# SCAMPS: Synthetics for Camera Measurement of Physiological Signals

**Daniel McDuff**
Microsoft
Redmond, WA, USA

**Miah Wander**
Microsoft
Redmond, WA, USA

**Xin Liu**
UW
Seattle, WA, USA

**Brian L. Hill**
UCLA
Los Angeles, CA, USA

**Javier Hernandez**
Microsoft
Redmond, WA, USA

**Jonathan Lester**
Microsoft
Redmond, WA, USA
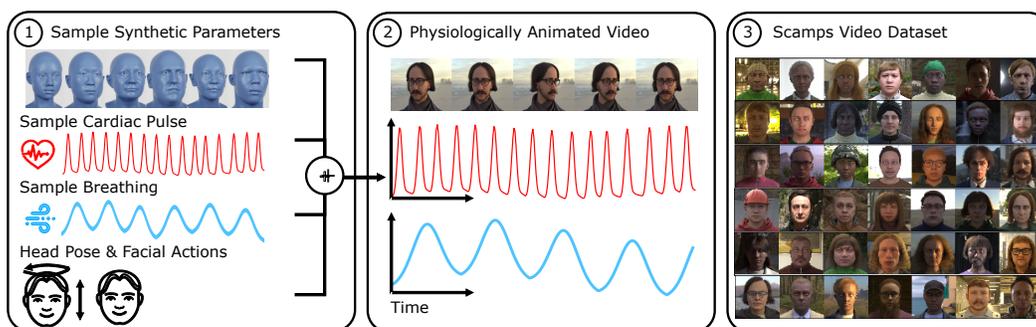
**Tadas Baltrusaitis**
Microsoft
Cambridge, UK

Figure 1: SCAMPS: A dataset of synthetic videos with aligned physiological and behavioral signals.

## Abstract

The use of cameras and computational algorithms for noninvasive, low-cost and scalable measurement of physiological (e.g., cardiac and pulmonary) vital signs is very attractive. However, diverse data representing a range of environments, body motions, illumination conditions and physiological states is laborious, time consuming and expensive to obtain. Synthetic data have proven a valuable tool in several areas of machine learning, yet are not widely available for camera measurement of physiological states. Synthetic data offer "perfect" labels (e.g., without noise and with precise synchronization), labels that may not be possible to obtain otherwise (e.g., precise pixel level segmentation maps) and provide a high degree of control over variation and diversity in the dataset. We present SCAMPS, a dataset of synthetics containing 2,800 videos (1.68M frames) with aligned cardiac and respiratory signals and facial action intensities. The RGB frames are provided alongside segmentation maps. We provide precise descriptive statistics about the underlying waveforms, including inter-beat interval, heart rate variability, and pulse arrival time. Finally, we present baseline results training on these synthetic data and testing on real-world datasets to illustrate generalizability.

Table 1: Summary of Public Camera Physiological Measurement Datasets.

| Dataset | Subjects | Videos | Gold Standard | Sub. Div. | Env. Div. | Free Access |
|---|---|---|---|---|---|---|
| MAHNOB [34] | 27 | 527 | ECG, EEG, Breath. | ✗ | ✗ | ✓ |
| BP4D [53] | 140 | 1400 | BP, AU | ✓ | ✗ | ✗ |
| VIPL-HR [26] | 107 | 3130 | PPG, HR, SpO$_2$ | ✗ | ✗ | ✓ |
| COHFACE [13] | 40 | 160 | PPG | ✗ | ✗ | ✓ |
| UBFC-RPPG [4] | 42 | 42 | PPG, PR | ✗ | ✗ | ✓ |
| UBFC-PHYS [25] | 56 | 168 | PPG, EDA | ✗ | ✗ | ✓ |
| RICE CamHRV [29] | 12 | 60 | PPG | ✗ | ✗ | ✓ |
| MR-NIRP [27] | 18 | 37 | PPG | ✗ | ✗ | ✓ |
| PURE [37] | 10 | 59 | PPG, SpO$_2$ | ✗ | ✗ | ✓ |
| rPPG [15] | 8 | 52 | PR, SpO$_2$ | ✗ | ✗ | ✓ |
| OBF [18] | 106 | 212 | PPG, ECG, BR | ✗ | ✗ | ✗ |
| PFF [14] | 13 | 85 | PR | ✗ | ✗ | ✓ |
| VicarPPG [41] | 20 | 10 | PPG | ✗ | ✗ | ✓ |
| CMU [7] | 140 | 140 | PR | ✓ | ✓ | ✓ |
| SCAMPS* | 2800 | 2800 | PPG, PR, Breath., BR, AU | ✓ | ✓ | ✓ |

ECG = Electrocardiogram waveform, EDA = Electrodermal activity, EEG. = Electroencephalogram waveforms, Breath = Breathing waveform, PPG = Photoplethysmogram waveform, BP = Blood pressure waveform, PR = Pulse rate, BR = Breathing rate, SpO$_2$ = Blood oxygenation, AU = Action Units. * SCAMPS is the only synthetic dataset.

