

# MOMBAT: Heart Rate Monitoring from Face Video using Pulse Modeling and Bayesian Tracking

Puneet Gupta<sup>a</sup>, Brojeshwar Bhowmick<sup>b</sup>, Arpan Pal<sup>b</sup>

<sup>a</sup>Department of Computer Science and Engineering, IIT Indore, Indore, India

<sup>b</sup>Embedded system and Robotics, TCS Research and Innovation, Kolkata-700106, India

---

## Abstract

A non-invasive yet inexpensive method for heart rate (HR) monitoring is of great importance in many real-world applications including healthcare, psychology understanding, affective computing and biometrics. Face videos are currently utilized for such HR monitoring, but unfortunately this can lead to errors due to the noise introduced by facial expressions, out-of-plane movements, camera parameters (like focus change) and environmental factors. We alleviate these issues by proposing a novel face video based HR monitoring method *MOMBAT*, that is, *MONitoring using Modeling and BAYesian Tracking*. We utilize out-of-plane face movements to define a novel quality estimation mechanism. Subsequently, we introduce a Fourier basis based modeling to reconstruct the cardiovascular pulse signal at the locations containing the poor quality, that is, the locations affected by out-of-plane face movements. Furthermore, we design a Bayesian decision theory based HR tracking mechanism to rectify the spurious HR estimates. Experimental results reveal that our proposed method, *MOMBAT* outperforms state-of-the-art HR monitoring methods and performs HR monitoring with an average absolute error of 1.329 beats per minute and the Pearson correlation between estimated and actual heart rate is 0.9746. Moreover, it demonstrates that HR monitoring is significantly improved by incorporating the pulse modeling and HR tracking.

*Keywords:* Heart rate monitoring, Face video, Remote PPG, Heart rate tracking, Pulse modeling

---

*Email addresses:* puneet@iiti.ac.in (Puneet Gupta), b.bhowmick@tcs.com (Brojeshwar Bhowmick), arpan.pal@tcs.com (Arpan Pal)

## 1. Introduction

Heart rate (HR) is given by the total number of times a heart contracts or beats per minute. It can assess the human pathological and physiological parameters [30], thus it has attracted the fields of: i) healthcare; ii) psychology understanding of stress and mental state; iii) affective computing for understanding human emotion; and iv) biometrics for liveness and spoof detection. These fields can be benefited if HR monitoring is accurate and acquired from an inexpensive sensor in a user-friendly and non-contact manner. This motivates us to propose an accurate HR monitoring method using face videos in this paper.

HR can be measured using contact or non-contact mechanisms. Contact mechanisms require the sensors like electrocardiography (ECG) or photo-plethysmography (PPG). They can mitigate illumination artifacts and provide synchronized multi-modal physiological data, but they should be properly placed on the body. In real-world scenarios, motion can change the contact area between the sensor and human skin, which eventually results in spurious HR monitoring [55]. Since these methods require the contact between the user and the sensor for large duration, they restrict unobtrusive monitoring and require a dedicated sensor for single user monitoring. Also, maintaining the sensor contact is cumbersome for: i) neonates surveillance; ii) analyzing sleep quality; iii) exercise monitoring during rehabilitation etc. [43] [9]; and iv) observing skin damaged patients. These issues can be handled by performing HR monitoring using non-contact mechanisms, which allow the monitoring anytime and anywhere with minimal user involvement. These non-contact mechanisms can also be used for covert monitoring and thereby utilized for sleep monitoring, lie detection [31] and stress monitoring [41]. Due to these advantages, non-contact based HR monitoring is proliferating.

Traditional non-contact mechanisms require bulky, expensive and dedicated sensors like Microwave Doppler and laser for HR monitoring [23]. Modern non-contact mechanisms employ inexpensive and portable camera sensors for HR estimation. They are based on the phenomenon that heart beats generate the cardiovascular pulse which propagates in the entire human body. It introduces color variations in the reflected light

[14] and micro-movements in the face [1]. Both these contain the cardiovascular pulse information and are imperceptible to the human eye, but they can be analyzed using the camera for estimating the HR.

Existing face based HR methods analyze the micro-motion or color variations across time and refer to them as temporal signals [14]. The cardiovascular pulse is estimated from the temporal signals and it is eventually used for HR estimation. Along with the subtle pulse signal, the temporal signal constitutes prominent noise originated from: i) facial expression; ii) eye blinking; iii) face movements; iv) respiration; v) camera parameters (for example, change in focus); and vi) environmental factors (for example, illumination variations). Extraction of HR signal from such a noisy temporal signal is thus a challenging problem. In this paper, we alleviate these issues to improve the face videos based HR monitoring by introducing a novel method *MOMBAT*, that is, MOnitoring using Modeling and BAYesian Tracking. The main research contributions of our proposed method, *MOMBAT* are: i) it introduces a novel quality estimation mechanism that adapts according to the out-of-plane face movements and provide quality of each frame, unlike existing mechanisms that provide single quality for the entire video; ii) it initiates the utilization of Fourier basis based pulse modeling for reconstructing the pulse signals at the poor quality video frames using the pulse signals at the good quality video frames; and iii) it presents a novel Bayesian decision framework for rectifying the spurious HR estimates.

The paper is organized as follows. The preliminaries required for better understanding of our method, *MOMBAT* are discussed in Section 2 and *MOMBAT* is presented in Section 3. The experimental results are analyzed in Section 4 followed by conclusions in the last section.

## **2. Preliminaries**

### *2.1. Face HR Estimation*

Typically, any face videos based HR estimation method consists of the following three stages; preprocessing, HR estimation, and post-processing.

