

Meta-rPPG: Remote Heart Rate Estimation Using a Transductive Meta-Learner

Eugene Lee Evan Chen Chen-Yi Lee

Institute of Electronics, National Chiao Tung University
{eugenelet.ee06g, evanchen.ee06}@nctu.edu.tw, cylee@si2lab.org

Abstract. Remote heart rate estimation is the measurement of heart rate without any physical contact with the subject and is accomplished using remote photoplethysmography (rPPG) in this work. rPPG signals are usually collected using a video camera with a limitation of being sensitive to multiple contributing factors, e.g. variation in skin tone, lighting condition and facial structure. End-to-end supervised learning approach performs well when training data is abundant, covering a distribution that doesn't deviate too much from the distribution of testing data or during deployment. To cope with the unforeseeable distributional changes during deployment, we propose a transductive meta-learner that takes unlabeled samples during testing (deployment) for a self-supervised weight adjustment (also known as transductive inference), providing fast adaptation to the distributional changes. Using this approach, we achieve state-of-the-art performance on MAHNOB-HCI and UBFC-rPPG.

Keywords: Remote heart rate estimation, rPPG, meta-learning, transductive inference

1 Introduction

Remote photoplethysmography (rPPG) is useful in situations where conventional approaches for heart rate estimation like electrocardiogram (ECG) and photoplethysmogram (PPG) that requires physical contact with the subject is infeasible. The acquisition of rPPG signal is useful for the estimation of physiological signal like heart rate (HR) and heart rate variation (HRV) which are important parameters for remote health-care. rPPG is usually obtained using a video camera and a growing number of studies have used rPPG for HR estimation [40,39,1,10,24,52,55].

The adoption of deep learning techniques for the measurement of rPPG from face images is not novel and is supported by numerous studies [8,61,35,60]. All of the prior work involving deep learning uses an end-to-end supervised learning approach where a global model is deployed during inference (to the best of our knowledge), also known as inductive inference. It has been pointed out by the recent advances in deep learning that such approach doesn't perform well when there are changes in the modeled distribution when moving from training dataset to the real world [2,15]. The unpredictability of the changes in environment and

test subjects (skin-tone, facial structure, etc.) hinders the performance of remote heart rate estimation. To cope with the dynamical changes of the environment and subjects, we propose a transductive meta-learner that is able to perform fast adaption during deployment. Our algorithm introduces a warm-start time frame (2 seconds) for adaptation (weight update) to cope with the distributional changes, resulting in better performance during remote heart rate estimation.

Current advances in meta-learning [14,33,44] have shed light on the techniques of designing and training a deep neural network (DNN) that is able to adapt to new tasks through minimal weight changes. As prior work in meta-learning is built on well-defined tasks consisting of the classification of a few labeled examples (shots) provided as support set during test time for fast adaptation, it can't be directly applied to our context as labeled data are unobtainable during deployment. In our context, adaptation has to be done in a self-supervised fashion, hence we incorporate transductive inference techniques [18,26] into the design of our learning algorithm. The application of meta-learning and transductive inference to rPPG estimation is not trivial since we have to consider the modeling of both spatial and temporal information in the formulation of our meta-learner. In our work, we split our DNN into two parts: a feature extractor and a rPPG estimator modeled by a convolutional neural network (CNN) and long short-term memory (LSTM) [17] network respectively. We introduce a synthetic gradient generator modeled by a shallow Hourglass network [32] along with a novel prototypical distance minimizer to perform transductive inference when labeled data are not available during deployment. Intuitively, the variation in distribution is caused by the visual aspect of the incoming signal, hence adaptation (weight changes) is only performed on the feature extractor while the rPPG estimator's parameters will be frozen during deployment. In summary, our main contributions are as follows:

1. In Section 3, we propose a meta-learning framework that exploits spatiotemporal information for remote heart rate estimation, specifically designed for fast adaptation to new video sequences during deployment.
2. In Section 3.2, we introduce a transductive inference method that makes use of unlabeled data for fast adaptation during deployment using a *synthetic gradient generator* and a novel *prototypical distance minimizer*.
3. In Section 3.3, we propose the formulation of rPPG estimation as an ordinal regression problem to cope with the mismatch in temporal domain of both visual input and collected PPG data, as studied in [11].
4. In Section 4, we validate our proposed methods empirically on two public face heart rate dataset (MAHNOB-HCI [48] and UBFC-rPPG [3]), trained using our collected dataset. Our experimental results demonstrate that our approach outperforms existing methods on remote heart rate estimation.

2 Related Work

Remote Heart Rate Estimation. The study of the measurement of heart rate using a video camera is not novel and has been supported by existing works.

The remote measurement of heart rate is accomplished through the calculation of the cardiac pulse from the extracted rPPG signal which embeds the blood volume pulse (BVP). rPPG signal is extracted from the visible spectrum of light acquired by a video camera [7]. Contact PPG sensors found in wearable devices and medical equipment estimates the heart rate by acquiring the light reflected from the skin using a light-emitting diode (LED) as its source. The constituents of the reflected light is composed of the amount of light absorbed by the skin (constant) and a time varying pulse signal contributed by the absorption of light by blood capillaries beneath the skin. Since existing contact PPG methods are non-intrusive, light sources with specific wavelengths are used, e.g. green light (525 nm) and infrared light (880 nm), having the properties of good absorption by blood and has minimal overlap with the visible light spectrum [27].

The feasibility of remote heart rate estimation was first proven by Verkruyse et al. [53], showing that PPG signals can be extracted from videos collected under an ambient light setting using webcams. To deal with noise, Poh et al. [40] performed blind source separation on webcam filmed videos to extract pulse signal. A study is done in [51] revolving the comparison of the absorptivity of independent channels from the RGB channels by the human skin. It is concluded that the green channel is easily absorbed by the human skin, giving high signal-to-noise ratio for PPG signal acquisition. Based on the results in [51], Li et al. [24] used only the green channel for rPPG measurement and introduced an adaptive filtering technique (normalized least mean square) to eliminate motion artifact during deployment.

As it is hypothesized that all three channels (RGB) contain considerable information for the estimation of the rPPG signals, different channels are weighted and summed to receive better results when compared to using only a single channel (green) in [30]. CHROM [10] estimates the rPPG by using a method that linearly combines the RGB channels using the knowledge of a skin reflection model to separate pulse signal from motion-induced distortion. The POS [55] method and the SB [56] model used the same skin reflection model as CHROM, but made a different assumption on the distortion model where another projection direction is used to separate the pulse and the distortions. All the mentioned methods are based on the calculation of the spatial mean of the entire face, assuming the contribution of each pixel to the estimation of rPPG is equal. Such assumption is not resilient to noise, rendering it infeasible in extreme conditions.

Recently, deep learning approaches have been introduced for rPPG signal estimation. DeepPhys [8] is an end-to-end supervised method implemented using a feed-forward CNN. Instead of using a Long Short-Term Memory (LSTM) cell to model the temporal information, an attention mechanism is used to learn the difference between frames. rPPGNet [61] considers the possibility of different types of video compression affecting the performance of rPPG estimation and proposed a network that handles both video quality reconstruction and rPPG signal estimation trained in an end-to-end fashion. RhythmNet [35] trained a CNN-RNN model based on a training set that contains diverse illumination and pose to enhance performance on public dataset. PhysNet [60] constructed

both a 3DCNN-based network and a RNN-based network and compared their performance.

Meta-Learning. The motivation for introducing meta-learning to our rPPG estimation framework is to perform fast adaptation of weights when our network is deployed in a setting that is not covered by our training distribution. Most of the studies revolving meta-learning are in few-shot classification, which have well-defined training and testing tasks. The structure of such well-defined tasks is exploited in [29,45,54] where the structure of the network is designed to take both training and test samples into consideration during inference. Gradient-based few-shot learning has also been proposed in [14,33,41] with a limitation that the network capacity has to be small to prevent overfitting due to the small number of training samples. Few-shot learning based on the metric space has also been studied in [21,54,47]. Meta-learning has also been applied to tasks beyond few-shot classification [12,59,58], showing convincing results.