# 1 Introduction

Camera physiological measurement is a rapidly growing field of computer vision and computational photography that leverages imaging devices, signal processing and machine learned models to perform non-contact recovery of vital physiological processes [21]. Data plays an important role in both training and evaluating these models. However, generalization can be weak if the training data are not representative and systematic evaluation can be challenging if testing data do not contain the variations and diversity necessary. Public datasets (e.g., [53, 26, 4]) have contributed significantly to the understanding of algorithmic performance in this domain. These datasets are time consuming to collect, contain highly personally identifiable and sensitive biometrics (including facial videos and physiological waveforms). It is difficult to collect datasets that contain a well distributed set of examples across multiple cardiac and pulmonary parameters (e.g., heart and breathing rates and variabilities, pulse arrival times, waveform morphologies). Furthermore, almost all of these datasets are collected in a single location, with limited diversity in subject appearance, ambient illumination, context and behaviors. Table 1 summarizes some of the properties of these datasets, including whether they are freely (i.e., at no cost) available to researchers in both industry and academia. Finally, at the time of writing, neural architectures [5, 19, 52] provide the state-of-the-art performance for camera measurement of physiology. These models are "data hungry" and often this performance is primarily a function of the availability and quality of the training dataset.

Synthetics have proven valuable in several areas of computer vision, particularly face and body analyses. In training, synthetics have been used successfully to create models for landmark localization and face parsing [50], body pose estimation [33] and eye tracking [51]. Although not completely representative of real observations, synthetics are also valuable in testing (e.g., for face detection [22] or eye tracking [38]). Parameterized computer graphics simulators are one way of testing vision models [44, 45, 46, 43, 31]. Generally, it has been proposed that graphics models be used for performance evaluation [12, 31, 22]. However, increasingly synthetics are also being used to help address shortcomings in performance, such as biases. Kortylewaski et al. [16, 17] show that the damage of real-world dataset biases on facial recognition systems can be partially addressed by pre-training on synthetic data. To address the issue of the lack of representation of skin type in camera physiology datasets computational techniques have been employed to translate real videos from light-skin subjects to dark-skin subjects while being careful to preserve the cardiac signals [1]. A neural generator was used in that work to simulate changes in melanin, or skin tone. However, this approach does not simulate other changes in appearance that might also be correlated with skin type. Nowara et al. [28] used video magnification for augmenting the spatial appearance of videos
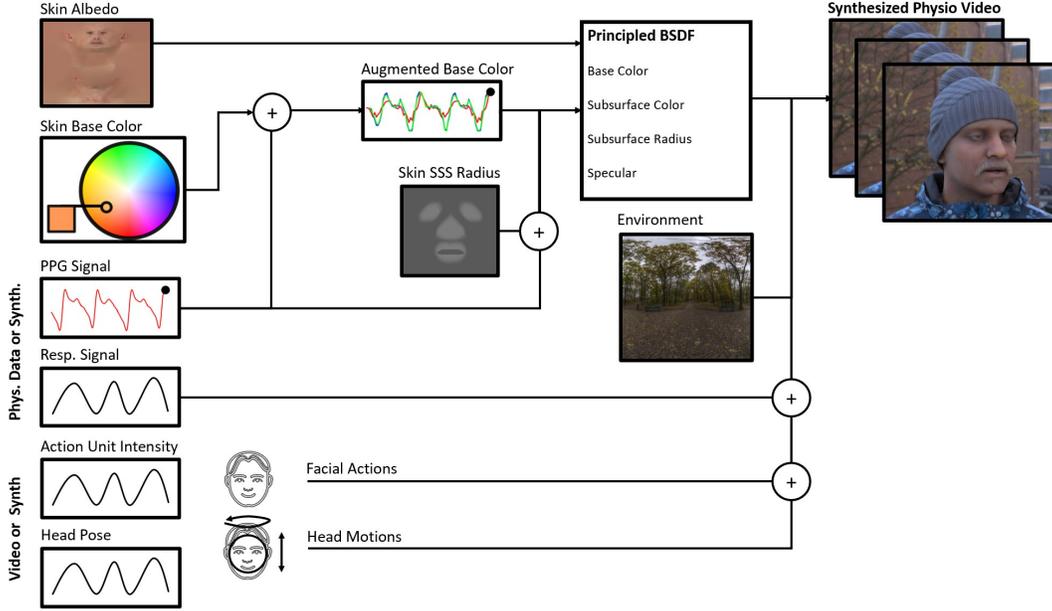
Figure 2: The synthetic videos were created using a graphics pipeline. We use a model of facial blood flow by adjusting properties of the physically-based shading material we use for the skin and breathing by controlling the motion of the head and torso. Facial actions and head motions are added to create realism and variability.

and the temporal magnitude of changes in pixels. These augmentations help in the learning process, ultimately leading to the model learning better representations.

Wood et al. [50] recently presented a sophisticated facial synthetics pipeline that produced high-fidelity data. They were able to successfully train state-of-the-art landmark localization and face parsing models. However, creating high fidelity 3D assets for simulating many different facial appearances (e.g., bone structures, facial attributes, skin tones etc.) is time consuming and expensive. The data that these pipelines can create will then not necessarily be available broadly to researchers. Therefore, in this paper we present a new dataset (SCAMPS) of high fidelity synthetic human simulations that will be made publicly available. These data are designed for the purposes of training and testing camera physiological measurement methods. To summarize our contributions: 1) We present the first public synthetic dataset for camera physiological measurement. 2) These data include precisely synchronized multi-parameter physiological ground-truth waveforms (cardiac, breathing) alongside facial action and head pose. 3) Results illustrating baseline performance training on the SCAMPS dataset and testing on two public datasets (UBFC-rPPG [4] and MMSE-HR [53]). We hope that this dataset allows researchers to explore the potential for synthetics in the domain of camera physiological measurement, including but not limited to: addressing the simulation-to-real (sim2real) generalization gap, how to leverage very precisely aligned segmentation maps and physiological waveforms for learning models, multimodal learning combining estimation of physiological (e.g., HR) and behavioral (e.g., AUs) signals, and using synthetic data to help address bias in camera physiological measurement models.