### *2.1.1. Preprocessing*

During preprocessing, a region of interest (ROI) containing useful pulse information is detected. Skin pixels contain pulse information, thus face detection followed by removing non-skin pixels are performed for ROI extraction [37]. Usually, a face is detected using Viola-Jones [49] or model based [3] face detectors. Subsequently, non-skin pixels due to background and hairs, are removed by applying skin color discrimination techniques. Inevitable movements (like eye blinking) near the eye areas can degrade the HR estimation [15]. Thus, the eye areas are detected by employing facial geometry heuristics or trained classifiers [51] and then these eye areas are removed for the better estimation. The remaining face area is used to define the region of interest (ROI). Some commonly used ROI are full face, forehead region or cheek areas [37]. ROI locations can be shifted by the facial movements in z-direction known as out-of-plane transformations or movements in x and y-dimensions known as in-plane transformations. Both these transformations can result in spurious HR estimation due to ROI shifting. Hence, these transformations are minimized using face registration for improving the HR estimation [17]. One can use mobile based 3D depth estimation also to get the depth of landmark [6] [28] for compensating out-of-plane movements. We use simple distance between the eyes is used by [17] for the registration, but it can be spurious due to eye-blinking. These transformations can be accurately measured by wearables [7], but it requires human contact and thus, avoided for non-contact face video based HR.

### *2.1.2. Temporal Signal Extraction*

Micro-motion and subtle color variations in the face video can be determined using Lagrangian [1] and Eulerian techniques [34] respectively. These variations across different frames provide temporal signals. In Lagrangian techniques, discriminating features are extracted from the ROI and they are explicitly tracked in the subsequent frames for determining the temporal signals [45]. This tracking is not only time-consuming, but also spurious due to improper illumination. Alternatively, temporal signals can be determined using Eulerian techniques, where color variations are examined in the fixed ROI [34]. The Eulerian techniques are less time-consuming than the Lagrangian tech-

niques, but they are applicable only when small variations are present [53]. It requires fixed ROI and hence, altered tremendously even if the face is slightly moved [16].

Eulerian temporal signals are given by the color variations in the face video, having RGB color channels. Amongst these channels, the green channel contains the strongest photo-plethysmographic signal because: i) haemoglobin absorbs green light better than red, which makes green light less susceptible to motion noise as compared to red light; and ii) green light penetrates sufficiently deeper into the skin as compared to blue light [48]. It is apparent that better performance can be expected by fusing all RGB color channels. Model based methods utilize optical and physiological properties of skin reflection to perform such a fusion. Unfortunately, such methods are not applicable in all possible scenarios. For example, well known model based methods, CHROM [11] and POS [52] do not provide correct HR estimation when pulse signal and noise share similar amplitudes [52]. Furthermore, POS fails when face videos are acquired under inhomogeneous illumination conditions, that is, when faces are illuminated by multiple light sources [52].

### *2.1.3. HR Estimation*

Pulse signal is estimated from the temporal signal using statistical learning. As an instance, periodicity analysis and blind source separation (BSS) techniques are used for the pulse signal estimation by [29] and [17] respectively. Usually, Fast Fourier Transform (FFT) is applied to the pulse signal and the frequency corresponding to the maximum amplitude in the pulse spectrum corresponds to the HR [8]. But when the temporal signal is contaminated with noise, several spurious peaks are generated and the actual HR may not correspond to the maximum amplitude peak. An example is shown in Figure 1 where several spurious peaks are generated due to facial movements. Several filtering techniques can be employed to remove the noise in the temporal signals and thereby improve HR estimation. For example, Detrending filter is applied to alleviate the non-stationary trend in the pulse signal [46]. Some spectrum subtraction methods that mitigate the noise from the pulse signal are proposed by [22, 27]; and [26]. The noise due to motion artifact is estimated by [22] using facial boundary tracking. Such tracking is spurious due to facial pose variations. Similarly, background

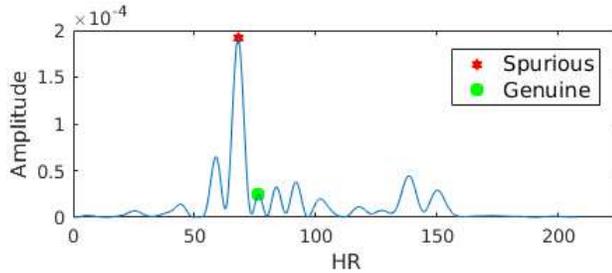


Figure 1: Spectrum of the temporal signal containing noise. HR (in bpm) and their corresponding amplitude are depicted using X and Y-axes respectively.

variations and brightness are estimated by [27] and [26] respectively, for estimating the noise due to illumination variations. But they are highly dependent on the background characteristics and the distance between the background and face [25].

## 2.2. HR Monitoring

HR monitoring continuously performs HR estimations at different small time intervals and, usually, facial deformations affect a small number of frames. Due to this, some HR estimates in the monitoring can be spurious due to the inevitable facial deformations. Better estimation can be expected when a large number of frames are considered [47]. But it results in the loss of HR variations which is highly useful for medical purposes [5]. Furthermore, it restricts user-friendliness due to high wait time. Typically, the number of frames in a time interval is chosen such that the cardiovascular pulse wave can complete at least two cycles.

HR monitoring is performed by [39] using green channel variations and band-pass filtering. Likewise, methods [47] and [35] perform HR monitoring using matrix completion and convolution neural network (CNN) respectively. Erroneous HR estimates are rectified by [17] to improve the HR monitoring using image registration and global HR. The global HR is estimated from all the video frames and thus, it can be spurious when temporal signals contain noise.

### 2.3. Constrained Local Neural Field (CLNF)

Detecting the discriminatory facial features is an extensively studied research topic. These are referred to as facial landmark points [50]. Usually, they are located around face boundaries, eyes, eyebrows, mouth and nose. Constrained Local Model (CLM) is highly useful for landmark detection. It consists of: i) point distribution model (PDM) that uses rigid and non-rigid shape transformations for modelling the global location of discriminatory points; ii) patch experts which models the behaviour of a landmark by analysing the appearance around its local neighbourhood; and iii) joint optimization which aims to fit PDM and the experts in the best possible way [42]. The unknown shape parameter  $\mathbf{p}$  is estimated by the joint optimization, which is given by

$$\mathbf{p}^* = \arg \min_{\mathbf{p}} \left[ R(\mathbf{p}) + \sum_{i=1}^n D_i(\mathbf{x}_i, I) \right] \quad (1)$$

where  $R$  is the regularization term which restricts the introduction of unlikely shapes;  $D_i$  is the misalignment in the location of  $i^{th}$  landmark in the image  $I$ ; and  $\mathbf{x}_i$  is the  $i^{th}$  landmark location in 3-D which is given by

$$\mathbf{x}_i = \mathbf{s} \cdot \mathbf{R} \cdot (\bar{\mathbf{x}}_i + \phi_i \mathbf{q}) + \mathbf{t} \quad (2)$$

where  $\bar{\mathbf{x}}_i$  denotes the mean value of  $i^{th}$  feature given by PDM;  $\phi_i$  is the component matrix; and vector  $\mathbf{q}$  is used to control the non-rigid shape [2]. Remaining parameters scaling  $\mathbf{s}$ , translation  $\mathbf{t}$  and rotation  $\mathbf{R}$  controls the rigid shape. In essence, shape parameters are given by  $\mathbf{p} = [\mathbf{s}, \mathbf{t}, \mathbf{R}, \mathbf{q}]$ . The performance of CLM heavily relies on PDM, patch expert and joint optimization. In CLNF [3], patch experts are given by local neural field for modeling spatial relationships between pixels, while non-uniform regularized landmark mean-shift is proposed for joint optimization by taking into account the reliability of patch experts.

