of avatars and cardiac signals; however, their computer
124 graphics pipeline had an extremely high overhead. Creating
125 a pipeline for generating videos of avatars required
126 years worth of investment, procuring assets (e.g., 3D fac-
127 ial scans and environments), building a parameterized ren-
128 dering pipeline, and then generating each video frame-by-
129 frame. Wang et al. [43] used a learning based method to
130 generate synthetic videos given a reference image and tar-
131 get PPG signal. Their creative approach successfully incor-
132 porated PPG signals producing videos that benefited train-
133 ing. However, the videos created lacked the visual fidelity
134 of other synthetics or real video datasets, and their pipeline
involved several relatively complex components.
135

136 We question whether existing motion transfer algorithms
137 can be used effectively for augmenting rPPG video data
138 and explore what steps need to be taken to achieve opti-
139 mal results. Our main contributions are as follows: (1) We
140 perform a systematic investigation of the impact of motion
141 augmentation on the physiological information within the
142 video and in-turn the corresponding labels, (2) through ex-
143 perimentation, we provide quantitative, empirical evidence
144 that training with certain kinds of motion-augmented data is
145 effective for camera-based physiological measurement
146 algorithms, and (3) we demonstrate, through inter-dataset re-
147 sults, the usefulness of motion augmentation for improving
148 the generalization of models for camera-based physiolog-
149 ical sensing. We achieve state-of-the-art results on mul-
150 tiple public benchmark datasets, including those with signif-
151 icant motion. We summarize the key findings of this paper
152 about the effectiveness of motion transfer as a data aug-
153 mentation tool in Sec. 5. We provide our code for augmenting
154 datasets, training using these data, and pre-trained models
155 trained on motion-augmented data (all assets are released
156 with responsible use licenses [5]).
157

2. Background

158 **Generative Synthetics for Training Models:** Statistical
159 generative models [9, 14, 13, 33, 10, 6] capture a proba-
160 bilistic representation of a dataset from which samples can
161 be drawn. These models are typically trained to mimic the
distribution of the training set and can be trained without the

162 need for labels, allowing large sets of data to be used. Fa-
163 cial video generation using generative models has advanced
164 rapidly over recent years [15, 30]. Numerous image-driven
165 works have accomplished the ability to separate identity and
166 pose in source and driving images used for high quality, ro-
167 bust video generation using generative adversarial networks
168 (GANs) [49, 32, 41, 11]. Image-driven facial video genera-
169 tion methods attempt to preserve the identity of a given
170 source image while manipulating the pose based on a driv-
171 ing video to generate a new video. The identity from the
172 driving video is excluded with the help of a keypoint-based
173 motion transfer approach, where keypoints are predicted
174 for both a source image and a driving image in order to
175 model local motion using shifts in the corresponding key-
176 points [32, 41, 11]. Face video generation that is achieved
177 by using keypoints that take pose and expression into ac-
178 count can be successful for the task of head video genera-
179 tion, but can at times have a loss in source image iden-
180 tity and unwanted temporal artifacts [32, 49, 11]. Face-
181 Vid2Vid [41] utilizes canonical keypoints in addition to
182 source and driving image keypoints in order to capture a tar-
183 get person’s geometry signature, which includes the shape
184 of the target’s face, nose, and eyes. This allows for im-
185 proved head video generation that minimizes source iden-
186 tity loss while effectively transferring motion from a driving
187 video.
188

189 **rPPG Models:** The principle that photoplethysmogra-
190 phy could be performed with a camera and without contact
191 with the body was established by Blazek et al. [2] and repli-
192 cated in a series of following experiments [37, 39]. The
193 application of more advanced signal processing methods
194 helped make measurement somewhat more robust under
195 real-world conditions [28, 42], as did leveraging knowledge
196 of physiological and physical properties [42]. Yet, these
197 models were still very sensitive to body motions. Neural
198 data-driven models currently achieve state-of-the-art results
199 in most cases [4, 47, 16, 48, 17], but are a function of the
200 data used to train them. While intra-dataset performance
201 is generally strong, inter-dataset performance is often sub-
202 stantively worse. In order to alleviate the dependency on la-
203 beled data, several researchers have proposed unsupervised
204 learning procedures [8, 40, 46]. However, most require fine-
205 tuning on a labeled set and also reveal that supervised learn-
206 ing still holds some additional benefit. As an alternative or
207 a complement, generative methods have been suggested to
208 “create” data [21, 43].
209