## 2 Camera Physiological Measurement

Camera measurement of physiological signals involves analysis of subtle changes in light reflected from the body. In videos, the photoplethysmographic signal manifests as small skin pixel color changes over time. The breathing signal is observed as motion, particularly prominent around the chest. Blazek, Wu and Hoelscher [3] proposed the first imaging system for measuring cardiac signals. This computer-based CCD near-infrared (NIR) imaging system provided evidence that peripheral blood volume could be measured without contact using an imager. Successful replications of these experiments cemented the concept [49, 39, 47]. Applying machine learning tools and knowl-
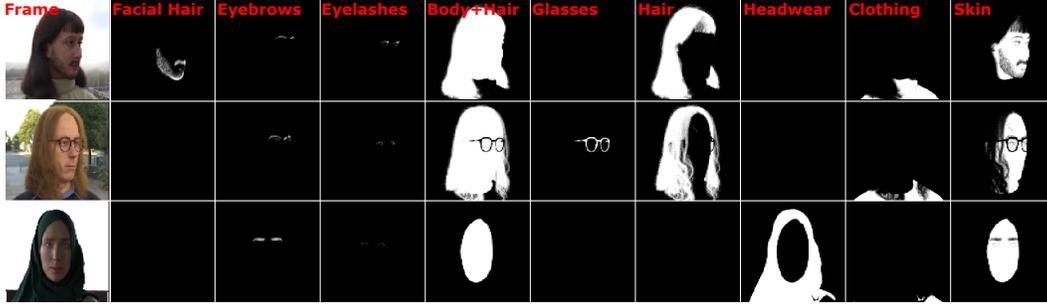
Figure 3: Each RGB frame is accompanied by segmentation masks for beard, eyelashes, eyebrows, glasses, hair, skin and clothing.

edge of physiological principals has helped to create more robust measurement methods [30, 48, 11]. With supervised methods, data soon becomes a limiting factor [36, 5, 35, 20, 52]. The significance of training data is increasing as large parameter models illustrate the potential for representation learning [52]. Work on body motion analysis from video, has found that to be a rich source of physiological information. enabling the recovery of breathing [40] and cardiac signals [2]. These methods do not require light to penetrate the skin but rather use optical flow and other motion tracking methods to measure, usually very small, motions. These subtle changes are easily swamped by larger body motions and facial expressions. Therefore, an algorithm needs to learn to successfully separate the sources from pixel changes both spatially and temporally. If we subscribe to the results of machine learning research, is likely that supervised models can learn to separate signals more effectively than handcrafted rules. For more comprehensive overviews of video physiological measurement see Chen et al. [6], Shao et al. [32] and McDuff [21].

## 3   Waveform Synthesis

Our synthesis pipeline starts with a module for generating the underlying physiologic and behavioral signals. These signals are then used to drive those properties of the synthetic humans providing precisely synchronized ground-truth labels.[1] Examples of the generated waveforms can be found in Fig. 4. To create physiological waveforms with variability we sampled several waveform parameters, such as heart rate variability standard deviation of NN intervals (HRV SDNN), relative amplitude of the systolic and dicrotic waves and the delay between the systolic and dicrotic waves from a set of uniform distributions. The bounds used for each of these parameters are specified below.

**Inter-beat Interval, PPG, ECG Waveforms.** The PPG and ECG signals were created to have the same underlying beat sequence. We first sample the beat sequence based on a heart rate (HR) frequency sampled uniformly from 40 to 150 beats/min. Heart rate variability is simulated by adding random perturbations to the beat timings. The standard deviation of these perturbations reflects the standard deviation of NN intervals (SDNN) and was sampled uniformly from 0.05 seconds to 8/HR seconds. We observed that it was important for the upper bound to be proportional to the heart rate (or mean NN interval) to create realistic variability.

For the purposes of this simulation, the morphology of the ECG wave is not relevant (e.g., we do not try to simulate a realistic QRS complex), only the timing. Thus, the ECG waveform is constructed as a time delayed series of impulses based on the NN intervals. We provide the interbeat intervals directly so that no peak detection is required for the ground-truth waveforms.

Given the beat timings and pulse arrival time (PAT) the PPG wave was then composed of a forward wave and dicrotic wave. The forward wave is created by convolving a Gaussian window with the beat impulse sequence. The leading slope of the dicrotic wave is created by convolving a Gaussian with a time lagged copy of the beat impulse sequence, the trailing slope is generated by performing the same convolution with a decaying exponential in place of the Gaussian window.

---

[1]It is important to note that the purpose of our waveform synthesis approach was not to create signals derived from a true physical model of arterial hemodynamics and tissue perfusion, but instead to develop a simple and efficient way to generate physiologically plausible waveforms.