### 2.4. Pulse Extraction using Kurtosis Optimization

Each temporal signal contains a pulse signal along with the noise. In case of multiple temporal signals, the pulse signal is extracted using blind source separation by estimating the individual source components [34]. Amplitudes of pulse signal and noise in the temporal signals depend on the facial structure, user characteristics (like

skin color) and environmental settings (like illumination). Hence, z-score normalization [40] is applied to normalize the temporal signals. Moreover, the temporal signal,  $F^i$  contains noise,  $\boldsymbol{\eta}$  and actual pulse signal,  $\mathbf{X}_a$  but modified by the facial structure. That is,

$$F^i(n) = A\mathbf{X}_a(n) + \boldsymbol{\eta}(n) \quad (3)$$

where  $n$  and  $A$  denote the time instant and matrix incorporating the effects of facial structure respectively. Further, the actual pulse signal is not known and it requires estimation from the temporal signal, that is,

$$\mathbf{X}_e(n) = BF^i(n) \quad (4)$$

where  $\mathbf{X}_e$  and  $B$  denote the estimated actual pulse and the transformation matrix respectively. It can be observed from Equations (3) and (4) that:

$$\mathbf{X}_e(n) = T\mathbf{X}_a(n) + \hat{\boldsymbol{\eta}}(n) \quad (5)$$

such that  $T = BA$  and  $\hat{\boldsymbol{\eta}} = B\boldsymbol{\eta}$ .

For accurate HR monitoring,  $\mathbf{X}_e$  should be similar to  $\mathbf{X}_a$ . That is, magnitude of  $T$  should be 1 and appropriate shape constraints should be imposed on the estimated pulse spectrum. Such shape constraints are imposed using higher order cumulants [38]. The highest order of cumulant is restricted to 4 because higher-order cumulants are easily affected by the tail of the distribution which makes them sensitive to outliers and they are slightly independent in the middle of the distribution containing useful information [24]. It is proved in [32] that constraints on cumulant similarities till fourth order can be fulfilled by defining the objective function as:

$$\max_T |K[\mathbf{X}_e]| \quad \text{subject to } T^*T = 1 \quad (6)$$

where  $K[\bullet]$  denotes the Kurtosis [12] while  $|\bullet|$  and  $*$  represent the absolute value and conjugate operations respectively. This Kurtosis based maximization is solved using [32] to obtain the estimated pulse signal because it quickly provides the global convergence.

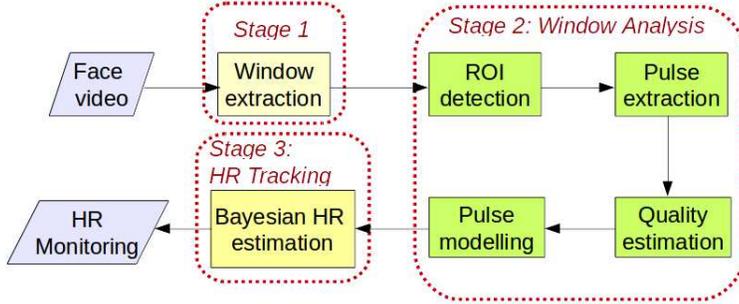


Figure 2: Flow-graph of our Proposed Method, *MOMBAT*

### 2.5. Quality estimation of pulse signal

Quality of pulse signal can be estimated using the peak signal to noise ratio (PSNR) [54]. Typically, the amplitude of the pulse spectrum obtained after converting the pulse signal into the frequency domain, should contain a peak at the HR frequency and negligible values at other frequencies. Unfortunately, in the pulse spectrum, the noise increases the amplitude at other frequencies. Thus, PSNR can be defined such that the signal can be interpreted as the amplitudes corresponding to HR frequency while noise can be thought of the amplitudes at the remaining frequencies. Mathematically, the quality given by PSNR,  $q$  is given by:

$$q = \frac{\sum_{i=\maxLoc(\mathbf{S}_p)-n_p}^{\maxLoc(\mathbf{S}_p)+n_p} \mathbf{S}_p(i)}{\text{sum}(\mathbf{S}_p) - \sum_{i=\maxLoc(\mathbf{S}_p)-n_p}^{\maxLoc(\mathbf{S}_p)+n_p} \mathbf{S}_p(i)} \quad (7)$$

where  $\mathbf{S}_p$  denotes the spectrum of the estimated pulse signal;  $\text{sum}$  performs the sum over all the frequencies;  $n_p$  represent the neighbourhood size; and  $\maxLoc$  returns the position containing the maximum value (thus, the location of HR frequency is given by  $\maxLoc(\mathbf{S}_p)$ ). In equation (7), signal (or numerator) is obtained by adding the amplitude of HR frequency and its few neighbourhoods while noise (or denominator) is obtained by adding the amplitude of the remaining frequencies.

## 3. Proposed Method

Our proposed face based HR monitoring method, *MOMBAT* is presented in this section. It consists of the following three stages: window extraction, window analysis

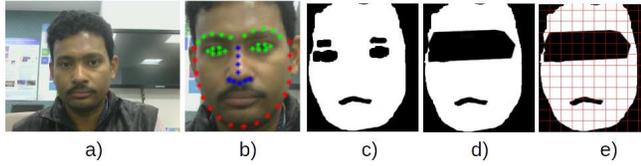


Figure 3: Steps in ROI Detection: a) Video frame; b) Landmarks on the face using CLNF; c) Detected skin mask; d) Resultant skin mask; and e) Selected ROIs from resultant mask. (Figures best viewed in colors)

and HR tracking. In the first stage, we divide the face video into several overlapping windows. In the next stage, we estimate the cardiovascular pulse and quality for each window. Subsequently, we introduce pulse signal modeling to obtain better HR estimates. In the last stage, we propose Bayesian tracking to improve the HR monitoring. Figure 2 illustrates the flow-graph of the proposed method, *MOMBAT*.