210 **rPPG Datasets:** As with many health applications,
211 those working in camera physiological measurement face
212 challenges associated with collecting and managing data.
213 Public datasets (such as UBFC-rPPG [3], PURE [35],
214 VIPL-HR [25]) are valuable resources. However, given the
215 challenging nature of the rPPG task researchers have col-
216 lected and released data under heavily constrained condi-
217

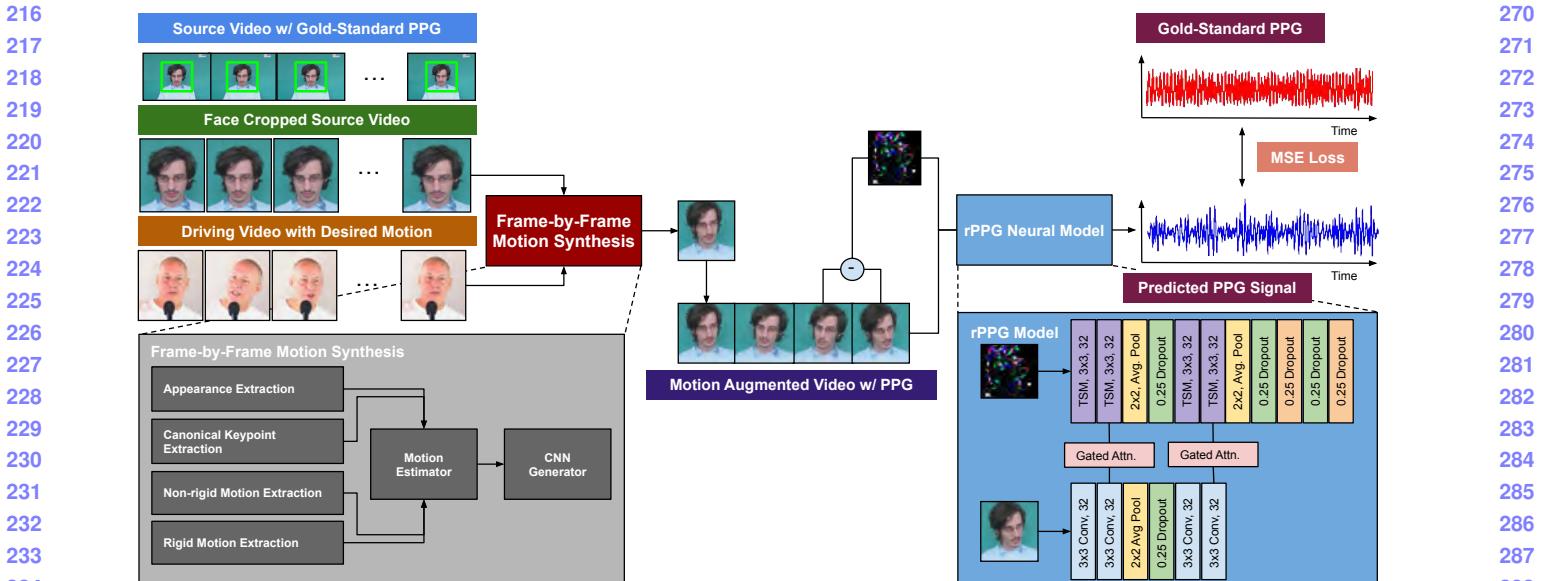


Figure 2: **Motion augmentation and training pipeline.** We augment each frames of a source video with corresponding frames of a randomly selected driving video to create an augmented video with the identity of the source video and motion of the target video. We then train a PPG estimation network on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

tions with very little physical motion. More recent datasets (such as UBFC-PHYS [24] and MMPD [38]) contain larger and more natural motions. However, the baseline results on these datasets are not very strong.