These waves are then summed together with a dicrotic amplitude factor. The forward and dicrotic waves are then superimposed, with parameterized attenuation of the dicrotic wave relative to the forward wave, to create a physiologically plausible PPG waveform.

This signal was then low pass filtered to clean up the edges of the Gaussians, using a filter cut-off frequency of 8 Hz. Finally, the signal was normalized to give a signal of maximum amplitude of 1. This process creates PPG waveforms with the characteristic profile of systolic peaks and smaller diastolic peaks or inflections, but also with variability in the form. Finally, a small baseline drift at the breathing frequency is applied to the PPG signal to capture the subtle variations observed with breathing.

**Breathing Waveforms.** Each breathing waveform was created using sequence of breathing times based on a breathing frequency sampled from 8 to 24 breaths/min. A Gaussian window was convolved with the resulting impulse sequence. This signal was then low pass filtered to clean up the edges of the Gaussians, using a filter cut-off frequency of 8 Hz. Finally, the signal was normalized to give a signal of maximum amplitude of 1.

**Facial Actions, Blinking and Head Pose.** Unlike the physiologic waveforms, facial actions (with the exception of perhaps blinking) are rarely periodic. Therefore, we adopt an event based model [42]. For each facial action the event signal was created by a set of ramped step functions. The minimum and maximum event durations were 1 and 4 seconds, respectively. Blinking was treated separately from the other facial actions as the behavior is relatively more frequent and repetitive. For blinks the min and max event durations were 0.3 and 1 second respectively.

In each video we generate action unit "events". The start time and duration since previous event govern when the events onset and the gap between two events of the same action unit. These were are sampled from uniform distributions with bounds [0.3, 18] seconds and [1, 18] seconds, respectively. As such, in videos with action unit events there are examples of the onset and offset of most actions, some multiple times. Because facial actions are sparse but blinking occurs frequently, we generated all videos with blinking (eyes closed) events but only a subset of videos with facial actions, more details are provided below.

## 4  Video Synthesis

**Identity.** We use a texture map transferred from a high-quality 3D facial scan as the albedo of the material for creating each face. These texture maps are sampled from a set of 511 facial scans of subjects including a range of skin types/tones, genders and ages. As only varying the blood flow signal in the skin is important for our use case the facial hair is removed from these textures by an artist. Then the skin properties can be easily manipulated. Hair (and clothing) are added back in later to create the final appearance. We want the renders to display both diffuse and specular reflection effects, the diffuse reflection is handled as described below when we simulate blood flow and the specular reflection is controlled with an artist-created roughness map. Specular reflections make some parts of the face (e.g. the lips) shinier than others.

**Photoplethysmography.** Changes in diffuse reflection due to blood flow are achieved by varying the surface color and subsurface scattering of the skin texture map. We simulate blood flow by adjusting properties of the physically-based shading material we use for the face. The synthesized PPG waveform is used to drive the temporal changes. We manipulate skin tone changes using the subsurface color parameters. The weights for this are derived from the absorption spectrum of hemoglobin and typical frequency bands from an exemplar digital camera[2] (Red: 550-700 nm, Green: 400-650 nm, Blue: 350-550 nm). We manipulate the subsurface radius for the channels to capture the changes in scattering as the blood volume varies within the skin. A subsurface scattering radius texture is used to spatially-weight these and simulate variations in the thickness of the skin across the face. The same relative weighting of the RGB channels (0.36, 0.41, 0.23) is used for the BSDF subsurface radii. In absence of a more complex temporal-spatial model, we vary the parameters across the skin pixels in the same way across all frames. We recognize this is unlikely to be optimal, but does limit blood flow changes to skin pixels. We hope to be able to introduce a more realistic spatial variation in future. We used relative subsurface scattering coefficients of 0.36 (+/- 0.1), 0.41 (+/- 0.1) and 0.23 (+/- 0.1) for the red, green and blue channels respectively.

---

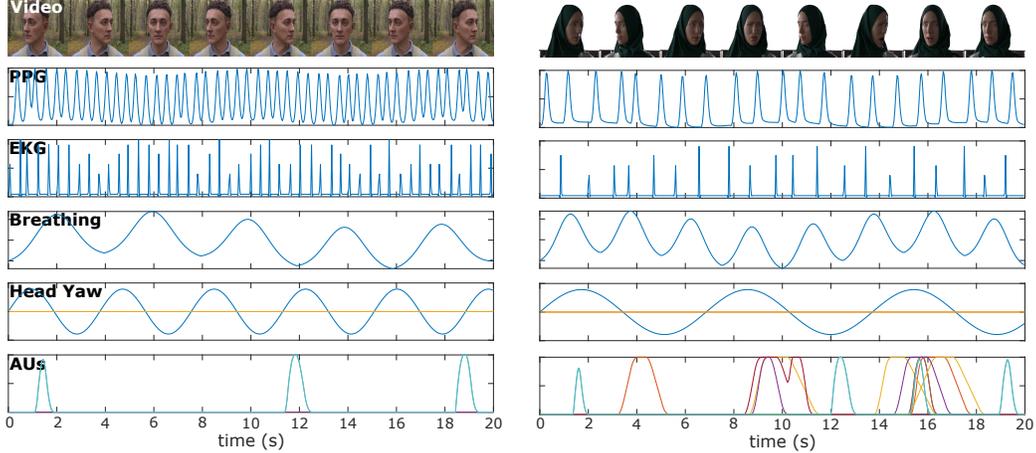[2]https://www.bnl.gov/atf/docs/scout-g_users_manual.pdf

Figure 4: Our synthetic videos are accompanied by frame-level PPG, pseudo ECG/interbeat intervals, breathing, head pose and action unit labels. Here we show examples of two videos with a subset of video frames for reference.