In our work, the parameters of our network is divided into two parts where one is responsible for fast adaptation and the other only responsible for estimation. Similar update methodology has been introduced in [31,16,18,44,62]. All the proposed methods have different weight update orders and network construction but they have the consensus on the effectiveness of maintaining two sets of weights in a few-shot learning setting.

3 Methodology

Given an input image from the camera, a face is first detected followed by the detection of facial landmarks. Pixel values within the landmarks are retained while the rest are filled with 0. The face image is then cropped and reshaped into a $K \times K$ image. We define the i -th video stream and PPG data containing a single subject as $\mathbf{x}^{(i)}$ and $\mathbf{t}^{(i)}$ respectively (we will drop the superscript i whenever the context involves a single video stream for brevity) where $\mathbf{x}^{(i)}$ and $\mathbf{t}^{(i)}$ are sampled from a distribution of tasks $(\mathbf{x}^{(i)}, \mathbf{t}^{(i)}) \sim p(\mathcal{T})$. Continuous T frames of face images are aggregated into a sequence, giving us the input data of our network, $\mathbf{x}_{t:t+T} \in \mathbb{R}^{K \times K \times T}$. The face sequences are paired with the PPG signal $\mathbf{t}_{t:t+T} \in \mathbb{R}^T$ collected using a PPG sensor placed beneath the index finger, where each sample \mathbf{t}_t is temporally aligned to each frame of the face sequence \mathbf{x}_t during the collection process. The output of our network is a rPPG signal $\mathbf{y} \in \mathbb{R}^T$ which is an estimation of the PPG signal having a small temporal offset caused by the carotid-radial pulse wave velocity (PWV) [11].

The training of our network is different from typical training practice found in an end-to-end supervised learning setting. We cast the learning of our network into a few-shot learning setting, involving both support set \mathcal{S} and query set \mathcal{Q} . In a few-shot learning setting, \mathcal{S} and \mathcal{Q} are sampled from a large pool of tasks distribution $\{\mathcal{S}_n, \mathcal{Q}_n\} \sim \mathcal{T}$ where the sampled \mathcal{S}_n and \mathcal{Q}_n share the same set of classes but have disjoint samples for a classification problem. Our learning setting differs from the existing few-shot classification setting in two ways: 1. we

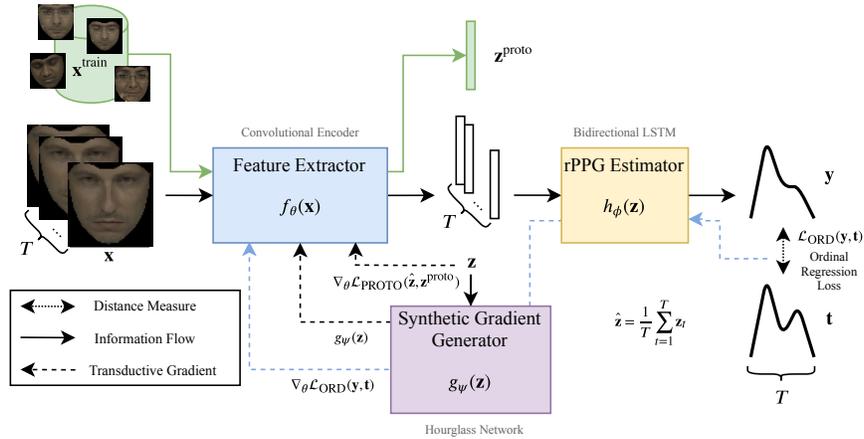


Fig. 1: Overview of our system for transductive inference. Three modules: feature extractor, rPPG estimator and synthetic gradient generator are found in our system. During inference, only the feature extractor and rPPG estimator will be used. Only the parameters of the feature extractor, θ , will be updated during transductive learning. Gradients used for the update of θ are shown in dashed lines where the blue dashed line is only present during the training phase of meta-learning. $\mathbf{z}^{\text{proto}}$ is generated using the training set, $\mathbf{x}^{\text{train}}$.

are solving a regression problem, 2. our input is a sequence of images instead of still images that are distributed independent and identically.

In comparison to the typical few-shot learning setting, instead of sampling a set of samples originating from disjoint classes, we sample a sequence of image, $\mathbf{x}_{t:t+T}$, from disjoint video streams, e.g. videos containing different subject or background. During each sampling process, we sample N independent video streams from our collected pool of video streams also known as distribution of tasks, $p(\mathcal{T})$. For each video stream, we split it into shorter sequences, where each sequence is further split into V and W frames. In correspondence to the few-shot learning setting, V is the support set \mathcal{S} used for adaptation and W is the query set \mathcal{Q} used to evaluate the performance of the model.

3.1 Network Architecture

Our meta-learning framework for remote heart rate estimation consists of three modules: convolutional encoder, rPPG estimator and synthetic gradient generator. To infer the rPPG signal, only the convolutional encoder and rPPG estimator will be used whereas the synthetic gradient generator will only be used during transductive learning. Our network is designed to exploit spatiotemporal information by first modeling visual information using a convolutional encoder followed by the modeling of the estimation of PPG signal using a LSTM rPPG estimator. We name our approach Meta-rPPG where detailed configuration of

our architecture is shown in Table 1. An overview of our proposed framework is shown in Fig. 1.

Convolutional Encoder. To extract latent features from a stream of images, we use a feature extractor modeled by a CNN, $f_\theta(\cdot)$. We use a ResNet-like structure as the backbone of our convolutional encoder. Given an input stream of T frames, the convolutional encoder is shared across frames:

$$p_\theta(\mathbf{z}_i|\mathbf{x}_i) = f_\theta(\mathbf{x}_i), \quad (1)$$

giving us T independent distribution of latent features, $\{p_\theta(\mathbf{z}_i|\mathbf{x}_i)\}_{i=t}^{t+T}$, that will be fed to a rPPG estimator to model the temporal information for the estimation of rPPG signal.

rPPG Estimator. The latent features extracted from the input image stream are then passed to rPPG estimator modeled by a LSTM-MLP module, $h_\phi(\cdot)$. The intuition behind the split is to separate the parameters responsible for visual modeling from the parameters responsible for the estimation of rPPG signal, accomplished via the temporal modeling of the latent features:

$$p_\phi(\mathbf{y}|\mathbf{z}_t, \mathbf{z}_{t+1}, \dots, \mathbf{z}_{t+T}) = h_\phi(\{p_\theta(\mathbf{z}_i|\mathbf{x}_i)\}_{i=t}^{t+T}). \quad (2)$$

The LSTM module is designed to model the temporal information of a fixed sequence of T latent features. The output of each step of the LSTM is followed by a MLP module which is responsible for the estimation of rPPG signal. As we model the estimation of rPPG signal as an ordinal regression task, we have a multitask (S tasks) output. Details are deferred to Section 3.3.

Synthetic Gradient Generator. For the fast adaptation of the parameters of our model during deployment, we introduce a synthetic gradient generator modeled by a shallow Hourglass network [32], $g_\psi(\cdot)$. The idea of using a synthetic gradient generator for transductive inference is not novel as it was first introduced in [19] to parallelize the backpropagation of gradient and to augment backpropagation-through-time (BPTT) of long sequences found in LSTM. It is then applied to a few-shot learning framework in [18] to generate gradient for unlabeled samples, giving a significant boost in performance. Our synthetic gradient generator attempts to model the gradient at \mathbf{z} backpropagated from the ordinal regression loss from the output, $\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})$ (defined in (13)):

$$g_\psi(\mathbf{z}) = \nabla_{\mathbf{z}} \mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t}). \quad (3)$$

Our synthetic gradient generator will be used during transductive inference for the fast adaptation to new video sequences and at the adaptation phase during training, using the support set \mathcal{S} .

Table 1: Conv2DBlocks are composed of Conv2D, Batchnorm, average pooling, and ReLU. Conv1DBlocks are composed of Conv1D, Batchnorm and ReLU. Shortcut connections are added between Conv2DBlocks for the Convolutional Encoder. The synthetic gradient generator is designed as a Hourglass network. \checkmark indicates the information the layer acts upon. Output size is defined as $T \times \text{Channels} \times K \times K$ for Convolutional Encoder and $T \times \text{Channels}$ for the rest.