3. Motion Augmentation 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 breath-

ing 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 [41], 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 redesign

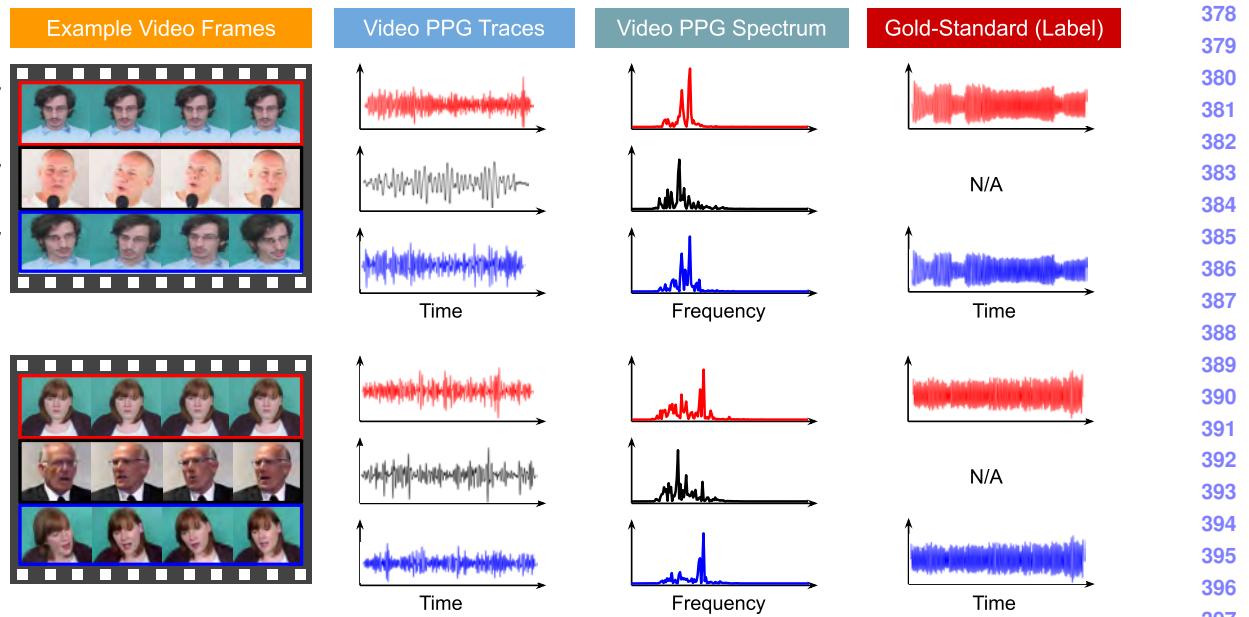


Figure 3: **Preserving physiological signals in motion augmented videos.** We show that applying neural motion transfer preserve the physiological signal corresponding to the heart-rate present in the peak of the frequency spectrum of the source and the augmented video.

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 [35] 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 [41] 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 [41] 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 [16] to predict the 1st-order derivative of the PPG signal after training on videos

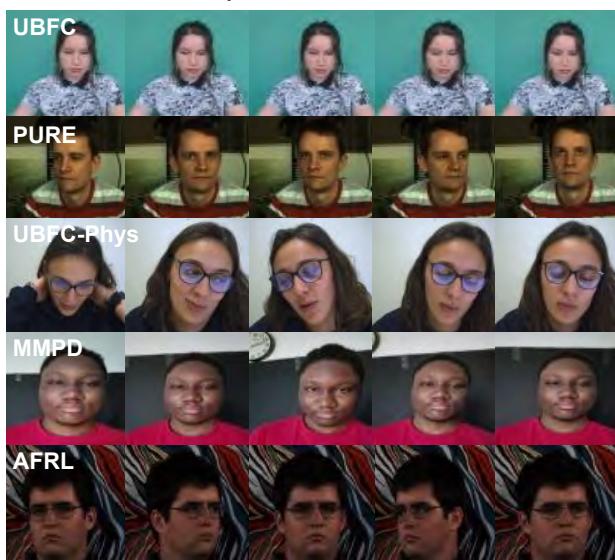
augmented with motion. We also use DeepPhys [4] and EfficientPhys [17] to highlight the consistent benefits of motion augmentation across different neural models.

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 [41] 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 [42], 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

432 Table 1: A Summary of the rPPG Benchmark Datasets.
433

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

457 different appearances and motion conditions, both in the
458 source videos and the driving videos. We present additional
459 results in the supplementary materials. Thus, we can effec-
460 tively claim that motion transfer algorithms like Face-
461 Vid2Vid [41] do preserve the underlying physiological sig-
462 nal, like heart rate, and they can be a very effective tool
463 for large scale augmentation of training videos for PPG es-
464 timation tasks. Our quantitative experimental results show
465 that deep neural networks for camera physiological mea-
466 surement can take advantage of this to significantly improve
467 model performance by training on motion-augmented data.

