

EfficientPhys: Enabling Simple, Fast and Accurate Camera-Based Cardiac Measurement

Xin Liu¹, Brian Hill², Ziheng Jiang¹, Shwetak Patel¹, Daniel McDuff³

Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, WA, USA¹

Department of Computer Science, University of California, Los Angeles, CA, USA²

Microsoft Research, Redmond, WA, USA³

{xliu0, shwetak, ziheng}@cs.washington.edu

brian.l.hill@cs.ucla.edu, damcduff@microsoft.com

Abstract

Camera-based physiological measurement is a growing field with neural models providing state-of-the-art performance. Prior research has explored various “end-to-end” architectures; however these methods still require several preprocessing steps and are not able to run directly on mobile and edge devices. The operations are often non-trivial to implement, making replication and deployment difficult and can even have a higher computational budget than the “core” network itself. In this paper, we propose two novel and efficient neural models for camera-based physiological measurement called EfficientPhys that remove the need for face detection, segmentation, normalization, color space transformation or any other preprocessing steps. Using an input of raw video frames, our models achieve strong accuracy on three public datasets. We show that this is the case whether using a transformer or convolutional backbone. We further evaluate the latency of the proposed networks and show that our most lightweight network also achieves a 33% improvement in efficiency.

1. Introduction

Camera-based physiological measurement is a non-contact approach for capturing cardiac signals via light reflected from the body [20]. The most common such signal is the blood volume pulse (BVP) measured via the photoplethysmogram (PPG). From this, heart rate [33, 36], respiration rate [26] and pulse transit times [29] can be derived. Furthermore, there is promising evidence that the PPG signals can be used to measure signs of arterial disease [34]. Neural models are the current state-of-the-art for camera PPG measurement [5, 10, 16]. These networks can learn strong feature representations and effectively disentangle the subtle changes in pixels due to underlying physiological

processes from those due to body motions, lighting changes and other sources of “noise”.

While prior research has framed architectures as “end-to-end” methods, those that achieve state-of-the-art performance actually require several preprocessing steps before data is used as input to the network. For example, [5] and [10] use hand-crafted normalized difference frames and normalized appearance frames as input to their convolutional attention network. [24] and [19] use a complex schema to create feature maps called “MSTmaps”, their process includes facial landmark detection, extraction of several regions of interest (ROI) using these landmarks, and then averaging pixel values in both the RGB and YUV color spaces.

These preprocessing steps have several drawbacks: 1) They make assumptions about optimal normalization or representation without allowing the network to learn these features in a data-driven manner. 2) They are computationally costly and in many cases add a significant number of operations to the video processing pipeline. There are several reasons why running camera-based physiological sensing on-device is desirable: privacy preservation, the ability to use raw (i.e., uncompressed) video and data cost and bandwidth savings. Any additional computation needs to be justified by improving model accuracy. Moreover, since camera-based physiological sensing is a privacy-sensitive application, it is preferable to store the data on local devices instead of streaming both video and physiological data to the cloud. The overhead from processing is not acceptable if we aim to make the system accessible to low-end mobile devices. 3) Many of these steps are non-trivial to implement and optimize in and of themselves. This makes it harder to deploy real-time systems and to replicate the implementation on different platforms. For instance, implementing existing methods on Android, iOS, or in JavaScript requires a significant amount of effort. Some libraries, such as facial landmark detection, are not even available on every platform.

Thus, the last mile engineering using the existing methods becomes especially challenging.

Ideally, a video-based physiological measurement method would be able to run at a high frame rate even on mobile devices, be simple to implement across different platforms, and achieve state-of-the-art performance. Addressing the aforementioned challenges would help achieve these properties. We propose a truly end-to-end network, EfficientPhys, for which the input is unprocessed video frames without requiring accurate face cropping (see Fig. 1). Due to recent advancements in visual transformers, we propose both a convolutional and visual transformer architecture and compare and contrast the performance of these two.

In summary, our key contributions are to: 1) propose two novel one-stop neural architectures, a visual transformer and a convolutional network, which do not require any preprocessing steps, 2) evaluate the proposed methods on three popular benchmark datasets, 3) evaluate on-device latency across both state-of-the-art machine learning-based approaches as well as signal processing-based techniques. To the best of our knowledge, this is the first paper that explores the visual transformer in camera-based physiological measurement and its comparison with convolutional networks. This is also the first paper exploring a completely end-to-end on-device neural architecture for mobile devices.

2. Related Work

Camera-based Vital Measurement. There is a growing community studying the use of cameras to sense physiological vitals signs [38, 36, 12]. Prior work established the fundamentals of how RGB images could be used to extract the pulse signal using signal source separation techniques (e.g., ICA) [26]. Other methods derived these parameters from physically-based models to achieve elegant and fast demixing (e.g., Plane Orthogonal-to-Skin (POS))[37]. By calculating a projection plane orthogonal to the skin-tone based on optical and physiological principles, the authors were able to achieve a stronger BVP signal-to-noise ratio (SNR).

Since the underlying relationship between the pulse and skin pixels is complex, deep convolutional neural networks have shown superior performance over the traditional source separation algorithms. DeepPhys [5] was the first paper that demonstrated that a deep neural network outperforms all the traditional signal processing approaches. Liu et al. have also proposed an on-device efficient neural architecture called MTTs-CAN for camera-based physiological sensing, which leverages a tensor-shift module and 2D-convolutional operations to perform efficient spatial-temporal modeling [10]. More recently, an adversarial learning approach, called Dual-GAN, has also been studied to learn noise-resistant mappings from video frames to pulse waveform and noise distributions [19]. With two generative-adversarial networks,

they can promote each adversarial network’s representation and further improve the feature disentanglement between pulse and various noise sources.