**Breathing.** Inhaling and exhaling cause motions of the head and chest. To capture this in the avatars we use an approximation by controlling pitch of the chest and head using the synthesized breathing input signal. The amplitude of the head and chest motions were subtle and when combined with the head rotations and facial expressions are often difficult to see; however, prior validation has shown the models trained on similar synthesized data can generalise to real videos.

**Facial Actions** Facial expressions are controlled using blendshapes that map approximately to 10 facial action units [9]: outer brow raise (AU2), brow lowerer (AU4), lid tightener (AU7), lip corner puller (AU12), lip corner depressor (AU15), chin raiser (AU17), lip puckerer (AU18), jaw drop (AU26), mouth stretch (AU27) and eyes closed (AU43). The facial action coding system is a widely used and relatively objective method for quantifying facial movements. The goal of controlling these actions is to create upper and lower facial motions. We recognize that the behaviors do not necessarily simulate realistic talking or expressions, as the dynamics of these are difficult to simulate.

## 5 Dataset

We created a dataset of 2,800 video sequences. Each video has frame level ground-truth labels for PPG, inter-beat (RR) intervals, breathing waveform, breathing intervals and 10 facial actions. We also provide video level ground-truth labels for HRV SDNN, r-peak pulse arrival time (rPAT) and dicrotic wave amplitude. These parameters were used to generate a set of 20 second PPG waveforms at 300Hz. Finally, action unit intensities were generated. The ground-truth metrics are provided as both MAT and CSV files. Each video was then rendered using the corresponding waveforms and action unit intensities, and randomly sampled appearance properties, including skin texture, hair, clothing and environment.

Figure 6 shows the distribution of heart rates, HRV SDNNs, dicrotic wave amplitudes and breathing rates in the dataset. HR, rPAT and dicrotic wave amplitudes were sampled uniformly. HRV SDNN was not sampled uniformly, as qualitatively large HRV values, while interesting, could create quite extreme differences in interbeat intervals and we deemed it appropriate to create more examples with smaller variability.

To create a dataset that can be used for training and testing under a diverse range of conditions we synthesized videos while systematically changing different confounders: 1) head motions, 2) facial actions, and 3) dynamic illumination. A training, validation and test split of the data is provided on our project page as is a file indicating which confounders are present in each video. As each video was sampled with a different combination of appearance parameters, they all contain avatars with different appearance. However, some avatars may look similar if they have the same skin texture and hair style. Figure 1 and 5 both show a collage of frames from different videos illustrating the diversity in appearance. The video frame (RGB) come with segmentation maps (see Fig. 3) that

6

Figure 5: Example frames from the SCAMPS dataset showing the diversity in avatar appearance, behavior and environment.
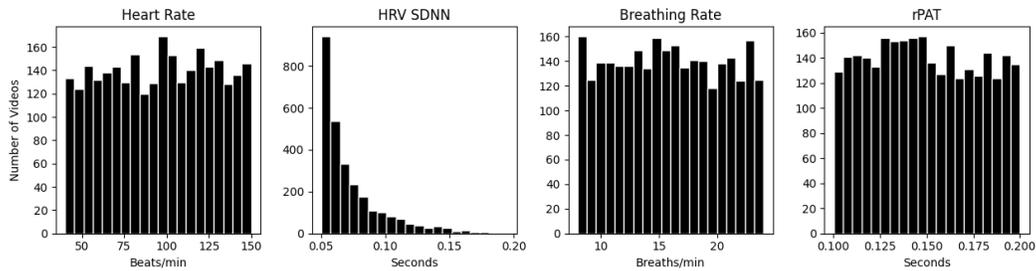


Figure 6: Examples of the distribution of heart rates, HRV SDNNs, breathing rates and dicrotic wave amplitudes in the SCAMPS dataset. An advantage of synthetic data pipelines is the ability to create a wide range of examples with specific distributions.

provide pixel level labels for beard, eyelashes, eyebrows, glasses, hair, skin and clothing. This is important as we know that the PPG signal will not be present in material that do not have blood flow (e.g., hair, clothing) and so we expect any supervised learning method to learn to segment skin as one of the operations. Therefore, we anticipate that segmentation maps will be useful to the community, both in training and in testing camera PPG methods.

**Head Motions.** Two thousand videos have rotation head motions and 800 have no head motion. Of the videos with head rotations, 1200 have smooth rotation (400 videos at 10, 20 and 30 degrees per second) and a further 800 have non-smooth head rotations in which the head was randomly positioned every second to a different angle. Ground-truth head angles are provided in the label files.

**Facial Actions.** Half of the videos (1,400) have facial actions generated with the event model described above, the other half have no facial actions. This enables training and/or testing systematically introducing the confounder of facial motions on the physiological measurement. The sequences and combinations of facial actions in each video were randomly sampled and therefore some of the facial expressions can look unnatural; however, this does provide a relatively dense set of examples of facial action onsets and offsets. We contrast this to many facial expression datasets in which facial actions are relatively sparse. We felt that more examples would generally be more useful for training models.