### 3.1. Window Extraction

HR monitoring requires the estimation of multiple HR at various time intervals and eventually concatenation of all HR estimates. Hence, just like the existing HR monitoring methods, we divide the face video into multiple overlapping windows [17].

### 3.2. Window Analysis

In this section, we analyze each extracted window to estimate the corresponding cardiovascular pulse, HR and quality. Initially, we detect ROIs from the window and mitigate in-plane face movements. Then, we extract the cardiovascular pulse from the ROIs using Eulerian technique followed by FFT based analysis. Subsequently, we estimate the quality of the pulse according to their out-of-plane deformations and utilize it to rectify the pulse using pulse modeling.

#### 3.2.1. ROI Detection

The facial skin area contains useful pulse information, hence we utilize it to define ROI. Initially, we detect the facial areas and landmarks using Constrained Local Neural Field (CLNF) model proposed by [4]. HR can be spurious when non-skin pixels (like beard) and eye areas are utilized for HR estimation [14]. Thus, we detect these areas and remove them. We utilize skin detection proposed by [33] to detect the non-skin

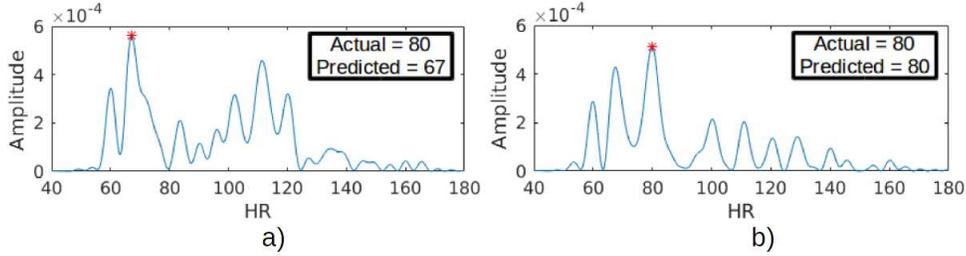


Figure 4: Usefulness of Image Registration: a) Pulse obtained without frame registration; and b) Pulse obtained after image registration. HR (in bpm) and their corresponding amplitude are depicted using X and Y-axes respectively.

pixels and we obtain the eye area by the convex hull of facial landmarks corresponding to the eyes and eyebrows. Furthermore, subtle motion in facial boundaries can significantly alter the temporal signals and thereby result in spurious HR. Hence, we remove the boundary pixels by performing morphological erosion [18]. Figure 3 illustrates these steps.

The translation and rotation of face in x and y-dimensions, known as in-plane transformations, can shift the location of the ROI in subsequent frames and thereby alter the Eulerian temporal signals and results in the spurious HR estimation. It motivates us to perform image registration between subsequent frames, so as to mitigate the in-plane transformations. An example depicting the applicability of image registration is shown in Figure 4. It shows the pulse spectrum before and after applying the image registration in Figures 4(a) and 4(b) respectively. It can be observed that HR can be correctly estimated after employing image registration. We perform the registration between subsequent frames by minimizing the deviation between nose landmark points because nose area is least affected by the facial expressions. Figure 3(b) shows the chosen landmark points in blue color. Mathematically, we first estimate the transformation matrix,  $\bar{T}$  between the current and previous frames using:

$$\bar{T} = \arg \min_T \left[ \sum_{i=1}^q \|\mathbf{F}_i - T(\mathbf{M}_i)\|_2 \right] \quad (8)$$

where  $\mathbf{M}_i$  and  $\mathbf{F}_i$  denote the positions of  $i^{th}$  nose landmark point in current and pre-

vious video frames respectively;  $q$  is the total number of nose landmark points; and  $T$  is the transformation matrix consisting of translations and rotation in 2-D, that is:

$$T = \begin{bmatrix} \cos(\theta) & -\sin(\theta) & t_x \\ \sin(\theta) & \cos(\theta) & t_y \\ 0 & 0 & 1 \end{bmatrix} \quad (9)$$

where  $\theta$  is the rotation angle while  $t_x$  and  $t_y$  are the translations in x and y directions respectively. It is important to note that  $M_i$  and  $F_i$  denotes the feature points positions in homogeneous coordinates, that is, the feature at  $[x, y]^T$  is represented by  $[x, y, 1]^T$ . We utilize Gradient Descent optimization to solve the Equation (8) [18]. The in-plane transformation is minimized by registering the current image using:

$$R = \bar{T}(I_M) \quad (10)$$

where  $R$  is the registered image and  $I_M$  is the current video frame. Thereafter, facial expressions can also result in spurious HR estimation. It can be mitigated by considering several face areas as different ROIs rather than considering full face as one ROI [14]. Hence, we utilize the method proposed by [14] for ROI extraction. For brevity, it divides the resultant registered face area into non-overlapping square blocks and considers them as ROIs. Also, it chooses the block-size such that the detected area should contain 10 blocks in the horizontal direction. An example is shown in Figure 3(e).

### 3.2.2. Pulse Extraction

We estimate the cardiovascular pulse using the method proposed by [14]. For brevity, it first extracts the Eulerian temporal signals from each ROI using the variations introduced in the average green channel intensities because the green channel contains the strongest plethysmographic signal amongst RGB color channels. Mathematically, the temporal signal  $S^i$  corresponding to  $i^{th}$  ROI is given by:

$$S^i = [s_1^i, s_2^i, \dots, s_{(f-1)}^i] \quad (11)$$

where  $f$  is the total number of frames and  $s_k^i$  representing the variations in  $k^{th}$  frame for  $i^{th}$  ROI is given by:

$$s_k^i = \sum_{(x,y) \in B^i} \left( F_{(k+1)}^g(x, y) - F_k^g(x, y) \right) \quad (12)$$

where  $B^i$  represents the  $i^{th}$  ROI;  $(x, y)$  denotes a pixel location; and  $F_k^g$  stores green channel intensities in  $k^{th}$  frame. The extracted temporal signals contain noise which is mitigated by utilizing a band-pass filter and a Detrending filter [14]. The cardiovascular pulse,  $\bar{X}_e$  is eventually extracted by applying the kurtosis based optimization proposed in [14].