4. Experiments

470 We consider five datasets for training and evalua-
471 tion, **UBFC-rPPG** [3], **PURE** [35], **UBFC-PHYS** [24],
472 **AFRL** [7], and **MMPD** [38] (see Table 1). They consist
473 of facial videos and corresponding gold-standard PPG sig-
474 nals. We use some of these datasets for augmentation
475 with neural motion transfer and training the rPPG models,
476 and use the rest to evaluate different aspects of the effective-
477 ness of neural motion augmentation. To our knowledge, we
478 perform the most extensive inter-dataset evaluation of PPG
479 estimation to date, testing on five independent test datasets.

480 **Implementation Details:** For all of the results reported
481 in this section, the predicted PPG signals were filtered using
482 a band-pass signal with cut-offs 0.75 Hz and 2.5 Hz. The
483 heart rate was calculated based on the PPG signal using the
484 Fast Fourier Transform (FFT), with a measurement window
485 of the video length for UBFC-rPPG, PURE, UBFC-PHYS,

486 and MMPD. To evaluate the AFRL dataset, a measurement
487 window of 30 seconds was utilized for heart rate calcula-
488 tions. All networks were trained using an NVIDIA RTX
489 A4500 and PyTorch [27] implementations in the publicly
490 available rPPG-Toolbox [18]. A cyclic learning rate sched-
491 ule was utilized with 30 epochs, a learning rate of 0.009,
492 and a batch size of 4 for both training and inference.

4.1. Training with Motion Augmented Data

494 In Table 2, we compare the performance of a supervised
495 PPG estimation network, TS-CAN [16], trained on existing
496 video datasets and motion-augmented versions of those
497 datasets. We also show the performance of unsupervised
498 methods for comparison. For the sake of space and clar-
499 ity, the following tables only show the mean absolute er-
500 rror (MAE) in heart rate estimation for each video (and the
501 corresponding mean absolute percentage error (MAPE)).
502 Equivalent tables with root mean squared error (RMSE) and
503 Pearson correlation metrics can be found in our supple-
504 mentary material. The driving videos used for augmentation in
505 Table 2 contain significant amounts of unconstrained mo-
506 tion – both rigid and non-rigid.

507 We observe that training TS-CAN on augmented videos
508 produces state-of-the-art (SOTA) performance in most
509 cases. Additionally, we observe that in most cases, the aug-
510 mented versions outperform the non-augmented versions,
511 with a gain in performance up to 75% and an average gain of
512 26%. However, when comparing the performance of MA-
513 PURE versus PURE when tested on UBFC-PHYS, we note
514 a minor drop in performance rather than an improvement
515 due to the difficulty in effectively augmenting the PURE
516 dataset. This is because the PURE dataset already contains
517 significant amounts of rigid motion, and when augmented,
518 it may provide training data with artifacts that make the
519 learned rPPG task less useful in the face of a highly un-
520 constrained dataset with natural rigid and non-rigid motion.

521 **Details:** We utilize all downloadable videos from the
522 TalkingHead-1KH [41] dataset as our driving videos for
523 augmenting various PPG estimation datasets with motion.
524 We analyze the videos using OpenFace [1] to obtain the in-
525 tensity (0 to 5) of 17 Facial Action Units (AUs) and the head
526 pose rotations R_x , R_y , and R_z in radians (rad). To generate
527 MAUBFC-rPPG, we choose driving videos from a pool of
528 60 driving videos with a range of mean standard deviation
529 in head pose rotations from 0.10 to 0.14 rad to augment as
530 much rigid motion as possible into a source video dataset
531 that has very little of both rigid and non-rigid motion. We
532 do not constrain for non-rigid motion in this case, so we
533 observe a wide range of mean standard deviation in facial
534 AUs from 0.15 to 0.5 intensity. To generate MAPURE, we
535 choose driving videos with a range of mean standard de-
536 viation in facial AUs from 0.45 to 0.55 intensity to aug-
537 ment as much non-rigid motion as possible into a source

540 video dataset that has very little non-rigid motion. We do
541 not constrain for rigid motion in this case, so we observe a
542 wide range of mean standard deviation in head pose rota-
543 tions from 0.03 to 0.14 rad.
544