Module	Layer	Output Size	Kernel Size	Spatial	Temporal
Convolutional Encoder	Conv2DBlock	$60 \times 32 \times 32 \times 32$	3×3	\checkmark	
	Conv2DBlock	$60 \times 48 \times 16 \times 16$	3×3	\checkmark	
	Conv2DBlock	$60 \times 64 \times 8 \times 8$	3×3	\checkmark	
	Conv2DBlock	$60 \times 80 \times 4 \times 4$	3×3	\checkmark	
	Conv2DBlock	$60 \times 120 \times 2 \times 2$	3×3	\checkmark	
	AvgPool	60×120	2×2	\checkmark	
rPPG Estimator	Bidirectional LSTM	60×120	-	\checkmark	\checkmark
	Linear	60×80	-	\checkmark	
	Ordinal	60×40	-	\checkmark	
Synthetic Gradient Generator	Conv1DBlock	40×120	3×3	\checkmark	\checkmark
	Conv1DBlock	20×120	3×3	\checkmark	\checkmark
	Conv1DBlock	40×120	3×3	\checkmark	\checkmark
	Conv1DBlock	60×120	3×3	\checkmark	\checkmark

3.2 Transductive Meta-Learning

During the deployment of Meta-rPPG, we have to consider the possibility of observing input samples that are not modeled by our model, also known as out-of-distribution samples. A possible solution is through the introduction of a pre-processing step that attempts to project the input data into a common distribution that is covered by our model. The consideration of an infinitely large distribution containing all possible scenarios during deployment is near impossible for any pre-processing technique. To cope with the shift in distribution, we propose two methods for transductive inference which provides gradient to our convolutional encoder $f_{\theta}(\cdot)$ during deployment. The first method is through the generation of synthetic gradients using a generator that is modeled during training. The second method attempts to minimize prototypical distance between different tasks which is based on the hypothesis that the rPPG estimator $h_{\phi}(\cdot)$ is only responsible for the estimation of rPPG signal. The modeling of the estimation of rPPG signal should not be affected by the visual input and is limited to a constrained distribution.

Generating Synthetic Gradient. To generate gradients for weight update of our model during inference, we introduce a synthetic gradient generator that models the backpropagated gradient from the final output that uses an ordinal loss $\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})$ from (13), up to the output of the feature extractor \mathbf{z} , which

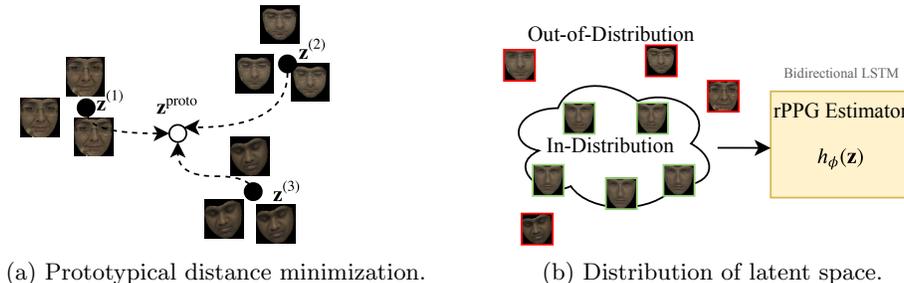


Fig. 2: (a) shows the high-level concept on the application of prototypical distance minimization on the latent space. (b) shows that the rPPG estimator only performs well on the in-distribution samples covered by the training dataset while performance on out-of-distribution samples is sub-optimal. In-distribution samples are outlined in green while out-of-distribution samples are outlined in red. The prototypical distance minimizer generates gradient that forces the out-of-distribution samples towards the center of the in-distribution samples.

can be simply put as $\nabla_{\mathbf{z}}\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})$. As our synthetic gradient generator $g_\psi(\mathbf{z})$ attempts to model $\nabla_{\mathbf{z}}\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})$, we can observe its role in the chain-rule of gradient update of the parameters of the feature extractor, θ :

$$\theta = \theta - \alpha \frac{\partial \mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \theta} \quad (4)$$

$$= \theta - \alpha g_\psi(\mathbf{z}) \frac{\partial \mathbf{z}}{\partial \theta}. \quad (5)$$

During the learning phase of training, we can update the weights of our synthetic gradient generator by minimizing the following objective function:

$$\mathcal{L}_{\text{SYN}}(g_\psi(\mathbf{z}), \nabla_{\mathbf{z}}\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})) = \|g_\psi(\mathbf{z}) - \nabla_{\mathbf{z}}\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})\|_2^2, \quad (6)$$

where the weight update of ψ is given as:

$$\psi = \psi - \eta \nabla_{\psi} \mathcal{L}_{\text{SYN}}(g_\psi(\mathbf{z}), \nabla_{\mathbf{z}}\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})). \quad (7)$$

Minimizing Prototypical Distance. As rPPG estimation is based on a visual input, there’s no guarantee that the samples collected for training is consistent with the samples used during testing due to the broad distribution in the visual space. In a statistical viewpoint, data provided by the training set modeled by our network can be viewed as in-distribution samples whereas data collected under a different setting, e.g. lighting condition, subject and camera settings, are considered as out-of-distribution samples. It is not surprising that deep neural networks doesn’t perform well on out-of-distribution samples as studied in [25,42]. To overcome this limitation, we propose a prototypical distance minimization technique that can be applied on video sequences having a property

where the statistical information that needs to be modeled by a neural network doesn't vary too much over time. We show a conceptual illustration in Fig. 2.

First, we consider each video sequences as separate task that needs to be modeled. We define the prototype of the latent variable of a specific task as:

$$\mathbf{z}^{(i)} = \frac{1}{T} \sum_{t=1}^T p_{\theta}(\mathbf{z}_t^{(i)} | \mathbf{x}_t^{(i)}). \quad (8)$$

Here, T consecutive samples from task or video i is sampled and the average is taken across the latent variable generated by the convolutional decoder from (1). We then obtain our first global latent variable prototype as:

$$\mathbf{z}^{\text{proto}} = \mathbb{E}_{\mathbf{x}^{(i)} \sim p(\mathcal{T})} \frac{1}{T} \sum_{t=1}^T p_{\theta}(\mathbf{z}_t^{(i)} | \mathbf{x}_t^{(i)}) \quad (9)$$

$$= \mathbb{E}_{\mathbf{x}^{(i)} \sim p(\mathcal{T})} \frac{1}{T} \sum_{t=1}^T f_{\theta}(\mathbf{x}_t^{(i)}). \quad (10)$$

As mentioned earlier, we perform Monte Carlo sampling of N tasks or videos for every training iteration, hence we are unable to obtain the statistical mean of the entire dataset in one shot. A more computationally feasible approach is to update our global latent variable prototype in an iterative fashion as:

$$\mathbf{z}^{\text{proto}} = \gamma \mathbf{z}^{\text{proto}} + (1 - \gamma) \mathbb{E}_{\mathbf{x}^{(i)} \sim p(\mathcal{T})} \frac{1}{T} \sum_{t=1}^T f_{\theta}(\mathbf{x}_t^{(i)}), \quad (11)$$

which can be understood as the weighted average between the old term and the newly sampled global latent variable prototype via the introduction of the hyperparameter γ . The global prototype in (11) is obtained during the learning phase of training and transductive gradient is generated by minimizing the loss:

$$\min_{\theta} \mathcal{L}_{\text{PROTO}}(\mathbf{z}, \mathbf{z}^{\text{proto}}) = \min_{\theta} \mathbb{E}_{\mathbf{x}^{(i)} \sim p(\mathcal{T})} \frac{1}{T} \sum_{t=1}^T \|p_{\theta}(\mathbf{z}_t^{(i)} | \mathbf{x}_t^{(i)}) - \mathbf{z}^{\text{proto}}\|_2^2. \quad (12)$$

Training of Meta-Learner. A meta-learner will perform well during testing if the setting during testing is similar to the setting during training. As mentioned earlier, we split the N tasks samples from all our collected video into sequences that are further split into V and W frames. We use V for the update of θ phrased as the *adaptation phase* and W for the update of ϕ , ψ , θ and $\mathbf{z}^{\text{proto}}$ phrased as the *learning phase*. In general, the split of V and W is put as $W > V$ and $W \cap V = \emptyset$ with the intuition that we attempt to minimize the frames required for adaptation and perform learning on the adapted space or distribution.