### 3.2.3. Quality estimation

Just like in-plane deformations, out-of-plane deformations caused by facial movements in z-direction, can shift the ROI and result in the spurious HR estimation. We introduce a novel quality measure which incorporates these out-of-plane movements to measure the confidence in the correct estimation of pulse signal at each frame. It is defined using the 3-D facial landmarks that we have detected by applying Constrained Local Neural Field (CLNF) model [4] in Section 3.2.1. Amongst these, we utilize only the 3-D facial landmark points corresponding to the face boundary for detecting the out-of-plane movements because the face boundary is highly affected by the motion in the z-direction. These selected landmark points are shown in red color in Figure 3(b). Out of these chosen points, the points containing the largest deviation in z-direction are used for the quality estimation. It can be observed that yaw head motion can move some boundary points in positive and some in negative z-directions, thus we evaluate the deviation using maximum absolute change in z-direction. In essence, the deviation in the  $k^{th}$  frame,  $d_{(k-1)}$  is given by:

$$d_{(k-1)} = \max \left( \left| l_k^j - l_{(k-1)}^j \right| \right) \quad (13)$$

where  $\max$  is the maximum operator;  $|\bullet|$  is the absolute operator; while  $l_k^j$  and  $l_{(k-1)}^j$  denote the z-coordinate of  $j^{th}$  landmark in  $k^{th}$  and  $(k-1)^{th}$  frames respectively. After evaluating the deviations for all the frames, except the last frame, we compute the quality at  $(k-1)^{th}$  frame,  $\hat{q}_{(k-1)}$  using:

$$\hat{q}_{(k-1)} = 1 - \left( \frac{d_{(k-1)} - \min(\mathbf{d})}{\max(\mathbf{d}) - \min(\mathbf{d})} \right) \quad (14)$$

where  $\min$  is the minimum operator and  $\mathbf{d}$  stores all the computed deviations, that is

$$\mathbf{d} = [d_1, d_2, \dots, d_{f-1}] \quad (15)$$

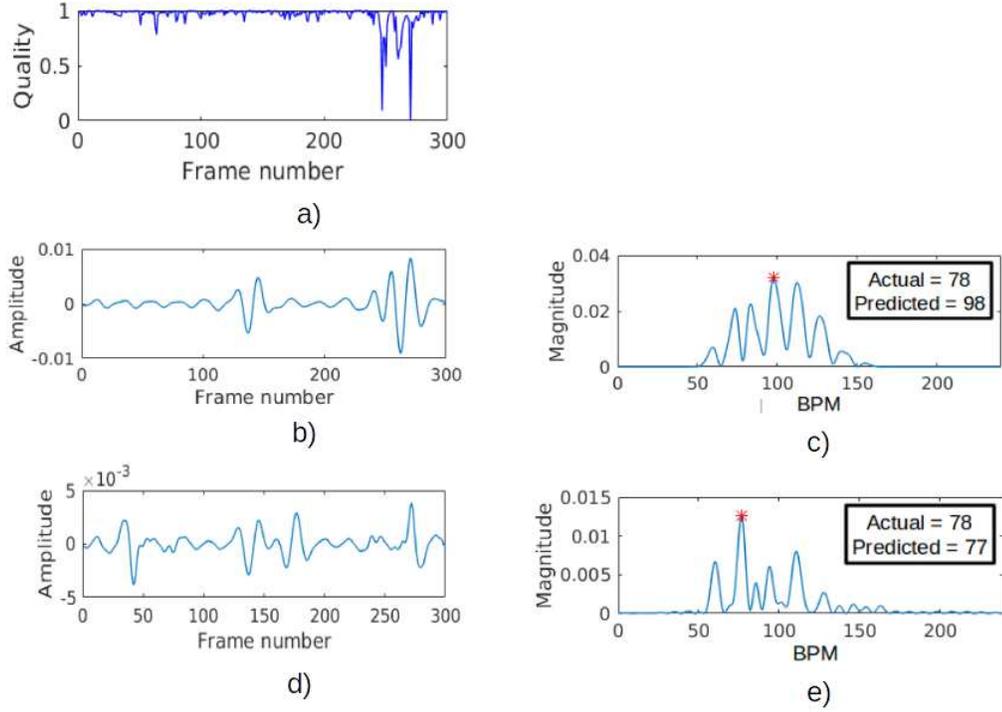


Figure 5: Steps in post-processing: a) Quality,  $\hat{Q}$ ; pulse signal,  $\bar{X}_e$  and its spectrum without post-processing are shown in b) and c) respectively; while pulse signal,  $\hat{X}_e$  and its spectrum after post-processing are shown in d) and e) respectively. For illustrating the spectrum in c) and e), we depict the HR (in bpm) and their corresponding amplitude using X and Y-axes respectively.

where  $f$  is the number of frames. In essence,  $d_{(k-1)}$  in Equation (14) is first normalized to  $[0, 1]$  and then modified to define quality such that low and high deviations corresponds to high and low quality values respectively. Thus, the quality due to out-of-plane movements,  $\hat{Q}$  is given by:

$$\hat{Q} = [\hat{q}_1, \hat{q}_2, \dots, \hat{q}_{(f-1)}] \quad (16)$$

An example of the quality estimation using out-of-plane movements is shown in Figure 5(a).