545 4.2. Effect of Motion Types

546 A key question in designing a motion augmentation strat-
547 egy is deciding what type of motion should be applied to
548 obtain the best performance on a certain evaluation dataset.
549 To answer this question, we separately analyze two types
550 of motion: rigid and non-rigid, by augmenting training data
551 with different magnitudes of motion. Rigid motion refers
552 to head pose rotation, while having minimal change in
553 facial action units or expressions. Non-rigid motion refers to
554 changes in facial expression, i.e. motion in facial action
555 units for various tasks like talking, while having minimal
556 head pose rotation.
557

558 **Rigid Motion:** For rigid motion, we consider UBFC-
559 rPPG as training data, which has very little head motion
560 and AFRL as test data which has large variations in rigid
561 head motion. We classify videos in the AFRL dataset into
562 different rigid head motion categories: ‘very small motion’,
563 ‘small motion’ (10 deg rotation per sec), and ‘large motion’
564 (30 deg rotation per sec). Based on this categorization, we
565 also select driving videos from our captured CDVS to have
566 ‘small motion’ and ‘large motion’ using the mean standard
567 deviation in estimated head pose rotations across all the
568 frames of a video. Specifically, for ‘small motion’ we used
569 mean standard deviation between 0.03 to 0.07 rad and for
570 ‘large motion’ between 0.10 to 0.14 rad. These parameters
571 are chosen to roughly match the distribution of head pose
572 rotation in ‘small motion’ and ‘large motion’ categories of
573 AFRL. We then use these videos from the CDVS dataset to
574 augment the source videos of UBFC-rPPG to create 3 sepa-
575 rate categories of augmented videos for ‘very small motion’
576 (which is the original UBFC-rPPG dataset), ‘small motion’,
577 and ‘large motion’ respectively. We then train TS-CAN on
578 augmented data in each category and test on the same cat-
579 egories of the AFRL dataset. We present these results in
580 Table 3.

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

588 **Non-rigid Motion:** For non-rigid motion, we also con-
589 sider UBFC-rPPG as training data since it has very little
590 motion, and the speech task of the PURE dataset [35] as the
591 test data which has significant non-rigid head motion. We
592 also augment the UBFC-rPPG dataset with non-rigid head
593 motion from our captured CDVS with ‘small’ and ‘large’

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

607 4.3. Effect of Multiple Augmentations

608 We consider whether it is plausible to augment the same
609 source video with multiple driving videos using neural mo-
610 tion transfer. Thus, the newly augmented dataset has the
611 same number of identities as the original dataset but a sig-
612 nificantly larger variation in motions. Our goal is to analyze
613 how many times one can augment a single source video be-
614 fore the performance starts to saturate or drop. We con-
615 sider UBFC-rPPG as training data that we augment with
616 randomly sampled driving videos from the TalkingHead-
617 1K dataset to produce MAUBFC-rPPG. We augment the
618 same source video from 1 to 4 times with different driv-
619 ing videos and evaluate on the PURE [35] dataset and the
620 UBFC-PHYS [24] dataset and report the results in Table 5.
621 We notice that the results saturate pretty quickly and can
622 start to decline after augmenting more than 2 times.
623

624 4.4. Synthetic vs Naturalistic Head Motion

625 In order to further evaluate the impact of motion transfer
626 as a data augmentation technique, we explore whether data
627 augmented with natural head motion using a neural motion
628 transfer algorithm is better than augmenting data with syn-
629 thetically generated motion using parametric motion anima-
630 tion, as used in the SCAMPS dataset [23]. The SCAMPS
631 dataset consists of synthetic human heads that can be rigged
632 to induce parametric motion. We consider 200 such sam-
633 ples from the SCAMPS dataset that consist of significant
634 synthetically generated rigid and non-rigid head motion (ID
635 1801 to 2000) as SCAMPS-200 (Motion). We then take in-
636 stances from the SCAMPS dataset with no head motion (ID
637 1 to 200) and augment them with naturalistic head motion
638 using our motion synthesis pipeline and a subset of driv-
639 ing videos from the TalkingHead-1KH dataset to produce
640 MASCAMPS-200. We choose driving videos with a range
641 of mean standard deviation in AUs from 0.35 to 0.40 inten-
642 sity and a range of mean standard deviation in head pose
643 rotations from 0.05 to 0.125 rad. Note that both SCAMPS-
644 200 (Motion) and MASCAMPS-200 consist of synthetics
645 with the same number of identities, with the only difference
646