However, DeepPhys and MTTs-CAN both require a few preprocessing steps including calculating difference frames and performing image normalization. Dual-GAN has a even more complex preprocessing module called MSTMaps proposed by [24]. The MSTMaps are a set of multi-scale spatial temporal maps created by 1) cropping the facial region, 2) extracting facial landmarks, 3) performing average pooling for every color channel and every ROI combination for each frame, 4) generating ROI combinations using all the detected ROI regions and landmarks, 5) multiplying each item in all ROI combinations with six channels. The final size of the MSTMap is $(2^n - 1) \times T \times 6$ where T is the number of frames and n is the number of ROI regions. Such a preprocessing module not only consumes large amounts of memory but also introduces a large computational burden to the entire pipeline. Moreover, stacking all of these extra procedures makes development and deployment much more difficult. Unlike these methods, the goal of our proposed method EfficientPhys is to create a preprocessing-free neural architecture that is simple to use and deploy, efficient on mobile devices, and accurate on settings with various types of noise.

Visual Transformers. Although convolutional neural networks have been widely studied and used in many computer vision applications, vision transformers were first successfully used on the task of image classification. By training on larger datasets, vision transformer (ViT) attains excellent performance and can be used in downstream fine-tuning with fewer amounts of data [7]. More recently, the state-of-the-art Swin vision transformer was proposed as a way to construct hierarchical feature maps and improve computational efficiency by using a hierarchical representation and limiting self-attention computation to non-overlapping local windows while allowing for cross-window connection [14]. However, transformer architectures have been barely studied in the field of camera-based vitals measurement. The closest work used transformers to detect remote photoplethysmography (rPPG) for attack/spoofing detection [39]. However, this paper did not evaluate the proposed vision transformer in the task of heart rate estimation using any public datasets, which is considered as the gold-standard benchmark for the field of camera-based vital measurement. More recently, Yu *et al.* recently proposed Physformer [41] which is also a visual transformer based architecture¹. To our best knowledge, our proposed vision transformer was the first architecture in camera-based heart rate measurement with detailed evaluation on multiple public datasets.

¹Physformer was published after this manuscript was submitted for review

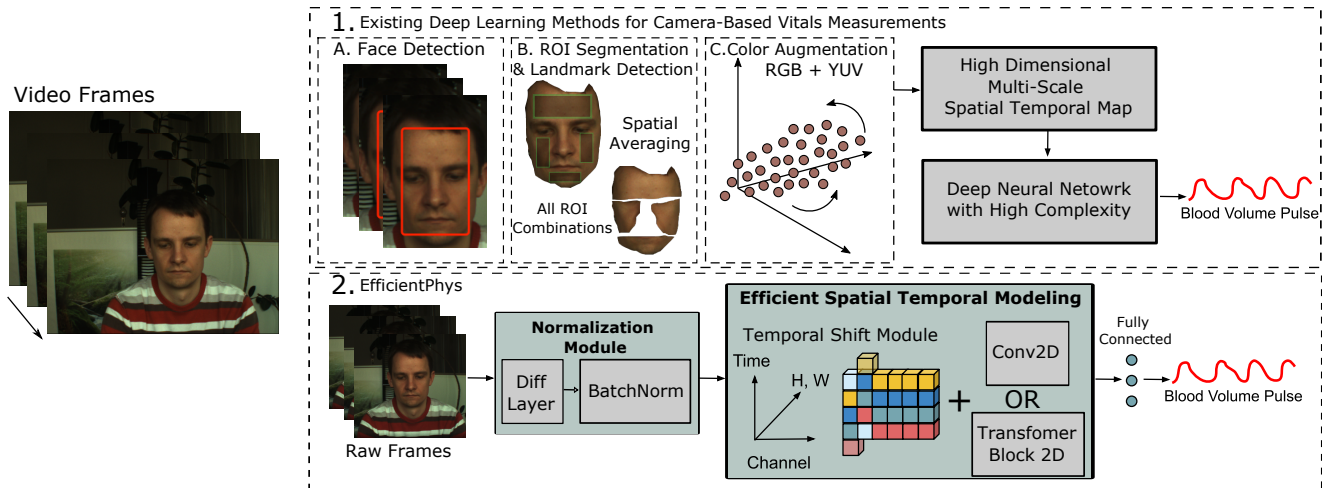


Figure 1. A high-level comparison of EfficientPhys and existing deep learning approaches for camera-based vitals measurement

3. Method

3.1. Convolution-based EfficientPhys

To enable simple, fast and accurate real-time on-device camera-based vitals measurement, we propose a one-stop solution architecture that takes raw video frames as the input to the network and outputs a first-derivative PPG signal. The convolution-based EfficientPhys is a one-branch network that contains a custom normalization layer, self-attention module, tensor-shift module and 2D convolution operation to perform efficient and accurate spatial-temporal modeling while making it simple to deploy.

Normalization Module. Existing neural methods all require different levels of preprocessing before providing the visual representation to the network to learn the underlying relationship between skin pixels and cardiac pulse signal. For instance, The state-of-the-art networks Dual-GAN [19] and CVD [24] proposed a hand-crafted spatial-temporal representations called STMaps. These preprocessed representations are generated for each video frame and includes steps of detecting 81 facial landmark points, extracting a set of region of interest (ROI) combinations ($2^n - 1$ where n is the number of ROIs, $n=6$) using these landmarks, and averaging pixel values in both the RGB and YUV color spaces, multiplying the 63 ROI combinations with the six channels. These modules not only add significant computational burden (Table 4 shows that Dual-GAN’s preprocessing module takes 275ms per frame) but also make the system more challenging to implement and deploy on real-world computing systems such as mobile devices.