**Background Motion and Dynamic Illumination.** A set of 400 of the videos have dynamic illumination and background motion created by simulating the subject turning around in the environment. Half of these 400 videos have facial actions and half have head motions in addition to the background motion.

7

Table 2: Cross-dataset heart rate evaluation on UBFC and MMSE-HR (beats per minute).

| Method | UBFC [4] | | | MMSE-HR [53] | | |
|---|---|---|---|---|---|---|
| | MAE↓ | RMSE↓ | ρ↑ | MAE↓ | RMSE↓ | ρ↑ |
| DeepPhys[5] (trained on SCAMPS) | 5.42 | 13.1 | 0.72 | 4.59 | 8.89 | 0.81 |
| POS[48] | 3.52 | 8.38 | 0.90 | 3.90 | 9.61 | 0.78 |
| CHROM[8] | 3.10 | 6.84 | 0.93 | 3.74 | 8.11 | 0.82 |
| ICA[30] | 4.39 | 11.60 | 0.82 | 5.44 | 12.0 | 0.66 |

MAE = Mean Absolute Error in HR estimation, RMSE = Root Mean Square Error in HR estimation, ρ = Pearson Correlation in HR estimation.

## 6 Baselines

One might ask the question "how well does a model trained on synthetic data generalize to real videos?" While there is some precedent for using synthetics for heart and breathing rate estimation [23, 24], those works did not use the SCAMPS dataset. To illustrate how this specific dataset can be used for video physiological measurement and provide initial baseline results, we performed experiments training with the SCAMPS dataset and testing on two public benchmark video datasets. The code and resulting trained models used to generate these results can be found on our project page.

**Model.** Our goal here is not to provide an exhaustive list of results on different model architectures, but a representative baseline for researchers to compare to. We do not argue that this is the current state-of-the-art but rather is a reasonable starting point for future research with synthetic data in the field of camera physiological measurement. We implemented DeepPhys [5] as the baseline supervised model due to its relative simplicity. We trained on frames with resolution 72x72 pixels. First, we cropped the center 240x240 pixel region of each 320x240 pixel raw images. We then down sample these to 72x72 using a bilinear downsampling method. Difference frames were computed by performing a difference operation on successive frames. The resulting appearance and difference frames were normalized consistent with the method in Chen and McDuff [5]. These frames are then used for training the supervised model. We used a learning rate of 0.0001 and the ADAM optimizer. We trained the model using videos from the SCAMPS training set for 10 epochs. We validated on a real video dataset (PURE [37]) as the testing sets are also real videos. The model from the epoch with lowest mean absolute error (MAE) heart rate estimation was selected and then we evaluated this model on the test sets. A Butterworth filter was applied to all model outputs (cut-off frequencies of 0.7 and 2.5 Hz) before computing the frequency spectra and heart rate.

**Results.** The results reported here are on the UBFC-rPPG [4] and MMSE-HR [53] datasets. Table 2 shows the mean absolute error (MAE), root mean squared error (RMSE) and correlation ($\rho$) in heart rate estimation compared to the gold-standard measures from each of the datasets. The results on both datasets show that the synthetic data are sufficient to train a reasonable supervised model. The trained model does not necessarily exceed the performance of the existing unsupervised methods and is in some cases a little worse. However, as first baselines these numbers do demonstrate that generalization from synthetic video to real ones is possible and also that there is room for improvement. By releasing the SCAMPS dataset we hope that researchers can design methods that bridge the sim-to-real gap that exists.

## 7 Access and Usage

The data may be used for research purposes and any images from the dataset can be used in academic publications. Researchers may redistribute the SCAMPS dataset, so long as they include all credit or attribution information and that the terms of redistribution require any recipient to do the same. The license agreement details the permissible use of the data and the appropriate citation, it can be found at: https://github.com/danmcduff/scampsdataset. Use of the dataset for commercial purposes is strictly prohibited, although research use at commercial companies is permissible. The authors commit to maintaining the dataset and ensuring access is available to the research community.

Some of our rendered faces may be close in appearance to the faces of real people. Any such similarity is naturally unintentional, as it would be in a dataset of real images, where people may appear

similar to others unknown to them. As such there is no personally identifiable data or biometrics contained within the data, but the authors bear responsibility in case of any violation of rights that might occur.

## 8 Transparency and Broader Impacts

This dataset was created for research and experimentation on camera measurement of physiological signals. While the dataset is useful for testing models, was not designed as a test set for evaluating the clinical efficacy of a model, just because a model performs well on synthetic data does not mean it will generalize to videos of real people. The SCAMPS dataset was not designed for computer vision tasks such as face recognition, gender recognition, facial attribute recognition, or emotion recognition. We do not believe this dataset would be suitable for these applications without further validation.

We have tried to make this dataset representative of a diverse population. However, it still does not capture a uniform distribution of skin types and other appearance characteristics. We are working on addressing these limitations. When using this dataset, as with others, one should be careful to pay attention to biases that might exist. Please see the SCAMPS dataset datasheet [10] included in the supplementary material and linked from our project page for more details.