648 Table 2: **Evaluation across all datasets.** We motion-augment two training datasets, UBFC-rPPG and PURE, to create MAUBFC-rPPG
 649 and MAPURE, respectively. We observe that the motion-augmented versions produce significant improvements (shown in bold).
 650

651 Training Set	652 Method	653 UBFC-rPPG		654 PURE		655 Testing Set		656 UBFC-PHYS		657 AFRL		658 MMPD	
		659 MAE \downarrow	660 MAPE \downarrow	661 MAE \downarrow	662 MAPE \downarrow	663 MAE \downarrow	664 MAPE \downarrow	665 MAE \downarrow	666 MAPE \downarrow	667 MAE \downarrow	668 MAPE \downarrow	669 MAE \downarrow	670 MAPE \downarrow
654 Unsupervised	Green	19.82	18.78	10.09	10.28	13.45	16.00	7.01	9.24	16.27	20.09	702	703
	ICA	14.70	14.34	4.77	4.47	8.00	9.48	6.77	8.96	13.10	16.33	704	705
	CHROM	3.98	3.78	5.77	11.52	4.68	6.20	5.41	7.95	8.85	11.93	706	707
	POS	4.00	3.86	3.67	7.25	4.62	6.29	6.93	10.00	8.18	11.12	708	709
655 UBFC-rPPG	TS-CAN	-	-	4.55	4.67	5.56	7.25	4.24	5.84	8.74	10.51	710	711
656 MAUBFC-rPPG	TS-CAN	-	-	1.14	1.30	3.93	5.24	2.67	3.65	6.80	7.97	712	713
PURE	TS-CAN	1.34	1.55	-	-	4.43	5.89	2.63	3.51	8.96	10.33	714	715
MAPURE	TS-CAN	1.03	1.17	-	-	4.39	5.90	2.37	3.26	8.08	9.54	716	717
MAUBFC-rPPG vs. UBFC-rPPG		-	-	+74.95%	+72.16%	+29.32%	+27.72%	+37.03%	+37.50%	+22.20%	+24.17%	718	719
MAPURE vs. PURE		+23.13%	+24.52%	-	-	+0.90%	-0.17%	+9.89%	+7.12%	+9.82%	+7.65%	720	721

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

664
 665 Table 3: **Effect of Motion Types – Rigid.** We augment
 666 UBFC-rPPG with various types of rigid head motions and test on
 667 AFRL [7].

668 Training Set	669 Rigid Motion	670 Testing Set			
		671 No Motion	672 Small Motion	673 Large Motion	674 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%			

675 Table 4: **Effect of Motion Types – Non-rigid.** We augment
 676 UBFC-rPPG with various types of non-rigid motions (expressions)
 677 and test on the speech task, in PURE [35].

678 Training Set	679 Non-Rigid Motion	680 Testing Set			
		681 Non-rigid Motion Task	682 MAE \downarrow	683 MAPE \downarrow	684
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%			

691 being synthetic and naturalistic head motion, respectively.

692 We train TS-CAN on both SCAMPS-200 (Motion) and
 693 MASCAMPS-200, and evaluated its performance on PURE
 694 and AFRL, as shown in Table 6. We observed that adding
 695 naturalistic motion improved performance by 13.2% on
 696 PURE and 31.1% on AFRL compared to synthetically gen-
 697 erated motion. It is worth noting that the average time taken
 698 to add synthetic motion to each frame of a sequence is 37
 699 seconds, compared to only 1.2 seconds for adding nat-
 700 ualistic motion using the neural motion transfer algorithm.
 701 For comparison, we also included real-world training data,

702
 703 Table 5: **Effect of Multiple Augmentations.** Augmenting each
 704 source video of UBFC-rPPG 1x, 2x, 3x, and 4x, we test on PURE
 705 and UBFC-PHYS datasets. The best results are shown in bold.