One of the goals of EfficientPhys is to remove these preprocessing modules completely and provide a one-stop solution. To achieve such simplicity and deployability, we propose a custom normalization module, which can perform motion modeling between every two consecutive RGB raw

frames and normalization to reduce the lighting and motion noise. More specifically, the proposed normalization module includes a difference layer and a batchnorm layer. The difference layer (e.g., torch.diff) computes the first forward difference along the temporal axis of the raw video frames, by subtracting every two adjacent frames. Performing motion modeling between every two consecutive frames and normalization is more like a high-pass filtering and can help reduce the global noise from lighting and motion noise, while maintaining the subtle changes from PPG. To provide optical basis in our work, equation 1 illustrates the optical grounding of difference frame where $D_k(t)$ of every two consecutive frames. $I(t)$ is the luminance intensity which is modulated by the specular reflection $v_s(t)$ and the diffuse reflection $v_d(t)$ as well as the optical sensor’s quantization noise $v_n(t)$.

$$D_k(t) = (I(t) \cdot (v_s(t) + v_d(t)) + v_n(t)) - (I(t-1) \cdot (v_s(t-1) + v_d(t-1)) + v_n(t-1)) \quad (1)$$

However, difference frames could be dramatically different in scale and make it hard for the network to learn meaningful feature representations, especially when the signal of interest is hidden in subtle pixel changes along the temporal axis and noise artifacts can cause significantly larger relative changes. To address this, we add a batch-normalization (batchnorm) layer following the difference layer. Adding a batchnorm layer provides two benefits: 1) it normalizes the difference frames to the same scale within the batch during training, 2) unlike fixed normalization in previous work [5, 10], batchnorm provides two learnable parameters β and γ for scaling (to a different variance) and shifting (to a different mean) and two non-trainable parameters which are the mean μ and the standard deviation σ . Through the learning process, the batchnorm layer can learn the best parameters

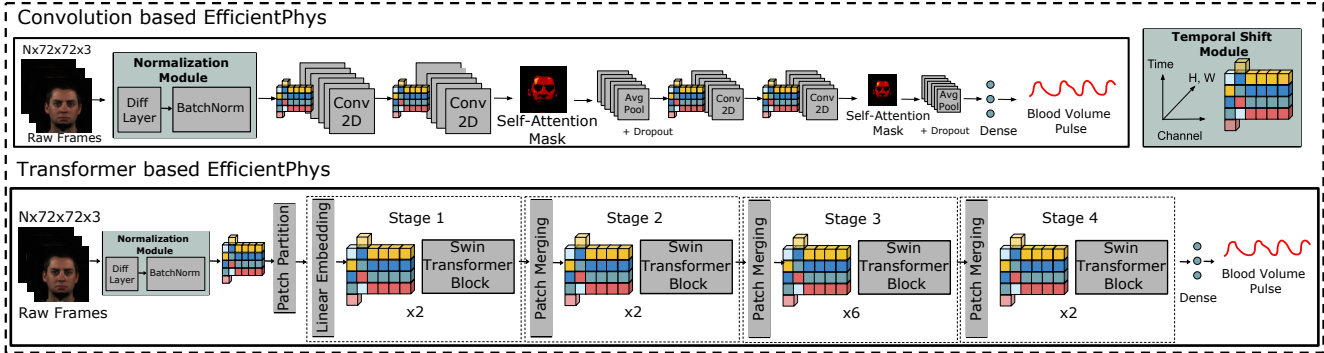


Figure 2. We present two novel architectures to enable simple, fast, and accurate camera-based vitals measurement: Convolution-based EfficientPhys and Transformer-based EfficientPhys. N is the number frames of video clip inputting to the network.

for amplifying the pixel changes while minimizing the noise as Equation 2 and Fig.3 show. Without a batchnorm layer, directly applying a difference layer means the frames appear “black”; because the subtle changes of skin pixels in every two consecutive frames are relatively very small. On the other hand, adding a follow-up batchnorm layer will help it learn the normalization function to magnify the subtle changes of skin pixels substantially. The result is not simply a magnification of values but a normalization and magnification. Moreover, we also compare the output batchnorm layer to the hand-crafted normalized frame as shown in Fig.3. The output of batchnorm layer contains more information and qualitative analysis suggests it should be a better tool for skin segmentation after the learning process.

$$N_k(t) = \frac{(\beta_t * \mathbf{D}_k(t) + \gamma_t) - \mu_{\mathbf{D}_k}}{\sigma_{\mathbf{D}_k}} \quad (2)$$

Self-Attention-Shifted Network. To efficiently capture the rich spatial-temporal information, we propose a self-attention-shifted network (SASN). SASN is built on top of the previous state-of-the-art method for on-device spatial-temporal modeling in optical cardiac measurement - temporal-shift convolutional attention network (TS-CAN) [10]. TS-CAN has two convolutional branches, one of which takes a preprocessed difference frame representation and one of which takes a normalized appearance frame. The motion branch performs the main spatial-temporal modeling and estimation, and the appearance branch provides attention masks to guide the motion branch to better isolate the pixels of interest (e.g., skin pixels). However, we argue that the attention masks do not have to be obtained through a separate appearance branch and they can be also learned with a single branch end-to-end network. As Fig. 2 illustrates, our proposed self-attention-shifted network starts with the custom normalization module discussed in the previous section and then continues with two tensor-shifted convolutional operations. After the second and fourth tensor-shifted 2D convolutional layers, we add a self-attention module to

help the network minimize the negative effects introduced by temporal shifting as well as motion and lighting noise. The self-attention layers are softmax attention layers with 1D convolutions followed by a sigmoid activation function. Then, normalization is applied to remove the outlier values in the attention mask, and the final normalized attention mask is element-wise multiplied with the output from the tensor-shifted convolution. Equation 3 summarizes how our self-attention mechanism works where $ts(\cdot)$ denotes temporal shift operation, ω_c^t denotes the 2D convolutional kernel followed by the temporal shift module, σ is the sigmoid activation and ω_a^t is the 1×1 convolutional kernel for self attention.