An advantage of camera physiological measurement is that contact with the body is not required and that cameras are ubiquitous sensors. However, these advantages can cause problems. Unobtrusive measurement from small, ubiquitous sensors makes measurement without a subject's knowledge simpler. It is important that norms and regulations that govern on-body physiological measurement devices are extended to camera measurement systems. Consent should always be obtained from subjects before measuring physiologic data of this kind.

## 9 Conclusions

The SCAMPS dataset contains high-fidelity simulations designed for training and testing camera-based physiological sensing algorithms. The dataset was designed to capture a diverse range of appearances, environments and lighting conditions. Synchronized ground-truth signals include interbeat and breath intervals and PPG, ECG and breathing waveforms precisely aligned with each video frame. Facial actions, blinking and head pose labels are also provided. Benchmark experiments show that it is possible to train models only with these synthetic data that generalize to real videos. We hope that this dataset helps support research towards more robust and fair vision-based physiological sensing models.

## References

[1] Y. Ba, Z. Wang, K. D. Karinca, O. D. Bozkurt, and A. Kadambi. Overcoming difficulty in obtaining dark-skinned subjects for remote-ppg by synthetic augmentation. *arXiv preprint arXiv:2106.06007*, 2021.

[2] G. Balakrishnan, F. Durand, and J. Guttag. Detecting pulse from head motions in video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3430–3437, 2013.

[3] V. Blazek, T. Wu, and D. 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.

[4] S. Bobbia, R. Macwan, Y. Benezeth, A. Mansouri, and J. Dubois. Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters*, 124:82–90, 2019.

[5] W. Chen and D. McDuff. Deepphys: Video-based physiological measurement using convolutional attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 349–365, 2018.

[6] X. Chen, J. Cheng, R. Song, Y. Liu, R. Ward, and Z. J. Wang. Video-based heart rate measurement: Recent advances and future prospects. *IEEE Transactions on Instrumentation and Measurement*, 68(10):3600–3615, 2018.

[7] A. Dasari, S. K. A. Prakash, L. A. Jeni, and C. S. Tucker. Evaluation of biases in remote photoplethysmography methods. *NPJ digital medicine*, 4(1):1–13, 2021.

[8] G. De Haan and V. Jeanne. Robust pulse rate from chrominance-based rppg. *IEEE Transactions on Biomedical Engineering*, 60(10):2878–2886, 2013.

[9] P. Ekman, W. V. Friesen, and J. Hager. *Facial action coding system: A technique for the measurement of facial movement*. Research Nexus, Salt Lake City, UT, 2002.

[10] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92, 2021.

[11] A. Gudi, M. Bittner, R. Lochmans, and J. van Gemert. Efficient real-time camera based estimation of heart rate and its variability. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[12] R. M. Haralick. Performance characterization in computer vision. In *BMVC92*, pages 1–8. Springer, 1992.

[13] G. Heusch, A. Anjos, and S. Marcel. A reproducible study on remote heart rate measurement. *arXiv preprint arXiv:1709.00962*, 2017.

[14] G.-S. Hsu, A. Ambikapathi, and M.-S. Chen. Deep learning with time-frequency representation for pulse estimation from facial videos. In *2017 IEEE international joint conference on biometrics (IJCB)*, pages 383–389. IEEE, 2017.

[15] M. Kopeliovich and M. Petrushan. Color signal processing methods for webcam-based heart rate evaluation. In *Proceedings of SAI Intelligent Systems Conference*, pages 703–723. Springer, 2019.

[16] A. Kortylewski, B. Egger, A. Schneider, T. Gerig, A. Morel-Forster, and T. Vetter. Empirically analyzing the effect of dataset biases on deep face recognition systems. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2093–2102, 2018.

[17] A. Kortylewski, B. Egger, A. Schneider, T. Gerig, A. Morel-Forster, and T. Vetter. Analyzing and reducing the damage of dataset bias to face recognition with synthetic data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.

[18] X. Li, I. Alikhani, J. Shi, T. Seppanen, J. Junttila, K. Majamaa-Voltti, M. Tulppo, and G. Zhao. The obf database: A large face video database for remote physiological signal measurement and atrial fibrillation detection. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 242–249. IEEE, 2018.

[19] X. Liu, J. Fromm, S. Patel, and D. McDuff. Multi-task temporal shift attention networks for on-device contactless vitals measurement. *NeurIPS*, 2020.

[20] H. Lu, H. Han, and S. K. 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.

[21] D. McDuff. Camera measurement of physiological vital signs. *arXiv preprint arXiv:2111.11547*, 2021.

[22] D. McDuff, R. Cheng, and A. Kapoor. Identifying bias in ai using simulation. 2018.

[23] D. McDuff, J. Hernandez, E. Wood, X. Liu, and T. Baltrusaitis. Advancing non-contact vital sign measurement using synthetic avatars. *arXiv preprint arXiv:2010.12949*, 2020.

[24] D. McDuff, X. Liu, J. Hernandez, E. Wood, and T. 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.

[25] R. Meziatisabour, Y. Benezeth, P. De Oliveira, J. Chappe, and F. Yang. Ubfc-phys: A multi-modal database for psychophysiological studies of social stress. *IEEE Transactions on Affective Computing*, 2021.

[26] X. Niu, H. Han, S. Shan, and X. Chen. Vipl-hr: A multi-modal database for pulse estimation from less-constrained face video. *arXiv preprint arXiv:1810.04927*, 2018.