### 3.2.4. Pulse Modeling

The frames affected by out-of-plane movements provide spurious temporal signals and thereby results in an incorrect estimation of cardiovascular pulse,  $\bar{X}_e$ . An example of such pulse is shown in Figure 5(b) along with its corresponding pulse spectrum in Figure 5(c). It can be observed from Figure 5(c) that the predicted HR is deviated significantly from the actual HR. To improve the efficacy of the pulse signal, we introduce Fourier basis based modeling that aims to reconstruct the pulse signal at those frames which are affected by out-of-plane movements. We formulate the problem of noise reduction as a data fitting problem [20]. It consists of the following steps: i) defining appropriate basis functions; ii) parameter fitting; and iii) signal reconstruction. Mathematically,  $\bar{X}_e$  can be decomposed as:

$$\bar{X}_e(x) \approx \sum_{i=1}^{\alpha} a_i \phi_i(x) \quad (17)$$

where  $\alpha$  is the number of basis;  $a_i$  denotes the model parameter for  $i^{th}$  parameter;  $x$  denote the frame number; and  $\phi_i(x)$  is the  $i^{th}$  basis function. Parameter  $\alpha$  plays a crucial role in the modeling. Pulse reconstruction is spurious when  $\alpha$  is small and if it is set to high value, then even noise can be modeled. We describe the parameter selection of  $\alpha$  in Section 4.3. For simplicity, Equation (17) can be written in matrix form using:

$$\bar{X}_e \approx \mathbf{A}\Phi \quad (18)$$

where

$$\Phi = \begin{bmatrix} \phi_1(1) & \phi_2(2) & \cdot & \cdot & \phi_1(f-1) \\ \phi_2(1) & \phi_2(2) & \cdot & \cdot & \phi_2(f-1) \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \phi_\alpha(1) & \phi_\alpha(2) & \cdot & \cdot & \phi_\alpha(f-1) \end{bmatrix} \quad (19)$$

and

$$\mathbf{A} = [a_1 \ a_2 \ \cdots \ a_\alpha] \quad (20)$$

We define basis functions,  $\Phi$ , using the well known Fourier basis [20]. It is because i) these basis are orthogonal which is required to provide stability in the optimization by assuring low residual error; and ii) their amplitude lies in the range of  $[-1, 1]$

which helps in avoiding the problem of overflowing integer with polynomial basis. The Fourier basis for the order  $n$  are given by:

$$\phi_{(n)}(x) = \begin{cases} \sin\left(\frac{(n+1)}{2} \times x\right) & \text{if } n \text{ is odd} \\ \cos\left(\frac{n}{2} \times x\right) & \text{if } n \text{ is even} \end{cases} \quad (21)$$

This represents an overdetermined system of linear equations because the small number of unknown parameters,  $\alpha$  needs to be estimated from a large number of observations,  $(f - 1)$ . Furthermore, we aim to reconstruct the pulse at the frames containing large out-of-plane movements by utilizing the pulse information at the frames containing small out-of-plane movements. Thus, we solve this overdetermined system of linear equations using weighted least square estimation where weights are given by quality due to out-of-plane movements,  $\hat{Q}$  [20]. That is,

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \left\| \hat{Q} (\bar{\mathbf{X}}_e - \mathbf{A}\Phi) \right\|_2 \quad (22)$$

where  $\hat{\mathbf{A}}$  contains the estimated modeling parameters;  $\|\bullet\|$  contains the norm; and  $\hat{Q}$  (solved in Equation (16)) is the quality due to out-of-plane movements. The solution of Equation (22) is given by:

$$\hat{\mathbf{A}} = (\Phi^T \bar{Q} \Phi)^{-1} \Phi^T \bar{Q} \bar{\mathbf{X}}_e \quad (23)$$

where  $\bar{Q}$  is the diagonal matrix formed from  $\hat{Q}$  in the following manner:

$$\bar{Q} = \text{diag}(\hat{q}_1, \hat{q}_2, \dots, \hat{q}_{(f-1)}) \quad (24)$$

Modeled pulse signal,  $\hat{\mathbf{X}}_e$  is obtained by:

$$\hat{\mathbf{X}}_e = \hat{\mathbf{A}}\Phi \quad (25)$$

An example of the modeled pulse signal is shown in Figure 5(d) along with its corresponding pulse spectrum in Figure 5(e). It can be observed from the Figures 5(c) and 5(e) that the HR estimation can be improved significantly after incorporating the proposed pulse modeling.

### 3.3. HR tracking

The spectrum of noise-free pulse signal should contain maximum amplitude at the HR frequency, but it is violated when the pulse signal contains noise. The modeled

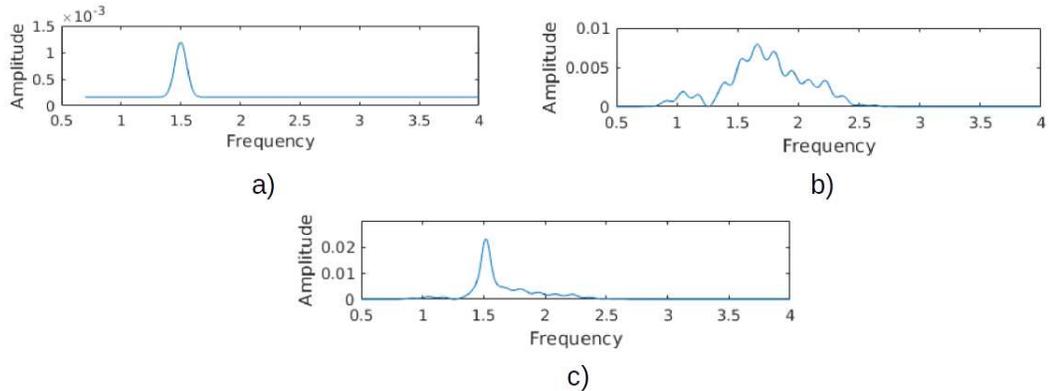


Figure 6: Example of pulse spectrum during tracking: a) Prior,  $P_e$  (for  $h_i = 1.5$  and  $\frac{c}{q_i} = 0.05$ ); b) Likelihood function,  $P_l$ ; and c) Posterior,  $P_p$ . These depict the HR frequency (in beats per second) and their corresponding amplitude using X and Y-axes respectively.

pulse signal obtained after applying the proposed pulse modeling technique, contains noise and thus, it can provide spurious HR. Usually, small HR change is observed between subsequent HR estimates. It motivates us to introduce the Bayesian framework which rectifies the spurious HR estimates. The framework consolidates the likelihood information derived from the current window and the prior information related to previously analyzed windows [13].