$$(\omega_c^t ts(N_k(t)) + b_c^t) \odot \frac{H_t W_t \cdot \sigma(\omega_a^t \mathbb{X}_\alpha^t + b_a^t)}{2 \|\sigma(\omega_a^t \mathbb{X}_\alpha^t + b_a^t)\|_1} \quad (3)$$

3.2. Transformer-based EfficientPhys

Efficient Spatial-Temporal Video Transformer. Due to the recent success of visual transformers for image and video understanding and the importance of attention mechanisms for this task [40, 5, 27, 10], we also present a visual transformer version of EfficientPhys. For this task, we need a visual transformer to learn both spatial and temporal representations. Several existing video-based visual transformers are based on 3D-embedding tokens and input all the frames into 3D encoder and spatial-temporal attention modules [1, 16]. However, the computational complexity makes these unfavourable for real-time efficient modeling on mobile devices. In the convolutional version we used tensor-shifted 2D convolutions which have been shown to achieve comparable performance as 3D convolutions [10]. Inspired by this, our proposed transformer-based EfficientPhys is based on a 2D visual transformer, Swin transformer [15], but with additional components that we will describe below.

Since the 2D Swin transformer is only able to learn spatial features that map raw RGB values to latent representa-

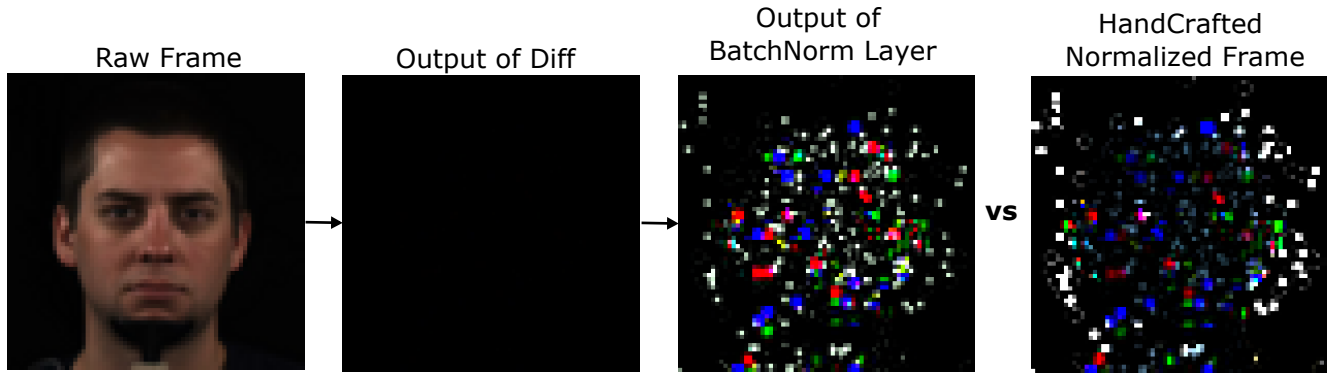


Figure 3. Outputs of diff and batchnorm layers and comparison with normalized frames generated via the hand-crafted process in prior work [5]. The output from the diff layer is almost black because the difference in skin pixels of consecutive frames is very subtle.

tions between a single frame and the target signal (pulse), it does not have ability to model temporal relationships beyond consecutive frames. One of the main contributions of the Swin transformer is the shifted window module which has linear computation complexity and allows cross-window connection by shifting the window partition and limiting self-attention computation to non-overlapping local windows. Inspired by the idea of shifting of spatial window partitions, we propose to add a tensor-shift module (TSM) [9] before every Swin transformer block to facilitate information exchange across the temporal axis. The TSM first splits the input tensor into three chunks, shifts the first chunk to the left by one place (advancing time by one frame) and shifts the second chunk to the right by one place (delaying time by one frame). All shifting operations are along temporal axis and are performed before the tensor is fed into each transformer block as shown in Fig. 2. By adding the TSM module to the Swin transformer, the new transformer architecture now has the ability to perform efficient spatial-temporal modeling and attention by combining shifting window partitions spatially and shifting frames temporally. It is worth noting that TSM does not introduce any learnable parameters, thus the proposed transformer architecture has the same number of parameters as the original Swin transformer. Finally, to enable truly end-to-end inference and learning, we also add the same normalization module proposed in the convolution EfficientPhys to this architecture.

In summary, the transformer-based EfficientPhys is the first end-to-end transformer architecture for camera-based cardiac pulse measurement that leverages tensor-shift modules and window-partition shift modules to perform efficient spatial-temporal modeling and attention to learn the underlying physiological signal from skin pixels.

4. Experiments

Training Data. To help create a robust and generalizable model for cross-dataset evaluation we use two datasets. The first is AFRL [8], which includes 300 videos from 25

subjects (17 males and 8 females). For each video, the raw resolution is 658x492 and the sampling rate of the synchronized pulse measurement is 30Hz. The dataset includes videos with a range of head motions. Every participant was instructed to remain stationary for the first two tasks, and then to perform head motions with increasing rotational velocity in the next four tasks (turning from left to right). Along with the AFRL dataset, we also leverage a synthetic avatar video dataset introduced by [21, 22] where each synthetic video is parameterized and generated with a custom pulse signal, background, facial appearance, and motion. More specifically, the input pulse signal is used to augment skin color and the subsurface radius of skin pixels to mimic the effect of the blood volume pulse on the skin’s appearance. Synthetic data such as this introduces greater diversity into the training set and has been shown to effectively help reduce disparities in performance by skin type. All the video data are resized into a resolution of 72×72 with resampling using pixel area relation ($cv2.INTER_AREA$).