The role of the adaptation phase is to train our convolutional encoder $f_{\theta}(\cdot)$ to map input image sequences to a representation or latent space that will perform well when fed to our rPPG estimator $h_{\phi}(\cdot)$. During training, gradients from three

sources will be used in the adaptation phase: synthetic gradient generator, prototypical distance minimizer and from the ordinal regression loss. During testing or deployment, we don't have any labeled data available, hence gradient from the ordinal regression loss is unattainable. Instead, we use gradients from our synthetic gradient generator and prototypical distance minimizer for adaptation. Note that the adaptation phase is run for L steps on the same V frames.

The role of the learning phase is to force our model to learn a suitable latent representation for rPPG signal estimation based on its image domain correspondence. Supervised learning is required in the learning phase, hence this phase will only be present during training. Here, θ and ϕ which corresponds to the convolutional encoder and rPPG estimator will be trained in an end-to-end fashion whereas the parameters of the synthetic gradient generator ψ will be updated using the gradient backpropagated from the ordinal regression loss. ψ is updated by minimizing the synthetic loss given in (6).

Before the inclusion of the adaptation phase during training, we first pre-train our network under the learning phase for R epochs. The reason is that the gradient backpropagated to the synthetic gradient generator and the prototype used for the prototypical minimizer rely on the task at hand (rPPG estimation) rather than using gradient and prototype based on a set of random weights, which could lead to instability. We summarize the training of our meta-learner in Algorithm 1.

3.3 Posing rPPG Estimation as an Ordinal Regression Task

Ordinal regression is commonly used in task that requires the prediction of labels that contains ordering information, e.g. age estimation [36,6], progression of various diseases [13,57,50,46], text message advertising [43] and various recommender systems [37]. The motivation of using ordinal regression in our work is because there's an ordering of rPPG value that can be exploited and there will always be a temporal and magnitude discrepancy between rPPG and PPG signal as they originate from different parts of the human body [11].

To cast the estimation of rPPG signal as an ordinal regression problem, we first normalize a segment of PPG signal to be within 0 and 1. We then quantize it uniformly into 40 segments where each segment represents a *rank*. Using the PPG segment $\mathbf{t}_{t:t+T}$ as an example, the quantized t -th sample will be categorized into one of the S segments $\{\tau_1 < \dots < \tau_S\}$. With a slight abuse of notations, if the categorized value for the t -th sample falls in the s -th segment, we denote it as $\mathbf{t}_{t,s} = \mathbb{1}\{\mathbf{t}_t > \tau_s\}$ to keep our formulation concise. $\mathbb{1}\{\cdot\}$ is an indicator function that returns 1 if the inner condition is true and 0 otherwise. As in [35], our rPPG estimator will have to solve S binary classification problem given as:

$$\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t}) = -\frac{1}{T} \sum_{t=1}^T \sum_{s=1}^S \mathbf{t}_{t,s} \log(p_{\phi}(\mathbf{y}_{t,s} | \mathbf{z}_{t:t+T})) + (1 - \mathbf{t}_{t,s}) \log(1 - p_{\phi}(\mathbf{y}_{t,s} | \mathbf{z}_{t:t+T})). \quad (13)$$

Note that in (13), the same notational abuse is applied to the output of our rPPG estimator, \mathbf{y} . During inference, we can obtain T rPPG samples corresponding to

Algorithm 1 Training of Meta-Learner

```

1: Input:  $p(\mathcal{T})$ : distribution of tasks
2: for  $i \leftarrow 1, R$  do ▷ Pre-train network in an end-to-end fashion for  $R$  epochs
3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$ 
4:   for  $(\mathbf{x}, \mathbf{t}) \sim \mathcal{T}_i$  do
5:     Update  $\theta$  and  $\phi$  by minimizing  $\mathcal{L}_{\text{ORD}}(\mathbf{y}, \mathbf{t})$  from (13)
6:   end for
7: end for
8: while not done do ▷ Begin transductive meta-learning
9:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$ 
10:  for  $(\mathbf{x}, \mathbf{t}) \sim \mathcal{T}_i$  do
11:     $\{\hat{\mathbf{x}}, \hat{\mathbf{t}}\}, \{\tilde{\mathbf{x}}, \tilde{\mathbf{t}}\} \leftarrow \mathbf{x}, \mathbf{t}$  ▷ Split into  $V$  and  $W$  consecutive frames
12:    for  $i \leftarrow 1, L$  do ▷ Adaptation phase (run  $L$  steps)
13:       $\theta \leftarrow \theta - \alpha(\nabla_{\theta} \mathcal{L}_{\text{proto}}(\hat{\mathbf{z}}, \hat{\mathbf{z}}^{\text{proto}}) + \nabla_{\theta} \mathcal{L}_{\text{ORD}}(\hat{\mathbf{y}}, \hat{\mathbf{t}}) + f_{\psi}(\hat{\mathbf{z}}))$ 
14:    end for
15:     $\psi = \psi - \eta \nabla_{\psi} \mathcal{L}_{\text{SYN}}(f_{\psi}(\hat{\mathbf{z}}), \nabla_{\hat{\mathbf{z}}} \mathcal{L}_{\text{ORD}}(\tilde{\mathbf{y}}, \tilde{\mathbf{t}}))$  ▷ Learning phase
16:     $\theta = \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{ORD}}(\tilde{\mathbf{y}}, \tilde{\mathbf{t}})$ 
17:     $\phi = \phi - \eta \nabla_{\phi} \mathcal{L}_{\text{ORD}}(\tilde{\mathbf{y}}, \tilde{\mathbf{t}})$ 
18:     $\mathbf{z}^{\text{proto}} = \gamma \mathbf{z}^{\text{proto}} + (1 - \gamma) \mathbb{E}_{\tilde{\mathbf{x}}^{(i)} \sim \tilde{\mathbf{x}} \frac{1}{T} \sum_{t=1}^T f_{\theta}(\tilde{\mathbf{x}}^{(i)})}$ 
19:  end for
20: end while

```

T consecutive frames, $\mathbf{x}_{t:t+T}$, fed to our model, giving us:

$$\mathbf{y}_{t:t+T} = h_{\phi}(\mathbf{z}_{t:t+T}) = \left\{ \sum_{s=1}^S \mathbb{1}\{p_{\phi}(\mathbf{y}_{i,s} = 1 | \mathbf{z}_{t:t+T}) > 0.5\} \right\}_{i=t}^T. \quad (14)$$

4 Experiments

We show the efficacy of our proposed method by performing empirical simulation MAHNOB-HCI [48] and UBFC-rPPG [3] using a model trained using our collected dataset. We also show how transductive inference helps in adapting to unseen datasets using a model trained using our own collected training dataset. Source code is available at <https://github.com/eugenelet/Meta-rPPG>.

4.1 Dataset and Experimental Settings

MAHNOB-HCI dataset [48] is recorded at 61 fps using a resolution of 780×580 and includes 527 videos from 27 subjects. Since our training dataset is collected at 30 fps, we downsample the videos to 30 fps by getting rid half of the video frames. To generate ground truth heart rate (HR) for evaluation, we use the EXG1 signals containing ECG signal. Peaks of ECG signal are found using `scipy.signal.find_peaks` and distance in samples between peaks is used for the calculation of HR. To compare fairly with previous works [8,34,49,60] [3, 12, 18], we follow the same routine in their works by using 30 seconds clip (frames 306 to 2135) of each video.

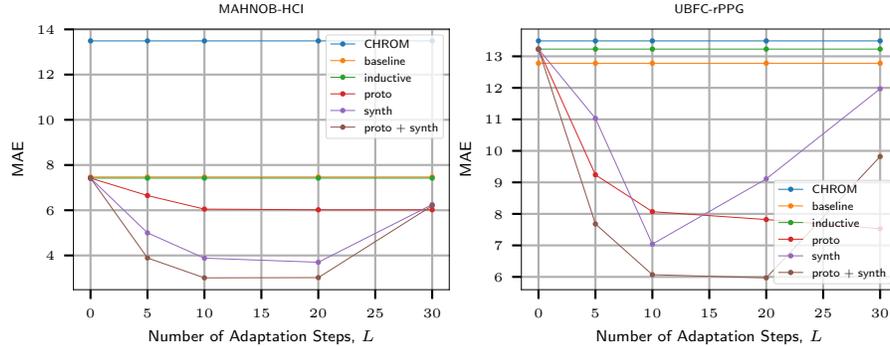


Fig. 3: MAE obtained using different rPPG estimation methods. Demonstrates how the number of adaptation steps, L , affects performance.