706 Training Set	707 Size	708 Subjects	709 Testing Set				710
			PURE	UBFC-PHYS	711	712	713
UBFC-rPPG	42	42	4.55	4.67	5.56	7.25	714
MAUBFC-rPPG	42	42	1.14	1.30	3.93	5.24	715
MAUBFC-rPPG 2x	84	42	1.12	1.29	3.90	5.22	716
MAUBFC-rPPG 3x	126	42	1.11	1.27	3.97	5.31	717
MAUBFC-rPPG 4x	168	42	1.19	1.33	4.10	5.40	718
OURS vs. BASELINE			+2.63%	+2.31%	+0.76%	+0.38%	719

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

725 Training Set	726 Size	727 Subjects	728 Testing Set				729
			PURE	AFRL	730	731	732
SCAMPS-200 (No motion)	10.29	11.09	7.75	10.54	37.00s	733	734
SCAMPS-200 (Motion)	5.38	5.42	7.25	10.20	37.00s	735	736
UBFC-rPPG	4.55	4.67	4.72	6.59	-	737	738
MASCAMPS-200	4.67	4.22	5.00	6.69	1.20s	739	740
MAUBFC-rPPG	1.14	1.30	3.24	4.37	2.39s	741	742
MASCAMPS vs. SCAMPS	+13.2%		+22.1%	+31.1%	+34.4%	+96.8%	743
MAUBFC vs. UBFC-rPPG	+74.4%		+72.2%	+31.4%	+33.7%	-	744

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

745 UBFC-rPPG, which showed that having real images signif-
 746 icantly improved performance over synthetic images. Fur-
 747 thermore, the only way to augment real images is to use the
 748 neural motion transfer algorithm, as parametric rigged head
 749 motion cannot be applied to real data.

4.5. Effect of PPG Estimation Models

750 It is important to decouple any data augmentation tech-
 751 nique from additional factors that affect its usefulness for
 752 a given set of training data. One such factor is the neu-
 753 ral network model used for training and evaluation. Thus,
 754

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
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 [17]	4.95	4.56
MAUBFC-rPPG	EfficientPhys	1.45	1.76
UBFC-rPPG	TS-CAN [16]	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 7. 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 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 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 3) and non-rigid (see Table 4) 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 [41]. 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 5). 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 6 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 [43] 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.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [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 Benetech, 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] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [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.
- [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.
- [13] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.
- [14] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- [15] 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.
- [16] Xin Liu, Josh Fromm, Shwetak Patel, and Daniel McDuff. Multi-task temporal shift attention networks for on-device contactless vitals measurement. *NeurIPS*, 2020.
- [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.
- [18] 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.
- [19] Daniel McDuff. Camera measurement of physiological vital signs. *ACM Computing Surveys (CSUR)*, 2021.
- [20] Daniel McDuff, Roger Cheng, and Ashish Kapoor. Identifying bias in ai using simulation. 2018.
- [21] 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.
- [22] 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.
- [23] 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.
- [24] Rita Meziatisabour, Yannick Benetech, 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.
- [25] 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.
- [26] 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.
- [27] 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.
- [28] 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.
- [29] 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