Testing Data. We use three popular benchmark datasets to evaluate the accuracy of the proposed EfficientPhys. UBFC [3] is a dataset of 42 videos from 42 subjects, and the raw resolution of each video is 640x480 in an uncompressed 8-bit RGB format. The sampling for synchronized pulse signal is 30 Hz. All of the tasks collected in UBFC are stationary. MMSE [42] is a dataset including 102 videos from 40 subjects, and the raw resolution of each video is 1040x1392. The ground-truth waveform for MMSE is blood pressure signal instead of blood-volume pulse signal, and the sampling rate is 25 Hz. It is worth noting that MMSE contains a diverse distribution of skin types in Fitzpatrick scale (II=8, III=11, IV=17, V+VI=4). PURE [32] is a dataset containing 60 videos from 10 subjects. The raw resolution of each video is 640x480, and the sampling rate of the ground-truth pulse signal is 60 Hz. PURE includes a diverse set of motion tasks such as steady, talking, slow/fast translation between head movements and the camera plane, and small/medium head rotation.

Implementation & Experiment Details. We implemented both convolution-based and transformer-based EfficientPhys in PyTorch [25]. We used an AdamW optimizer to train both networks instead of Adam by introducing additional regularization to reduce the effects of over-fitting through weight decay [17]. The learning rate we used for Convolutional model was 0.001 while the rate for transformer model was 0.0001. Based on empirical studies, we used the mean squared error (MSE) loss for training the transformer models and negative Pearson loss [35] for the convolutional model. We trained both models for ten epochs with a fixed random seed. We implemented TS-CAN based on the open-sourced code [11, 10] and used the Deep Physiological Sensing Toolbox [13] for the experiments on the UBFC and PURE datasets. To calculate the performance metrics, we first applied a band-pass filter to the signal with a cutoff frequency of 0.75 and 2.5 Hz (45 beats/minute to 150 beats/minute). We then followed Dual-GAN’s evaluation scheme using peak detection and FFT to get estimated heart rate on each video of UBFC and PURE datasets [18] and MetaPhys’s evaluation scheme on MMSE [11]. We conducted video-level evaluation where we calculated an averaged heart rate for each single video. We calculated three standard metrics for each video: mean absolute error (MAE), root mean squared error (RMSE) and Pearson correlation (ρ) in heart rate estimations and the corresponding ground-truth heart rates from the blood volume pulse collected via contact oximeter sensor.

To explore the efficiency of different architectures on mobile devices, we also conducted experiments on a quad-core Cortex-A72 Raspberry Pi 4B to evaluate the model’s performance on an edge device. We performed inference 10 times to get a reliable averaged on-device inference latency for EfficientPhys and TS-CAN. Due to the lack of open-source implementation of Dual-GAN, we were only able to find the implementation of STMaps which is the preprocessing module of Dual-GAN. Thus, we only evaluated the on-device latency for the preprocessing module in Dual-GAN. We also evaluated the latency of POS, CHROM, and ICA as they are traditional signal processing methods and don’t have a separate preprocessing module.

5. Results and Discussion

EfficientPhys vs. State-of-the-Art. In Table 1 and Table 2, we present results from our proposed EfficientPhys models and the current state-of-the-art neural and signal processing methods. The learning models are all trained on the same datasets (AFRL + Synthetic) and tested on three dataset (UBFC, PURE and MMSE) to test if the model can generalize to videos with a different facial appearance, background, and lighting. To investigate how the depth of the network impacts the Transformer architecture, we created two version of Transformer-based EfficientPhys: T1 and T2.

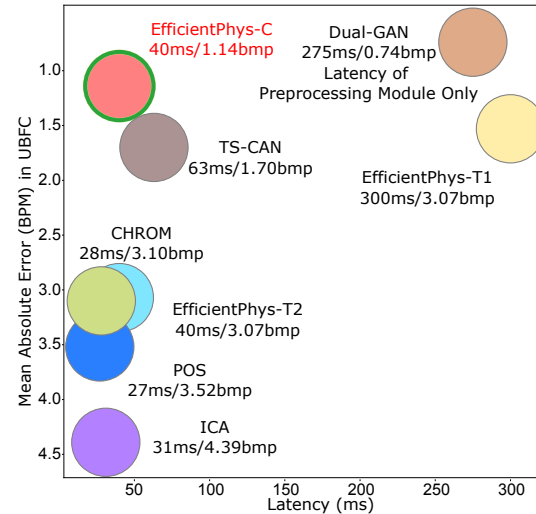


Figure 4. Accuracy-Latency Trade-off in eight different methods. Y-axis denotes the MAE error, and X-axis denotes the latency. The methods in the left-top corner have the best accuracy-latency Trade-off.

T1 uses the same depth as the Swin Transformer reported in [15] ([2, 2, 6, 2]). Each number indicates the number of Swin Transformer blocks as illustrated in Fig. 2. T2 is a more lightweight architecture to enable real-time on-device inference and has a depth of [2, 1]. EfficientPhys-C denotes the Convolution-based EfficientPhys as shown in the Fig.2. For UBFC and PURE, as Table 1 illustrates, EfficientPhys-C and EfficientPhys-T1 outperform all the existing methods. As Table 2 demonstrates, all the neural methods outperform the signal processing methods. EfficientPhys-T1 and TS-CAN achieved slightly better results than EfficientPhys-C and EfficientPhys-T2. Unfortunately, due to the lack of open source implementation or released models (e.g., Dual-GAN [19]), we could not successfully replicate their complicated model architecture and conduct cross-dataset evaluation on this comparison.

To conduct a fair comparison with the current state-of-the-art methods, we followed Dual-GAN [19] to train our models only on PURE and to test on UBFC as Table 3 shows. Although Dual-GAN outperforms all of other methods, we argue that the margin is relatively small as both Dual-GAN and EfficientPhys-C achieve a Pearson correlation of 0.99. Moreover, according to American National Standards Institute (ANSI) and Consumer Technology Association’s standard [2], MAPE of ± 5 is an acceptable error rate. Various studies have also used this standard to validate FDA approved sensors and systems [23, 30, 28, 4]. All the methods in Table 3 have met this recommended bar.