UBFC-rPPG dataset [3] is relatively new and contains 38 uncompressed videos along with finger oximeter signals. Ground truth heart rate is provided, hence it can directly be used for evaluation without any additional processing. This dataset has a wider range of HR collected induced by making participants play a time-sensitive mathematical game that supposedly raises their HR. Diffuse reflections was created by introducing ambient light in the experiment.

Collected Dataset. For the training of our model, we collected our own dataset. A RGB video camera is used for the collection of video at 30 fps using a resolution of 480×640 . PPG signals are collected using our self-designed device [22] where raw data can be sent to another host device via a UART port. Collection of video and PPG signal are synced by connecting both devices to a Nvidia Jetson TX2. Since our PPG sensor device collects PPG signal at 100 Hz, we downsample it to 30 Hz to match the visual stream. Video and PPG samples of length over 2 minutes are collected and are cropped to 2 minutes to match our proposed meta-learning framework. A total of 19 videos are collected where 18 are used for training and 1 for validation.

Experimental Settings. Facial videos and corresponding PPG signals are synchronized before training. For each video clip, we use the face detector found in `dlib` that uses [9] cascaded with a facial landmark detector implemented using [20]. To keep the results of our face detector consistent, we use a median flow tracker provided by OpenCV [5]. Pixels within the landmarks are retained and 5 pixels are used for zero-padding beyond the outermost pixels. The resulting face images are then resized to 64×64 . All experiments and network training are done on a single Nvidia GTX1080Ti using PyTorch [38]. The SGD optimizer is used. We set the learning rate, $\eta = 10^{-3}$, the adaptive learning rate, $\alpha = 10^{-5}$ and the prototype update weight, $\gamma = 0.8$. We train all the models for 20 epochs.

Performance Metrics. For the evaluation of public datasets, we report the error’s standard deviation (SD), mean absolute error (MAE), root mean square error (RMSE) and Pearson’s correlation coefficient (R).

4.2 Evaluation on MAHNOB-HCI and UBFC-rPPG

We evaluate 5 different configurations as part of the ablation study of the methods we introduced. The first configuration is the training of our proposed architecture in an end-to-end fashion and perform inductive inference on the test set, End-to-end (baseline). The second configuration is to train our network using all of our proposed methods but doesn’t perform adaptation during inference, Meta-rPPG (inductive). The third and fourth configuration perform transductive inference by using either prototypical distance minimizer, Meta-rPPG (proto only), or synthetic gradient generator, Meta-rPPG (synth only), respectively. The final configuration uses both proposed method for transductive inference, Meta-rPPG (proto+synth).

Results of average HR measurement for MAHNOB-HCI and UBFC-rPPG are shown in Table 2 and 3 respectively. We also show MAE results by the varying number of adaptation steps, L in Fig. 3. Since we use $L = 10$ during training, setting $L = 10$ during evaluation also gives us the best results, which agrees with the common practice used in meta-learning [14,33]. Note that our proposed architecture used for end-to-end (baseline) training is similar to [60] with a difference that we are using only 18 videos of length 2 minutes for training whereas [60] uses OBF dataset [23] that contains 212 videos of length 5 minutes. Considering the difference in magnitude of dataset size, it’s understandable that training end-to-end using our network doesn’t perform as well as in [60]. This also indicates that our performance can be further improved just by collecting more data. More experimental results are shown in the Supplementary Materials.

Table 2: Results of average HR measurement on MAHNOB-HCI.

Method	HR (bpm)			
	SD	MAE	RMSE	R
Poh2011 [39]	13.5	-	13.6	0.36
CHROM [10]	-	13.49	22.36	0.21
Li2014 [24]	6.88	-	7.62	0.81
SAMC [52]	5.81	-	6.23	0.83
SynRhythm [34]	10.88	-	11.08	-
HR-CNN [49]	-	7.25	9.24	0.51
DeepPhys [8]	-	4.57	-	-
PhysNet [60]	7.84	5.96	7.88	0.76
STVEN+rPPGNet [61]	5.57	4.03	5.93	0.88
End-to-end (baseline)	7.39	7.47	8.63	0.70
Meta-rPPG (inductive)	7.91	7.42	8.65	0.74
Meta-rPPG (proto only)	6.89	6.05	6.71	0.77
Meta-rPPG (synth only)	5.09	3.88	4.02	0.81
Meta-rPPG (proto+synth)	4.90	3.01	3.68	0.85

Table 3: Results of average HR measurement on UBFC-rPPG.

Method	HR (bpm)			
	SD	MAE	RMSE	R
GREEN [53]	20.2	10.2	20.6	-
ICA [39]	18.6	8.43	18.8	-
CHROM [10]	19.1	10.6	20.3	-
POS [55]	10.4	4.12	10.5	-
3D CNN [4]	8.55	5.45	8.64	-
End-to-end (baseline)	13.70	12.78	13.30	0.27
Meta-rPPG (inductive)	14.17	13.23	14.63	0.35
Meta-rPPG (proto only)	9.17	7.82	9.37	0.48
Meta-rPPG (synth only)	11.92	9.11	11.55	0.42
Meta-rPPG (proto+synth)	7.12	5.97	7.42	0.53

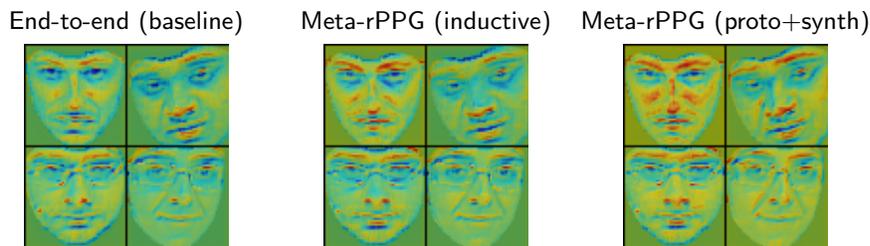


Fig. 4: Feature activation map visualization of 4 subjects using different training methods. Usage of transductive inference results in activations of higher contrast and covers larger region of facial features that contributes to rPPG estimation.

Visualization of feature activation map shows the importance of transductive inference during evaluation. In Fig. 10, we show feature activation maps visualization [28] generated using our proposed methods and supervised end-to-end trained model on images extracted from MAHNOB-HCI. Regions that are more likely to be useful for remote HR estimation are more pronounced when our approach is introduced. Results show that transductive inference is useful when applied to data excluded from the training distribution.

5 Conclusion

In this work, we introduce transductive inference into the framework of rPPG estimation. For transductive inference, we propose the use of a synthetic gradient generator and prototypical distance minimizer to provide gradient to our feature extractor when labeled data are unobtainable. By posing the learning of our network as a meta-learning framework, we see substantial improvements on MAHNOB-HCI and UBFC-rPPG dataset demonstrating state-of-the-art results.

Acknowledgements. This work is supported by Ministry of Science and Technology (MOST) of Taiwan: 107-2221-E-009 -125 -MY3.