To leverage the observation that there is a small HR change between subsequent HR estimates, we want to define our prior information such that the large confidence value is provided when the fluctuation between the current HR and previous HR is small while a small confidence value is provided when the fluctuation is large. Furthermore, we need to provide low prior information about the previous HR whenever they are spurious, otherwise the error can be propagated in the subsequent HR estimates. An important characteristic of spurious HR is that its corresponding pulse spectrum contains multiple peaks, thus the spectrum has low PSNR [54]. These conditions are met by defining the prior information for  $(i + 1)^{th}$  window using:

$$P_g(\boldsymbol{\theta}) = \sqrt{\frac{\gamma_i}{2\pi c}} \exp\left[-\frac{1}{2} \left(\frac{\gamma_i}{c} (\boldsymbol{\theta} - h_i)\right)^2\right] \sim \mathcal{N}\left(h_i, \frac{c}{\gamma_i}\right) \quad (26)$$

where  $h_i$  and  $\gamma_i$  denote the HR frequency and PSNR of the modeled pulse spectrum

in the previous window (that is,  $i^{th}$  window) respectively;  $c$  is a predefined constant;  $\theta$  is the set of all probable HR frequencies; and  $\mathcal{N}$  denote the normal distribution. It can be observed that the normal distribution is used in Equation (26) such that mean and variance are given by  $h_i$  and  $\frac{c}{\gamma_i}$  respectively. Thus, small PSNR results in high variance which in turn results in low prior knowledge. Also, it is suggested by [55] that the fluctuations between subsequent HR estimates usually lie within the range of -11bpm to +11bpm. Thus, we set  $c$  equal to 4, so that 3 times of the variance covers most of our permissible HR estimates.

Our definition of prior information in Equation (26) prohibits large fluctuation from the previous HR. Hence, if previous HR is spurious with low PSNR value, then the current HR values should be restricted with large range. But one should consider all the plausible HR frequencies when previous HR is spurious. To incorporate this intuition, we add a constant value in all the plausible HR frequencies, which are lying between 0.7 to 4Hz. That is, we modify the prior information using:

$$P_e(\theta) = \begin{cases} \frac{P_u(\theta) + P_g(\theta)}{(\sum_{\theta=0.7}^4 (P_u(\theta) + P_g(\theta)))}, & \text{if } 0.7 < \theta < 4Hz \\ 0, & \text{otherwise} \end{cases} \quad (27)$$

where  $P_e$  denotes the modified prior information;  $P_g$  is the distribution described in Equation (26); and  $P_u$  is given by:

$$P_u(\theta) = \begin{cases} \hat{c}, & \text{if } 0.7 < \theta < 4Hz \\ 0, & \text{otherwise} \end{cases} \quad (28)$$

In essence,  $P_u$  is the uniform distribution, defined in the HR frequency ranges of 0.7 to 4 Hz such that any frequency is equally probable with the value of  $\hat{c}$ . We describe the parameter selection of  $\hat{c}$  in Section 4. Furthermore, when the first window is analyzed  $P_g$  is set to zero for all the possible HR frequency ranges, so that all the values are equally likely and hence no useful prior information is utilized.

The likelihood function is denoted by  $P_l(\mathbf{S}_{i+1}|\theta)$  where  $\mathbf{S}_{i+1}$  denotes the spectrum of reconstructed pulse signal  $\hat{\mathbf{X}}_e$  corresponding to the  $(i+1)^{th}$  window. We estimate it using:

$$P_l(\mathbf{S}_{i+1}|\theta) = \frac{\mathbf{S}_{i+1}(\theta)}{\left(\sum_{\theta=0.7}^{4Hz} \mathbf{S}_{i+1}(\theta)\right)} \quad (29)$$

The posterior probability,  $P_p(\boldsymbol{\theta}|\mathbf{S}_{i+1})$  is evaluated by applying the Bayes rule [13], that is,

$$P_p(\boldsymbol{\theta}|\mathbf{S}_{i+1}) = \frac{P_l(\mathbf{S}_{i+1}|\boldsymbol{\theta}) P_e(\boldsymbol{\theta})}{P(\mathbf{S}_{i+1})} \quad (30)$$

where  $P(\mathbf{S}_{i+1})$  is the evidence factor. Equations (27), (29) and (30) can be combined in the following manner:

$$P_p(\boldsymbol{\theta}|\mathbf{S}_{i+1}) = \frac{\mathbf{S}_{i+1}(\boldsymbol{\theta}) (P_u(\boldsymbol{\theta}) + P_g(\boldsymbol{\theta}))}{Z_2} \quad (31)$$

where  $Z_2$  is a normalization coefficient given by:

$$Z_2 = \left( \sum_{\boldsymbol{\theta}=0.7}^4 (P_u(\boldsymbol{\theta}) + P_g(\boldsymbol{\theta})) \right) P(\mathbf{S}_{i+1}) \left( \sum_{\boldsymbol{\theta}=0.7}^{4Hz} \mathbf{S}_{i+1}(\boldsymbol{\theta}) \right) \quad (32)$$

An illustration of the prior information, likelihood function and their corresponding posterior probability is shown in Figure 6. We apply maximum a posteriori estimation for HR frequency estimation which provides the minimum-error-rate classifier based on zero-one loss function [13]. For brevity, the expected loss incurred on selecting a particular frequency,  $d$  is given by:

$$R(\boldsymbol{\theta} = d|\mathbf{S}_{i+1}) = \sum_{K=0.7}^4 L(d, \boldsymbol{\theta}) P_p(\boldsymbol{\theta} = K|\mathbf{S}_{i+1}) \quad (33)$$

where  $R$  is the incurred loss and  $L$  represent the loss function given by:

$$L(d, \boldsymbol{\theta}) = \begin{cases} 0, & \text{if } \boldsymbol{\theta} = d \\ 1, & \text{otherwise} \end{cases} \quad (34)$$

Further, it is obvious that the sum of likelihood function at all the possible values (which lies between 0.7 to 4Hz in our case) will be equal to one, that is,

$$\sum_{K=0.7}^4 P_p(\boldsymbol{\theta} = K|\mathbf{S}_{i+1}) = 1 \quad (35)$$

It can be seen by combining Equations (33), (34) and (35) that:

$$R(\boldsymbol{\theta} = d|\mathbf{S}_{i+1}) = 1 - P_p(\boldsymbol{\theta} = d|\mathbf{S}_{i+1}) \quad (36)$$