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

- 972 Chung-Wah Siu. Diagnostic assessment of a deep learning
973 system for detecting atrial fibrillation in pulse waveforms.
974 *Heart*, 104(23):1921–1928, 2018.
- 975 [30] Tong Sha, Wei Zhang, Tong Shen, Zhoujun Li, and Tao Mei.
976 Deep person generation: A survey from the perspective of
977 face, pose and cloth synthesis. *ACM Computing Surveys*,
978 2021.
- 979 [31] Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp,
980 Mark Finocchio, Richard Moore, Alex Kipman, and Andrew
981 Blake. Real-time human pose recognition in parts from sin-
982 gle depth images. In *CVPR 2011*, pages 1297–1304. Ieee,
983 2011.
- 984 [32] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov,
985 Elisa Ricci, and Nicu Sebe. First order motion model for
986 image animation. In *Conference on Neural Information Pro-
987 cessing Systems (NeurIPS)*, December 2019.
- 988 [33] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan,
989 and Surya Ganguli. Deep unsupervised learning using
990 nonequilibrium thermodynamics. In *International Confer-
991 ence on Machine Learning*, pages 2256–2265. PMLR, 2015.
- 992 [34] Radim Špetlík, Vojtech Franc, and Jirí Matas. Visual heart
993 rate estimation with convolutional neural network. In *Pro-
994 ceedings of the british machine vision conference, Newcas-
995 tle, UK*, pages 3–6, 2018.
- 996 [35] Ronny Stricker, Steffen Müller, and Horst-Michael Gross.
997 Non-contact video-based pulse rate measurement on a mo-
998 bile service robot. In *The 23rd IEEE International Sym-
999 posium on Robot and Human Interactive Communication*,
1000 pages 1056–1062. IEEE, 2014.
- 1001 [36] Lech Świński and Neil Dodgson. Rendering synthetic ground
1002 truth images for eye tracker evaluation. In *Proceedings of
1003 the Symposium on Eye Tracking Research and Applications*,
1004 pages 219–222, 2014.
- 1005 [37] Chihiro Takano and Yuji Ohta. Heart rate measurement
1006 based on a time-lapse image. *Medical engineering &*
1007 *physics*, 29(8):853–857, 2007.
- 1008 [38] Jiankai Tang, Kequan Chen, Yuntao Wang, Yuanchun Shi,
1009 Shwetak Patel, Daniel McDuff, and Xin Liu. Mmpd: Multi-
1010 domain mobile video physiology dataset, 2023.
- 1011 [39] Wim Verkruyse, Lars O Svaasand, and J Stuart Nelson. Re-
1012 mote plethysmographic imaging using ambient light. *Optics*
1013 *express*, 16(26):21434–21445, 2008.
- 1014 [40] Hao Wang, Euijoon Ahn, and Jinman Kim. Self-supervised
1015 representation learning framework for remote physiological
1016 measurement using spatiotemporal augmentation loss. In
1017 *Proceedings of the AAAI Conference on Artificial Intelli-
1018 gence*, volume 36, pages 2431–2439, 2022.
- 1019 [41] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot
1020 free-view neural talking-head synthesis for video conferenc-
1021 ing. *2021 IEEE/CVF Conference on Computer Vision and
1022 Pattern Recognition (CVPR)*, pages 10034–10044, 2020.
- 1023 [42] Wenjin Wang, Albertus C den Brinker, Sander Stuijk, and
1024 Gerard de Haan. Algorithmic principles of remote ppg. *IEEE
1025 Transactions on Biomedical Engineering*, 64(7):1479–1491,
1026 2017.
- 1027 [43] Zhen Wang, Yunhao Ba, Pradyumna Chari, Oyku Deniz
1028 Bozkurt, Gianna Brown, Parth Patwa, Niranjan Vaddi, Laleh
1029 Jalilian, and Achuta Kadambi. Synthetic generation of face
1030 videos with plethysmograph physiology. In *Proceedings of
1031 the IEEE/CVF Conference on Computer Vision and Pattern
1032 Recognition*, pages 20587–20596, 2022.
- 1033 [44] Erroll Wood, Tadas Baltrusaitis, Charlie Hewitt, Sebastian
1034 Dziadzio, Thomas J Cashman, and Jamie Shotton. Fake it
1035 till you make it: Face analysis in the wild using synthetic
1036 data alone. In *Proceedings of the IEEE/CVF International
1037 Conference on Computer Vision*, pages 3681–3691, 2021.
- 1038 [45] Erroll Wood, Tadas Baltrusaitis, Xucong Zhang, Yusuke
1039 Sugano, Peter Robinson, and Andreas Bulling. Rendering of
1040 eyes for eye-shape registration and gaze estimation. In *Pro-
1041 ceedings of the IEEE International Conference on Computer
1042 Vision*, pages 3756–3764, 2015.
- 1043 [46] Yuzhe Yang, Xin Liu, Jiang Wu, Silviu Borac, Dina Katabi,
1044 Ming-Zher Poh, and Daniel McDuff. Simper: Simple
1045 self-supervised learning of periodic targets. *arXiv preprint
1046 arXiv:2210.03115*, 2022.
- 1047 [47] Zitong Yu, Wei Peng, Xiaobai Li, Xiaopeng Hong, and
1048 Guoying Zhao. Remote heart rate measurement from highly
1049 compressed facial videos: an end-to-end deep learning so-
1050 lution with video enhancement. In *Proceedings of the
1051 IEEE/CVF International Conference on Computer Vision*,
1052 pages 151–160, 2019.
- 1053 [48] Zitong Yu, Yuming Shen, Jingang Shi, Hengshuang Zhao,
1054 Philip Torr, and Guoying Zhao. Physformer: Facial video-
1055 based physiological measurement with temporal difference
1056 transformer. *arXiv preprint arXiv:2111.12082*, 2021.
- 1057 [49] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and
1058 Victor Lempitsky. Few-shot adversarial learning of realistic
1059 neural talking head models. In *Proceedings of the IEEE/CVF
1060 international conference on computer vision*, pages 9459–
1061 9468, 2019.
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079