References

1. Balakrishnan, G., Durand, F., Guttag, J.: Detecting pulse from head motions in video. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3430–3437 (2013)
2. Bengio, Y., Deleu, T., Rahaman, N., Ke, R., Lachapelle, S., Bilaniuk, O., Goyal, A., Pal, C.: A meta-transfer objective for learning to disentangle causal mechanisms. arXiv preprint arXiv:1901.10912 (2019)
3. Bobbia, S., Macwan, R., Benezeth, Y., Mansouri, A., Dubois, J.: Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters* **124**, 82–90 (2019)
4. Bousefsaf, F., Pruski, A., Maaoui, C.: 3d convolutional neural networks for remote pulse rate measurement and mapping from facial video. *Applied Sciences* **9**(20), 4364 (2019)
5. Bradski, G.: The OpenCV Library. *Dr. Dobb’s Journal of Software Tools* (2000)
6. Cao, W., Mirjalili, V., Raschka, S.: Rank-consistent ordinal regression for neural networks. arXiv preprint arXiv:1901.07884 (2019)
7. Cennini, G., Arguel, J., Akşit, K., van Leest, A.: Heart rate monitoring via remote photoplethysmography with motion artifacts reduction. *Optics express* **18**(5), 4867–4875 (2010)
8. Chen, W., McDuff, D.: Deepphys: Video-based physiological measurement using convolutional attention networks. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 349–365 (2018)
9. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: 2005 IEEE computer society conference on computer vision and pattern recognition (CVPR’05). vol. 1, pp. 886–893. IEEE (2005)
10. De Haan, G., Jeanne, V.: Robust pulse rate from chrominance-based rppg. *IEEE Transactions on Biomedical Engineering* **60**(10), 2878–2886 (2013)
11. Digiglio, P., Li, R., Wang, W., Pan, T.: Microflotronic arterial tonometry for continuous wearable non-invasive hemodynamic monitoring. *Annals of biomedical engineering* **42**(11), 2278–2288 (2014)
12. Dou, Q., de Castro, D.C., Kamnitsas, K., Glocker, B.: Domain generalization via model-agnostic learning of semantic features. In: Advances in Neural Information Processing Systems. pp. 6450–6461 (2019)
13. Doyle, O.M., Westman, E., Marquand, A.F., Mecocci, P., Vellas, B., Tsolaki, M., Kłoszewska, I., Soininen, H., Lovestone, S., Williams, S.C., et al.: Predicting progression of alzheimers disease using ordinal regression. *PloS one* **9**(8) (2014)
14. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 1126–1135. JMLR. org (2017)
15. Finn, C., Rajeswaran, A., Kakade, S., Levine, S.: Online meta-learning. arXiv preprint arXiv:1902.08438 (2019)
16. Gidaris, S., Komodakis, N.: Dynamic few-shot visual learning without forgetting. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4367–4375 (2018)
17. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8), 1735–1780 (1997)
18. Hu, S.X., Moreno, P., Xiao, Y., Shen, X., Obozinski, G., Lawrence, N., Damianou, A.: Empirical bayes transductive meta-learning with synthetic gradients. In: International Conference on Learning Representations (ICLR) (2020), <https://openreview.net/forum?id=Hkg-xgrYvH>

19. Jaderberg, M., Czarnecki, W.M., Osindero, S., Vinyals, O., Graves, A., Silver, D., Kavukcuoglu, K.: Decoupled neural interfaces using synthetic gradients. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 1627–1635. JMLR. org (2017)
20. Kazemi, V., Sullivan, J.: One millisecond face alignment with an ensemble of regression trees. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1867–1874 (2014)
21. Koch, G., Zemel, R., Salakhutdinov, R.: Siamese neural networks for one-shot image recognition. In: ICML deep learning workshop. vol. 2. Lille (2015)
22. Lee, E., Hsu, T.J., Lee, C.Y.: Centralized state sensing using sensor array on wearable device. In: 2019 IEEE International Symposium on Circuits and Systems (ISCAS). pp. 1–5. IEEE (2019)
23. Li, X., Alikhani, I., Shi, J., Seppanen, T., Junntila, J., Majamaa-Voltti, K., Tulppo, M., Zhao, G.: 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). pp. 242–249. IEEE (2018)
24. Li, X., Chen, J., Zhao, G., Pietikainen, M.: Remote heart rate measurement from face videos under realistic situations. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4264–4271 (2014)
25. Liang, S., Li, Y., Srikant, R.: Enhancing the reliability of out-of-distribution image detection in neural networks. arXiv preprint arXiv:1706.02690 (2017)
26. Liu, Y., Lee, J., Park, M., Kim, S., Yang, E., Hwang, S.J., Yang, Y.: Learning to propagate labels: Transductive propagation network for few-shot learning. arXiv preprint arXiv:1805.10002 (2018)
27. Maeda, Y., Sekine, M., Tamura, T.: The advantages of wearable green reflected photoplethysmography. *Journal of medical systems* **35**(5), 829–834 (2011)
28. Menikdiwela, M., Nguyen, C., Li, H., Shaw, M.: Cnn-based small object detection and visualization with feature activation mapping. In: 2017 International Conference on Image and Vision Computing New Zealand (IVCNZ). pp. 1–5. IEEE (2017)
29. Mishra, N., Rohaninejad, M., Chen, X., Abbeel, P.: A simple neural attentive meta-learner. arXiv preprint arXiv:1707.03141 (2017)
30. Moço, A.V., Stuijk, S., de Haan, G.: Skin inhomogeneity as a source of error in remote ppg-imaging. *Biomedical optics express* **7**(11), 4718–4733 (2016)
31. Munkhdalai, T., Yu, H.: Meta networks. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 2554–2563. JMLR. org (2017)
32. Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: European conference on computer vision. pp. 483–499. Springer (2016)
33. Nichol, A., Achiam, J., Schulman, J.: On first-order meta-learning algorithms. arXiv preprint arXiv:1803.02999 (2018)
34. Niu, X., Han, H., Shan, S., Chen, X.: Synrhythm: Learning a deep heart rate estimator from general to specific. In: 2018 24th International Conference on Pattern Recognition (ICPR). pp. 3580–3585. IEEE (2018)
35. Niu, X., Shan, S., Han, H., Chen, X.: Rhythmnet: End-to-end heart rate estimation from face via spatial-temporal representation. *IEEE Transactions on Image Processing* (2019)
36. Niu, Z., Zhou, M., Wang, L., Gao, X., Hua, G.: Ordinal regression with multiple output cnn for age estimation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4920–4928 (2016)

37. Parra, D., Karatzoglou, A., Amatriain, X., Yavuz, I.: Implicit feedback recommendation via implicit-to-explicit ordinal logistic regression mapping. *Proceedings of the CARS-2011* **5** (2011)
38. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 32, pp. 8024–8035. Curran Associates, Inc. (2019), <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
39. Poh, M.Z., McDuff, D.J., Picard, R.W.: Advancements in noncontact, multiparameter physiological measurements using a webcam. *IEEE transactions on biomedical engineering* **58**(1), 7–11 (2010)
40. Poh, M.Z., McDuff, D.J., Picard, R.W.: Non-contact, automated cardiac pulse measurements using video imaging and blind source separation. *Optics express* **18**(10), 10762–10774 (2010)
41. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning (2016)
42. Ren, J., Liu, P.J., Fertig, E., Snoek, J., Poplin, R., Depristo, M., Dillon, J., Lakshminarayanan, B.: Likelihood ratios for out-of-distribution detection. In: *Advances in Neural Information Processing Systems*. pp. 14680–14691 (2019)
43. Rettie, R., Grandcolas, U., Deakins, B.: Text message advertising: Response rates and branding effects. *Journal of targeting, measurement and analysis for marketing* **13**(4), 304–312 (2005)
44. Rusu, A.A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., Hadsell, R.: Meta-learning with latent embedding optimization. *arXiv preprint arXiv:1807.05960* (2018)
45. Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., Lillicrap, T.: Meta-learning with memory-augmented neural networks. In: *International conference on machine learning*. pp. 1842–1850 (2016)
46. Sigrist, M.K., Taal, M.W., Bungay, P., McIntyre, C.W.: Progressive vascular calcification over 2 years is associated with arterial stiffening and increased mortality in patients with stages 4 and 5 chronic kidney disease. *Clinical Journal of the American Society of Nephrology* **2**(6), 1241–1248 (2007)
47. Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: *Advances in neural information processing systems*. pp. 4077–4087 (2017)
48. Soleymani, M., Lichtenauer, J., Pun, T., Pantic, M.: A multimodal database for affect recognition and implicit tagging. *IEEE transactions on affective computing* **3**(1), 42–55 (2011)
49. Špetlík, R., Franc, V., Matas, J.: Visual heart rate estimation with convolutional neural network. In: *Proceedings of the British Machine Vision Conference, Newcastle, UK*. pp. 3–6 (2018)
50. Streifler, J.Y., Eliasziw, M., Benavente, O.R., Hachinski, V.C., Fox, A.J., Barnett, H.: Lack of relationship between leukoaraiosis and carotid artery disease. *Archives of neurology* **52**(1), 21–24 (1995)
51. Takano, C., Ohta, Y.: Heart rate measurement based on a time-lapse image. *Medical engineering & physics* **29**(8), 853–857 (2007)
52. Tulyakov, S., Alameda-Pineda, X., Ricci, E., Yin, L., Cohn, J.F., Sebe, N.: Self-adaptive matrix completion for heart rate estimation from face videos under real-