Computational Cost and On-Device Latency. Fig. 4 and the Table 4 summarize the computational cost of the existing neural methods. Again, due to the lack of open source implementation and complex algorithm design, we

Table 1. Cross-dataset heart rate evaluation on UBFC and PURE (beats per minute).

Method	UBFC [3]				PURE [32]			
	MAE↓	MAPE↓	RMSE↓	ρ ↑	MAE↓	MAPE↓	RMSE↓	ρ ↑
EfficientPhys-C	1.14	1.16%	1.81	0.99	1.33	1.71%	5.99	0.97
EfficientPhys-T1	2.08	2.53%	4.91	0.96	1.11	1.30%	5.94	0.97
EfficientPhys-T2	3.07	3.41%	4.78	0.96	2.67	3.22%	9.08	0.92
TS-CAN[10]	1.70	1.99%	2.72	0.99	2.23	2.25%	3.71	0.98
POS[37]	3.52	3.36%	8.38	0.90	1.68	1.56%	9.60	0.92
CHROM[6]	3.10	3.83%	6.84	0.93	6.23	10.04%	17.18	0.71
ICA[26]	4.39	4.30%	11.60	0.82	5.70	5.69%	18.10	0.70

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

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

Method	MMSE [42]			
	MAE↓	MAPE↓	RMSE↓	ρ ↑
EfficientPhys-C	3.48	4.02%	7.21	0.86
EfficientPhys-T1	3.04	3.91%	5.91	0.92
EfficientPhys-T2	3.51	3.96%	6.98	0.88
TS-CAN[10]	3.04	3.41%	6.55	0.89
POS[37]	3.79	4.28%	8.47	0.82
CHROM[6]	3.61	4.50%	7.43	0.85
ICA[26]	7.96	9.20%	14.02	0.51

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

Table 3. Cross dataset evaluation with models trained on PURE only and tested on UBFC (beats per minute).

Method	MAE↓	MAPE↓	RMSE↓	ρ ↑
EfficientPhys-C	2.13	2.35 %	3.00	0.99
EfficientPhys-T1	3.83	4.32%	5.62	0.87
EfficientPhys-T2	3.97	4.35%	5.91	0.94
TS-CAN[10]	1.16	1.42%	2.78	0.99
Dual-GAN[19]	0.74	0.73%	1.02	0.99
PulseGAN[31]	2.09	2.23%	4.42	0.99

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

Table 4. On-Device data preprocessing latency and model inference latency per frame (ms).

Method	Preprocessing Model Total		
	(ms) ↓	(ms) ↓	(ms) ↓
EfficientPhys-C	0	40	40
EfficientPhys-T1	0	300	300
EfficientPhys-T2	0	40	40
TS-CAN[10]	3	60	63
Dual-GAN[19]	275	N/A	> 275
POS[37]	0	27	27
CHROM[6]	0	28	28
ICA[26]	0	31	31

ms = Preprocessing and model latency on Raspberry Pi 4B per frame.

were not able to replicate every architecture to benchmark its on-device latency. The results show that EfficientPhys-C

Table 5. Ablation study on EfficientPhys-C (Top) an EfficientPhys-T1 (Down). Models are trained only on PURE and tested on UBFC.

Self-Attention	Diff	BatchNorm	MAE
✓	✓	✓	2.13
✗	✓	✓	2.43
✓	✓	✗	16.06
✓	✗	✗	16.06

TSM	Normal. Module	MAE
✓	✓	3.83
✓	✗	16.10
✗	✓	11.52

only takes 40ms to process a single frame and it does not take any extra computational time to perform preprocessing. On the other hand, due to the complex model architecture and additional time for calculating hand-crafted normalized raw and difference frames, TS-CAN takes 63ms per frame. As mentioned earlier, Dual-GAN has a complicated preprocessing procedure for facial landmark detection, segmentation, color transformation and augmentation. We implemented this and benchmarked the preprocessing module on our platform, and it took 275ms per frame, which is already 7x than the entire computational time of EfficientPhys-C. The estimation network in Dual-GAN also includes 12 2D convolution operations and numerous 1D convolution operations. Thus, we believe it would add a significant amount of computational time on top of the 275ms preprocessing time per frame. The default Transformer-based EfficientPhys (T1) has an unfavorable inference time due to its deep architecture design and takes 300ms to process every single frame. After reducing the depth to EfficientPhys-T2, it can achieve the same inference time as the EfficientPhys-C. However, EfficientPhys-T2 has the poorest performance on all three benchmark datasets.

Convolution vs. Transformer in Camera-Based Vitals Measurement. Although visual transformers have begun to achieve state-of-the-art performance in some vision tasks, it is not the case for the task of video-based vitals measurement. Based on the results shown in Table 1 and Table 2,

Efficient-C outperforms both Efficient-T1 by 45% of MAE in UBFC and similar performance in MMSE and PURE, while Efficient-C is more than 7x faster in terms of latency. When we shrink the Transformer-based EfficientPhys to a similar complexity as Convolution-based EfficientPhys, the performance is significantly diminished. The errors from the lightweight Transformer-based EfficientPhys-T2 increased 48% of MAE in UBFC, 141% of MAE in PURE and 15% of MAE in MMSE. These results indicate a shallow transformer architecture struggles to model subtle changes of skin pixels in the video. These finding suggest two potential insights. First, further optimizations will be necessary for transformers to outperform, even relatively shallow, convolutional models in this domain, this is possibly especially true when there is not a large amount of high-quality data available. As previous studies have shown [7], Transformers usually require more pre-training samples to obtain state-of-the-art accuracy. Unfortunately, currently the amount of data in the field of camera-based vital measurement is limited compared to other visual tasks. Our experiments in Table 3 also support this hypothesis where EfficientPhys-C surpasses both EfficientPhys-T1 and T2 with training only on PURE. We believe synthetic data is one way to help address this issue. Second, the good accuracy-efficiency trade-off for visual transformer might not be scaled to on-device architectures without further work. Since many on-device neural networks require significantly less amount of computing resources to perform real-time operations, scaling the Transformer architecture down is not ideal as our experimental results of EfficientPhys-T2 have shown.