[27] E. M. Nowara, T. K. Marks, H. Mansour, and A. Veeraraghavan. Sparseppg: Towards driver monitoring using camera-based vital signs estimation in near-infrared. In *Computer Vision and Pattern Recognition (CVPR), 1st International Workshop on Computer Vision for Physiological Measurement*, 2018.

[28] E. M. Nowara, D. McDuff, and A. Veeraraghavan. Combining magnification and measurement for non-contact cardiac monitoring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3810–3819, 2021.

[29] A. Pai, A. Veeraraghavan, and A. Sabharwal. Camerahrv: robust measurement of heart rate variability using a camera. In *Optical Diagnostics and Sensing XVIII: Toward Point-of-Care Diagnostics*, volume 10501, page 105010S. International Society for Optics and Photonics, 2018.

[30] M.-Z. Poh, D. McDuff, and R. W. Picard. Advancements in noncontact, multiparameter physiological measurements using a webcam. *IEEE transactions on biomedical engineering*, 58(1):7–11, 2010.

[31] W. Qiu and A. Yuille. Unrealcv: Connecting computer vision to unreal engine. In *European Conference on Computer Vision*, pages 909–916. Springer, 2016.

[32] D. Shao, C. Liu, and F. Tsow. Noncontact physiological measurement using a camera: A technical review and future directions. *ACS sensors*, 6(2):321–334, 2020.

[33] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *CVPR 2011*, pages 1297–1304. Ieee, 2011.

[34] M. Soleymani, J. Lichtenauer, T. Pun, and M. Pantic. A multimodal database for affect recognition and implicit tagging. *IEEE transactions on affective computing*, 3(1):42–55, 2011.

[35] R. Song, H. Chen, J. Cheng, C. Li, Y. Liu, and X. Chen. Pulsegan: Learning to generate realistic pulse waveforms in remote photoplethysmography. *IEEE Journal of Biomedical and Health Informatics*, 25(5):1373–1384, 2021.

[36] R. Špetlík, V. Franc, and J. Matas. Visual heart rate estimation with convolutional neural network. In *Proceedings of the british machine vision conference, Newcastle, UK*, pages 3–6, 2018.

[37] R. Stricker, S. Müller, and H.-M. 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.

[38] L. Świrski and N. 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.

[39] C. Takano and Y. Ohta. Heart rate measurement based on a time-lapse image. *Medical engineering & physics*, 29(8):853–857, 2007.

[40] L. Tarassenko, M. Villarroel, A. Guazzi, J. Jorge, D. Clifton, and C. Pugh. Non-contact video-based vital sign monitoring using ambient light and auto-regressive models. *Physiological measurement*, 35(5):807, 2014.

[41] H. E. Tasli, A. Gudi, and M. Den Uyl. Remote ppg based vital sign measurement using adaptive facial regions. In *2014 IEEE international conference on image processing (ICIP)*, pages 1410–1414. IEEE, 2014.

[42] T. Vandal, D. McDuff, and R. El Kaliouby. Event detection: Ultra large-scale clustering of facial expressions. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 1, pages 1–8. IEEE, 2015.

[43] D. Vazquez, A. M. Lopez, J. Marin, D. Ponsa, and D. Geronimo. Virtual and real world adaptation for pedestrian detection. *IEEE transactions on pattern analysis and machine intelligence*, 36(4):797–809, 2014.

[44] V. Veeravasarapu, R. N. Hota, C. Rothkopf, and R. Visvanathan. Model validation for vision systems via graphics simulation. *arXiv preprint arXiv:1512.01401*, 2015.

[45] V. Veeravasarapu, R. N. Hota, C. Rothkopf, and R. Visvanathan. Simulations for validation of vision systems. *arXiv preprint arXiv:1512.01030*, 2015.

[46] V. Veeravasarapu, C. Rothkopf, and V. Ramesh. Model-driven simulations for deep convolutional neural networks. *arXiv preprint arXiv:1605.09582*, 2016.

[47] W. Verkruysse, L. O. Svaasand, and J. S. Nelson. Remote plethysmographic imaging using ambient light. *Optics express*, 16(26):21434–21445, 2008.

[48] W. Wang, A. C. den Brinker, S. Stuijk, and G. de Haan. Algorithmic principles of remote ppg. *IEEE Transactions on Biomedical Engineering*, 64(7):1479–1491, 2017.

[49] F. P. Wieringa, F. Mastik, and A. F. van der Steen. Contactless multiple wavelength photoplethysmographic imaging: A first step toward "spo 2 camera" technology. *Annals of biomedical engineering*, 33(8):1034–1041, 2005.

[50] E. Wood, T. Baltrusaitis, C. Hewitt, S. Dziadzio, T. J. Cashman, and J. 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.

[51] E. Wood, T. Baltrusaitis, X. Zhang, Y. Sugano, P. Robinson, and A. 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.

[52] Z. Yu, Y. Shen, J. Shi, H. Zhao, P. Torr, and G. Zhao. Physformer: Facial video-based physiological measurement with temporal difference transformer. *arXiv preprint arXiv:2111.12082*, 2021.

[53] Z. Zhang, J. M. Girard, Y. Wu, X. Zhang, P. Liu, U. Ciftci, S. Canavan, M. Reale, A. Horowitz, H. Yang, et al. Multimodal spontaneous emotion corpus for human behavior analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3438–3446, 2016.