- istic conditions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2396–2404 (2016)
53. Verkruyse, W., Svaasand, L.O., Nelson, J.S.: Remote plethysmographic imaging using ambient light. *Optics express* **16**(26), 21434–21445 (2008)
 54. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. In: Advances in neural information processing systems. pp. 3630–3638 (2016)
 55. Wang, W., den Brinker, A.C., Stuijk, S., de Haan, G.: Algorithmic principles of remote ppg. *IEEE Transactions on Biomedical Engineering* **64**(7), 1479–1491 (2016)
 56. Wang, W., den Brinker, A.C., Stuijk, S., de Haan, G.: Robust heart rate from fitness videos. *Physiological measurement* **38**(6), 1023 (2017)
 57. Weersma, R.K., Stokkers, P.C., van Bodegraven, A.A., van Hogezaand, R.A., Verspaget, H.W., de Jong, D.J., Van Der Woude, C., Oldenburg, B., Linskens, R., Festen, E., et al.: Molecular prediction of disease risk and severity in a large dutch crohns disease cohort. *Gut* **58**(3), 388–395 (2009)
 58. Wu, Y., Rosca, M., Lillicrap, T.: Deep compressed sensing. arXiv preprint arXiv:1905.06723 (2019)
 59. Yu, H., Sun, S., Yu, H., Chen, X., Shi, H., Huang, T.S., Chen, T.: Foal: Fast online adaptive learning for cardiac motion estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4313–4323 (2020)
 60. Yu, Z., Li, X., Zhao, G.: Remote photoplethysmograph signal measurement from facial videos using spatio-temporal networks. In: Proc. BMVC. pp. 1–12 (2019)
 61. Yu, Z., Peng, W., Li, X., Hong, X., Zhao, G.: Remote heart rate measurement from highly compressed facial videos: an end-to-end deep learning solution with video enhancement. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 151–160 (2019)
 62. Zintgraf, L.M., Shiarlis, K., Kurin, V., Hofmann, K., Whiteson, S.: Fast context adaptation via meta-learning. arXiv preprint arXiv:1810.03642 (2018)

Algorithm 2 Transductive Inference During Deployment

```

1: Input:  $\mathbf{x}$ : A single video stream
2:  $\hat{\mathbf{x}}, \tilde{\mathbf{x}} \leftarrow \mathbf{x}$  ▷  $\hat{\mathbf{x}}$ : first 2 seconds of video,  $\tilde{\mathbf{x}}$ : rest of video
3: for  $i \leftarrow 1, L$  do ▷ Adaptation phase (run  $L$  steps)
4:    $\theta \leftarrow \theta - \alpha(\nabla_{\theta} \mathcal{L}_{\text{proto}}(\hat{\mathbf{z}}, \hat{\mathbf{z}}^{\text{proto}}) + f_{\psi}(\hat{\mathbf{z}}))$ 
5: end for
6:  $\mathbf{y} \leftarrow h_{\phi}(f_{\theta}(\tilde{\mathbf{x}}))$  ▷ Estimation of rPPG signal using adapted feature extractor

```

A Performing Transductive Inference During Deployment

Here, we show the algorithm for transductive inference during deployment in Algorithm 2. The inference process is similar to the typical inference of a feed-forward deep neural network (the final line). The difference is in the inclusion of adaptation steps using the first 2 seconds of the video for transductive learning prior to actual inference.

B Performance Comparison Using Different Adaptation Steps

In this section, we study how the number of adaptation steps, L , used for transductive inference affects performance. We report results under different metrics, namely, mean absolute error (MAE), standard deviation of error (SD), root mean squared error (RMSE) and Pearson correlation coefficient (R). Tabulated performances for MAHNOB-HCI are shown in Table 4, 5, 6 and 7 for MAE, SD, RMSE and R respectively. In the same order, tabulated performances for UBFC-rPPG are shown in Table 8, 9, 10 and 11. Again, using the same order, we show comparison plots in Figure 5, 6, 7 and 8 for both MAHNOB-HCI and UBFC-rPPG.

From the results, we can deduce that the number of adaptation steps used during transductive inference should match the number of steps used during training. This rule only applies to the generation of synthetic gradients for transductive inference but not on the prototypical distance minimizer. We hypothesize that the prototypical distance minimizer will eventually converge to some value and doesn't hurt performance if run for infinite number of adaptation steps. This is intuitive as the idea of prototypical distance minimizer is to pull out-of-distribution samples towards the center of the distribution that is modeled by the rPPG estimator using the training data.

C Does Joint Adaptation of both Feature Extractor and rPPG Estimator Give Better Results?

We hypothesize that the estimation of rPPG signal is more efficient if we are able to update the weights of the feature extractor during testing for the adaptation to the new observed distribution. By doing so, we expect that the features

generated by the feature extractor fall within the distribution covered by the rPPG estimator, i.e. the weights of the rPPG estimator is obtained through the optimization on the training dataset. One might challenge that the joint adaptation of both feature extractor and rPPG estimator weights might result in better performance, contradicting our hypothesis. To demonstrate that our hypothesis holds, we perform an empirical study on whether the joint update approach or the sole update of the weights of the feature extractor performs better. The implementation is straight forward, the synthetic gradient generator is moved towards the output for the joint adaptation case. We show experiments on MAHNOB-HCI using an adaptation steps of 10 in Table 12. The empirical results support our hypothesis.

Table 4: Results of mean absolute error of HR measurement on MAHNOB-HCI using different adaptation steps, L .

Method	MAE of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	7.47	7.47	7.47	7.47	7.47
Meta-rPPG (proto only)	7.42	6.65	6.05	6.02	6.02
Meta-rPPG (synth only)	7.42	5.00	3.88	3.70	6.25
Meta-rPPG (proto+synth)	7.42	3.89	3.01	3.02	6.20

Table 5: Results of average standard deviation HR measurement error on MAHNOB-HCI using different adaptation steps, L .

Method	SD of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	7.39	7.39	7.39	7.39	7.39
Meta-rPPG (proto only)	7.91	6.08	6.95	6.89	6.86
Meta-rPPG (synth only)	7.91	5.89	5.09	4.96	7.72
Meta-rPPG (proto+synth)	7.91	4.90	3.68	4.95	7.81

D Visualization of Feature Activation Map Using Different Methods

In this section, we show the visualization for feature activation map using the methods we introduced and is compared with the baseline that uses an end-to-end supervised learning method. Ablation study is also performed by showing feature activation maps obtained using individual methods that we proposed.

Table 6: Results of root mean squared error of HR measurement on MAHNOB-HCI using different adaptation steps, L .

Method	RMSE of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	8.63	8.63	8.63	8.63	8.63
Meta-rPPG (proto only)	8.65	6.97	6.79	6.71	6.67
Meta-rPPG (synth only)	8.65	5.15	3.96	4.02	7.11
Meta-rPPG (proto+synth)	8.65	4.65	3.66	3.68	6.94

Table 7: Results of Pearson correlation coefficient of HR measurement on MAHNOB-HCI using different adaptation steps, L .

Method	R of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	0.70	0.70	0.70	0.70	0.70
Meta-rPPG (proto only)	0.74	0.77	0.77	0.77	0.79
Meta-rPPG (synth only)	0.74	0.79	0.81	0.81	0.77
Meta-rPPG (proto+synth)	0.74	0.83	0.85	0.85	0.75

Table 8: Results of mean absolute error of HR measurement on UBFC-rPPG using different adaptation steps, L .