Ablation Study. We provide ablation studies on various parameters in EfficientPhys-C and EfficientPhys-T1 in Table 5. Without the self-attention module, MAE of EfficientPhys-C is increased by 14%. Without the Normalization Module, in both EfficientPhys-C and EfficientPhys-T1, the MAEs increased by 753% and 420%. As Figure 3 illustrates, the output of the difference layer contains almost black pixels and these results indicate that neural methods are sensitive to the magnitude of pixel values and whether they are zero centered. Finally, without the tensor shift module (TSM) in transformer-based models, the error increased by 300% which indicates TSM plays an important role in exchanging temporal information and dynamics.

Simplifying Last-Mile ML Deployment. Numerous real-world applications are driven by novel machine learning algorithms. However, deploying these algorithms on different computing platforms has been extremely challenging for various reasons. One of these is that researchers sometimes only pay attention to the accuracy of the model and ignore the complexity of the last-mile engineering efforts. In this paper, we address this important issue through our one-stop architecture that takes the unprocessed raw frames and directly outputs the desired signal. This elegant and simple

design will not only reduce the burden of engineering required for cross-platform implementations, but also will help the research community to replicate and reproduce results.

Extensible to Other Signal. Finally, as another potential upside of our end-to-end design and the low latency, we envision EfficientPhys being applied to various other video-based applications. Since the input of our model is raw frames, we believe EfficientPhys can be easily extended to other tasks such as video-based blood pressure measurement and video understanding & recognition etc. On the other hand, most of the baseline methods we compared with (e.g., Dual-GAN, PulseGAN) require many custom preprocessing operations for video-based measurement which are less useful in other applications.

6. Broader Impacts and Ethics Statement

During the development of EfficientPhys, we wanted to ensure the innovations would not create larger disparities between different populations. By achieving state-of-the-art accuracy and efficiency as well as our simple and elegant design, we believe EfficientPhys will help make camera-based vitals measurement more widely available to the medical research community and broader community in computing. We also believe that this technology can have an especially strong impact in low-resource settings where there are greater barriers to accessing healthcare. We envision our proposed method could, with the appropriate clinical validation and regulatory approval, eventually be used in healthcare applications (e.g., real-time vitals measurement in telehealth appointments). During the COVID-19 pandemic the need for such technology has been clearly highlighted. We contextualize our contributions within the scope of democratizing technology for social good and helping to reduce health disparities with advanced AI technology. However, we are aware that machine learning systems are biased and can propagate inequalities. Before technology such as that presented in this paper is ready for deployment we need to make sure that that is not the case.

7. Conclusion

In this paper, we present a novel method called EfficientPhys to enable simple, fast, accurate camera-based contactless vitals measurement. We achieved strong performance with using significant less computational power. With the simple and elegant one-stop design, EfficientPhys also helps address the issue of last-time machine learning deployment and reduces health disparity.

References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. ViViT: A Video Vision