Method	MAE of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	12.78	12.78	12.78	12.78	12.78
Meta-rPPG (proto only)	13.23	9.24	8.07	7.82	7.53
Meta-rPPG (synth only)	13.23	11.03	7.04	9.11	11.97
Meta-rPPG (proto+synth)	13.23	7.68	6.07	5.97	9.82

Table 9: Results of average standard deviation HR measurement error on UBFC-rPPG using different adaptation steps, L .

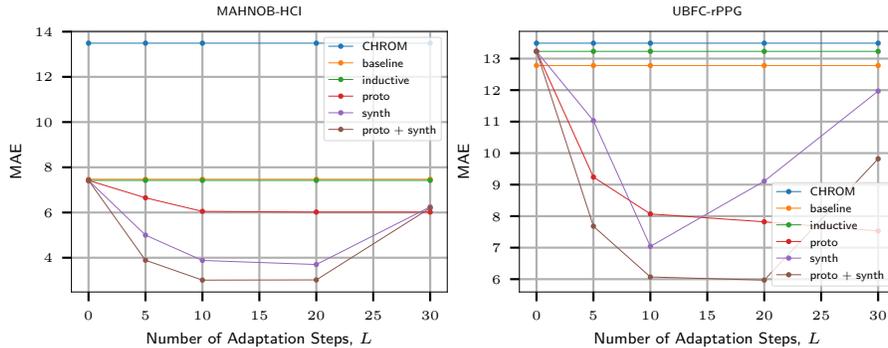
Method	SD of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	13.70	13.70	13.70	13.70	13.70
Meta-rPPG (proto only)	14.17	10.39	9.33	9.17	8.23
Meta-rPPG (synth only)	14.17	11.00	8.37	11.92	13.94
Meta-rPPG (proto+synth)	14.17	9.01	7.89	7.12	11.93

Table 10: Results of root mean squared error of HR measurement on UBFC-rPPG using different adaptation steps, L .

Method	RMSE of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	13.30	13.30	13.30	13.30	13.30
Meta-rPPG (proto only)	14.63	10.62	9.43	9.37	8.90
Meta-rPPG (synth only)	14.63	13.41	8.55	11.55	14.62
Meta-rPPG (proto+synth)	14.63	8.92	7.86	7.42	11.21

Table 11: Results of Pearson correlation coefficient of HR measurement on UBFC-rPPG using different adaptation steps, L .

Method	R of HR (bpm)				
	$L = 0$	$L = 5$	$L = 10$	$L = 20$	$L = 30$
End-to-end (baseline)	0.27	0.27	0.27	0.27	0.27
Meta-rPPG (proto only)	0.35	0.45	0.47	0.48	0.50
Meta-rPPG (synth only)	0.35	0.44	0.47	0.42	0.39
Meta-rPPG (proto+synth)	0.35	0.49	0.52	0.53	0.42

Fig. 5: Mean absolute error, MAE, obtained using different rPPG estimation methods. Demonstrates how the number of adaptation steps, L , affects performance.Table 12: Results of average HR measurement on MAHNOB-HCI comparing the difference between joint adaptation of *both feature extractor and rPPG estimator* and updating the *feature extractor only*.

Method	HR (bpm)		
	MAE	RMSE	R
Joint Adaptation	4.25	6.09	0.81
Extractor Only	3.01	3.68	0.85

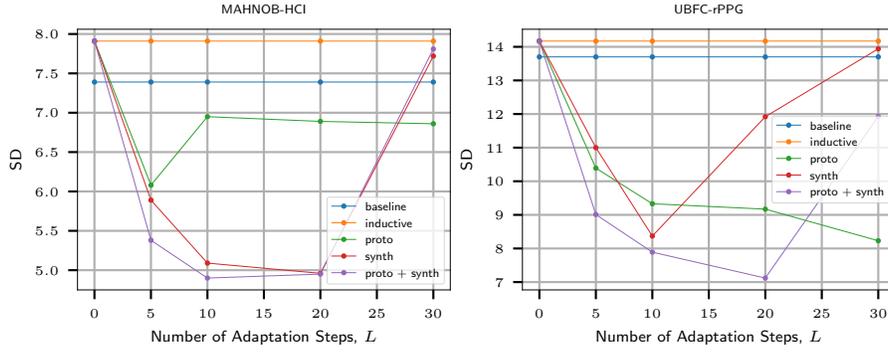


Fig. 6: Standard deviation of error, SD, obtained using different rPPG estimation methods. Demonstrates how the number of adaptation steps, L , affects performance.

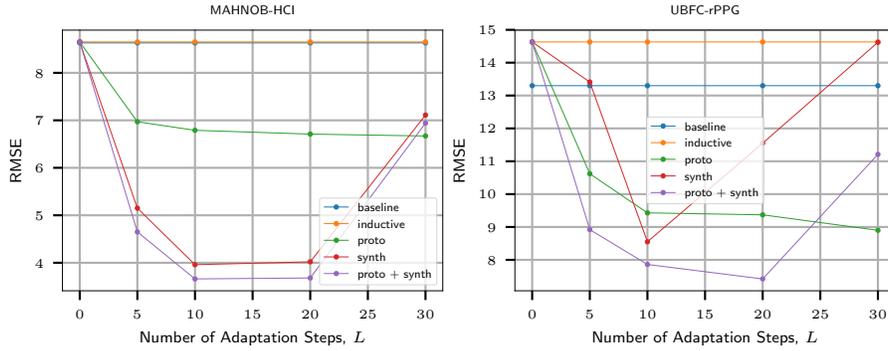


Fig. 7: Root mean squared error, RMSE, obtained using different rPPG estimation methods. Demonstrates how the number of adaptation steps, L , affects performance.

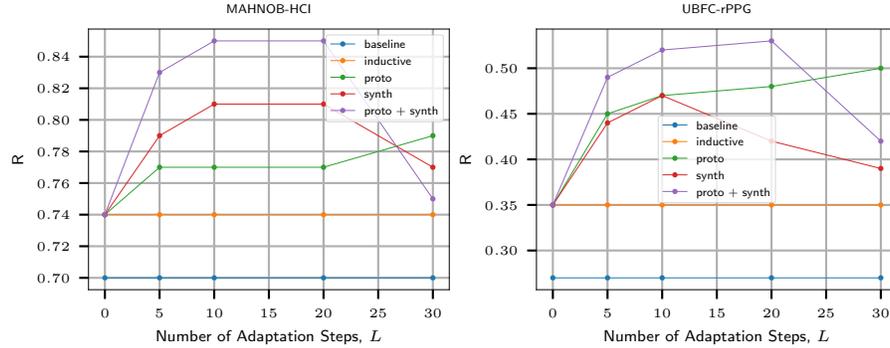


Fig 8: Pearson correlation coefficient, R , obtained using different rPPG estimation methods. Demonstrates how the number of adaptation steps, L , affects performance.

More subjects are also shown here to give a better understanding of the importance of transductive inference for rPPG estimation using a deep learning model.

E Demonstration Using Video

We show the implementation of our algorithm on videos extracted from MAHNOB-HCI to show its performance during deployment. We demonstrate our algorithm on videos of 3 subjects and a snapshot of a single frame from one of the video is shown in Fig. 9. From the video attached in the supplementary materials, it can be observed that the feature activation maps corresponding to our transductive inference method is relatively consistent when compared to a model trained in an end-to-end fashion. Please refer to our video for a better understanding on the improvements brought by the introduction of transductive inference to rPPG estimation.



Fig. 9: A frame extracted from a video from MAHNOB-HCI. From left to right: 1. Pre-processed face image (zero-ing of pixels outside facial landmarks), 2. feature activation maps corresponding to end-to-end trained model, 3. feature activation map corresponding to Meta-rPPG (proto+synth) and 4. plots containing rPPG signal (top), power spectral density of rPPG signal (middle), predicted (bottom-left) and ground truth heart rate (bottom-right) in beats-per-minute (BPM).



Fig. 10: Feature activation map visualization of 8 subjects (each row corresponds to one subject) using different training methods. Ablation study is performed by inspecting the feature map activation the results upon the application of every proposed transductive inference method. Usage of transductive inference results in activations of higher contrast and covers larger region of facial features that contributes to rPPG estimation.