- Transformer. *arXiv:2103.15691 [cs]*, Mar. 2021. arXiv: 2103.15691.
- [2] Consumer Technology Association. Physical activity monitoring for heart rate, ansi/cta-2065, 2018.
 - [3] Serge Bobbia, Richard Macwan, Yannick Benezeth, Alamin Mansouri, and Julien Dubois. Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters*, 124:82–90, 2019.
 - [4] Bryson Carrier, Brenna Barrios, Brayden D Jolley, and James W Navalta. Validity and reliability of physiological data in applied settings measured by wearable technology: A rapid systematic review. *Technologies*, 8(4):70, 2020.
 - [5] 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.
 - [6] Gerard De Haan and Vincent Jeanne. Robust pulse rate from chrominance-based rppg. *IEEE Transactions on Biomedical Engineering*, 60(10):2878–2886, 2013.
 - [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
 - [8] Justin R. Estep, Ethan B. Blackford, and Christopher M. Meier. Recovering pulse rate during motion artifact with a multi-imager array for non-contact imaging photoplethysmography. In *2014 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 1462–1469, Oct. 2014. ISSN: 1062-922X.
 - [9] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7083–7093, 2019.
 - [10] Xin Liu, Josh Fromm, Shwetak Patel, and Daniel McDuff. Multi-task temporal shift attention networks for on-device contactless vitals measurement. *Advances in Neural Information Processing Systems*, 33:19400–19411, 2020.
 - [11] Xin Liu, Ziheng Jiang, Josh Fromm, Xuhai Xu, Shwetak Patel, and Daniel McDuff. MetaPhys: Few-Shot Adaptation for Non-Contact Physiological Measurement. *arXiv:2010.01773 [cs]*, Mar. 2021. arXiv: 2010.01773.
 - [12] Xin Liu, Shwetak Patel, and Daniel McDuff. Camera-based physiological sensing: Challenges and future directions. *arXiv preprint arXiv:2110.13362*, 2021.
 - [13] 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.
 - [14] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*, 2021.
 - [15] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. *arXiv:2103.14030 [cs]*, Mar. 2021. arXiv: 2103.14030.
 - [16] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021.
 - [17] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
 - [18] Hao Lu, Hu Han, and S Kevin Zhou. Dual-GAN: Joint BVP and Noise Modeling for Remote Physiological Measurement. page 10.
 - [19] Hao Lu, Hu Han, and S Kevin Zhou. Dual-gan: Joint bvp and noise modeling for remote physiological measurement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12404–12413, 2021.
 - [20] Daniel McDuff. Camera measurement of physiological vital signs. *ACM Computing Surveys (CSUR)*, 2021.
 - [21] Daniel McDuff, Javier Hernandez, Erroll Wood, Xin Liu, and Tadas Baltrusaitis. Advancing Non-Contact Vital Sign Measurement using Synthetic Avatars. *arXiv:2010.12949 [cs]*, Oct. 2020. arXiv: 2010.12949.
 - [22] 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.
 - [23] Benjamin W Nelson and Nicholas B Allen. Accuracy of consumer wearable heart rate measurement during an ecologically valid 24-hour period: intraindividual validation study. *JMIR mHealth and uHealth*, 7(3):e10828, 2019.
 - [24] Xuesong Niu, Zitong Yu, Hu Han, Xiaobai Li, Shiguang Shan, and Guoying Zhao. Video-based remote physiological measurement via cross-verified feature disentangling. In *European Conference on Computer Vision*, pages 295–310. Springer, 2020.
 - [25] 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*, 32:8026–8037, 2019.
 - [26] Ming-Zher Poh, Daniel J McDuff, and Rosalind W Picard. Advancements in noncontact, multiparameter physiological measurements using a webcam. *IEEE transactions on biomedical engineering*, 58(1):7–11, 2010.
 - [27] Hua Qi, Qing Guo, Felix Juefei-Xu, Xiaofei Xie, Lei Ma, Wei Feng, Yang Liu, and Jianjun Zhao. DeepRhythm: Exposing deepfakes with attentional visual heartbeat rhythms. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 4318–4327, 2020.
 - [28] Joel D Reece, Jennifer A Bunn, Minsoo Choi, and James W Navalta. Assessing heart rate using consumer technology association standards. *Technologies*, 9(3):46, 2021.
 - [29] Dangdang Shao, Yuting Yang, Chenbin Liu, Francis Tsow, Hui Yu, and Nongjian Tao. Noncontact monitoring breathing pattern, exhalation flow rate and pulse transit time. *IEEE Transactions on Biomedical Engineering*, 61(11):2760–2767, 2014.
 - [30] Anna Shcherbina, C Mikael Mattsson, Daryl Waggott, Heidi Salisbury, Jeffrey W Christle, Trevor Hastie, Matthew T Wheeler, and Euan A Ashley. Accuracy in wrist-worn, sensor-based measurements of heart rate and energy expenditure in

- a diverse cohort. *Journal of personalized medicine*, 7(2):3, 2017.
- [31] Rencheng Song, Huan Chen, Juan Cheng, Chang Li, Yu Liu, and Xun Chen. PulseGAN: Learning to generate realistic pulse waveforms in remote photoplethysmography. *IEEE Journal of Biomedical and Health Informatics*, 25(5):1373–1384, 2021.
- [32] Ronny Stricker, Steffen Müller, and Horst-Michael Gross. Non-contact video-based pulse rate measurement on a mobile service robot. In *The 23rd IEEE International Symposium on Robot and Human Interactive Communication*, pages 1056–1062. IEEE, 2014.
- [33] Chihiro Takano and Yuji Ohta. Heart rate measurement based on a time-lapse image. *Medical engineering & physics*, 29(8):853–857, 2007.
- [34] Kenji Takazawa, Nobuhiro Tanaka, Masami Fujita, Osamu Matsuoka, Tokuyu Saiki, Masaru Aikawa, Sinobu Tamura, and Chiharu Ibukiyama. Assessment of Vasoactive Agents and Vascular Aging by the Second Derivative of Photoplethysmogram Waveform. *Hypertension*, 32(2):365–370, Aug. 1998.
- [35] Yun-Yun Tsou, Yi-An Lee, Chiou-Ting Hsu, and Shang-Hung Chang. Siamese-rppg network: Remote photoplethysmography signal estimation from face videos. In *Proceedings of the 35th annual ACM symposium on applied computing*, pages 2066–2073, 2020.
- [36] Wim Verkruyse, Lars O Svaasand, and J Stuart Nelson. Remote plethysmographic imaging using ambient light. *Optics express*, 16(26):21434–21445, 2008.
- [37] Wenjin Wang, Albertus C. den Brinker, Sander Stuijk, and Gerard de Haan. Algorithmic Principles of Remote PPG. *IEEE Transactions on Biomedical Engineering*, 64(7):1479–1491, July 2017. Conference Name: IEEE Transactions on Biomedical Engineering.
- [38] Ting Wu, Vladimir Blazek, and Hans Juergen Schmitt. Photoplethysmography imaging: a new noninvasive and noncontact method for mapping of the dermal perfusion changes. In *Optical Techniques and Instrumentation for the Measurement of Blood Composition, Structure, and Dynamics*, volume 4163, pages 62–70. International Society for Optics and Photonics, 2000.
- [39] Zitong Yu, Xiaobai Li, Pichao Wang, and Guoying Zhao. Transrppg: Remote photoplethysmography transformer for 3d mask face presentation attack detection. *IEEE Signal Processing Letters*, 2021.
- [40] Zitong Yu, Xiaobai Li, and Guoying Zhao. Remote photoplethysmograph signal measurement from facial videos using spatio-temporal networks. *arXiv preprint arXiv:1905.02419*, 2019.
- [41] Zitong Yu, Yuming Shen, Jingang Shi, Hengshuang Zhao, Philip HS Torr, and Guoying Zhao. Physformer: facial video-based physiological measurement with temporal difference transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4186–4196, 2022.
- [42] Zheng Zhang, Jeff M Girard, Yue Wu, Xing Zhang, Peng Liu, Umur Ciftci, Shaun Canavan, Michael Reale, Andy Horowitz, Huiyuan 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.