Hence, expected loss,  $R$  is minimized when  $\boldsymbol{\theta}$  is set to the value that maximizes the posterior probability  $P_p$ , that is,  $\boldsymbol{\theta}$  is set to HR frequency. Hence, we obtain the HR

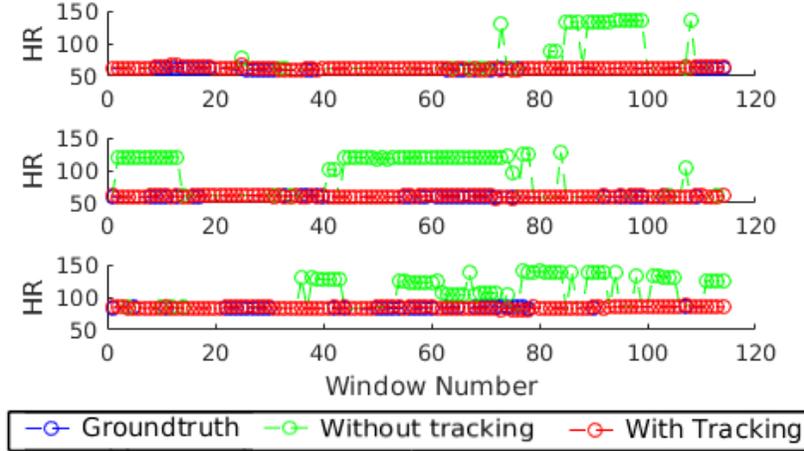


Figure 7: Examples of HR monitoring with and without tracking. X-axis denotes the window number and y-axis denotes the corresponding HR (in bpm). (Figures best viewed in colors)

frequency corresponding to  $(i + 1)^{th}$  window,  $h_{i+1}$  using:

$$h_{i+1} = \arg \max_{\theta} P_p(\theta | \mathcal{S}_{i+1}) \quad (37)$$

The corresponding HR is given by:

$$\hat{H}(i + 1) = \text{round}(h_{i+1} \times 60) \quad (38)$$

where *round* operator rounds off the value to the nearest integer. Some examples depicting the usefulness of the proposed HR tracking are shown in Figure 7. The figure depicts the actual HR monitoring along with the predicted HR monitoring when the proposed HR tracking is avoided and utilized. It demonstrates that the HR monitoring can be improved significantly when the proposed HR tracking is used.

## 4. Experimental Results

### 4.1. Data Recording of Our Dataset

The performance of our method, *MOMBAT* is evaluated on Intel i5-2400 CPU 3.10 GHz. Total 65 face videos have been collected from 65 different subjects (34 males and 31 females), out of which 15 videos are used for parameter selection (or

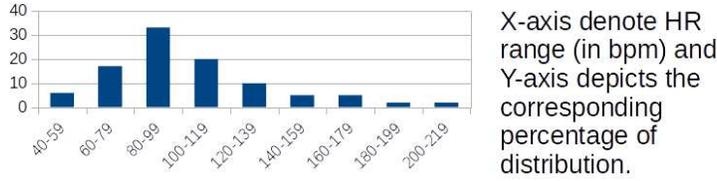


Figure 8: Distribution of HR in the database.

training) and the remaining 50 videos are used for performance evaluation (or testing). The videos are acquired from Logitech webcam C270 camera which is mounted on a laptop and the subjects are free to perform natural facial movements and head pose variations. The resolution of these acquired videos is  $640 \times 480$  pixels. Furthermore, we avoid any compression mechanism and save the videos in AVI raw format. These are acquired for 1 minute at 30 frames per second. The ground truth is obtained by simultaneously acquiring the actual pulse from the right index fingertip using CMS 50D+ pulse oximeter. The percentage of distribution of ground truth HR estimation from the acquired database is shown in Figure 8.

#### 4.2. Performance Measurement

The performance metrics used in our experiments are based on the predicted HR error,  $(\bar{P}(i, j) - \bar{A}(i, j))$ , where  $\bar{P}(i, j)$  and  $\bar{A}(i, j)$  denote the predicted and actual HR estimates respectively for  $i^{th}$  subject in  $j^{th}$  window. Accurate HR monitoring method requires that the prediction error is close to zero, alternatively, the mean  $\mu$  and standard deviation  $\sigma$  of the prediction error should be close to zero. Likewise, the percentage of samples with absolute error less than 5 bpm,  $err_5$  should be close to 100% for correct HR monitoring. Another metric employed for the evaluation is mean average error,  $MAE$  of all the subjects which is given by:

$$MAE = \frac{\sum_{i=1}^z \sum_{j=1}^{n_i} |\bar{P}(i, j) - \bar{A}(i, j)|}{\sum_{i=1}^z n_i} \quad (39)$$

where  $|\bullet|$  is the absolute operator;  $n_i$  represents the number of windows for  $i^{th}$  subject; and  $z$  is the total number of subjects. Lower value of  $MAE$  indicates that the predicted and estimated HR estimates are close to each other. Similarly, we also use total time,  $t_s$  required for HR monitoring in seconds as a performance metric. Furthermore, we used

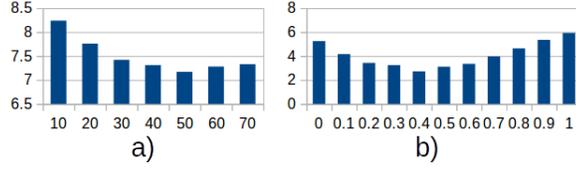


Figure 9: Performance analysis for: a)  $\alpha$  and b)  $\hat{c}$ . X and Y-axes denote the parameter value and  $MAE$  (in bpm) respectively.

the Pearson correlation coefficient,  $\rho$  to evaluate the similarity between two variables in terms of linear relationship. It lies between -1 to 1 and is given by:

$$\rho = \frac{\text{cov}(\bar{P}, \bar{A})}{\sigma(\bar{P}) \times \sigma(\bar{A})} \quad (40)$$

where  $cov$  and  $\sigma$  are the covariance and standard deviation operator respectively. Better HR estimation requires high similarity between the predicted and actual HR, that is, high  $\rho$ .

#### 4.3. Parameter Selection

Our proposed method *MOMBAT* requires proper selection of four parameters, which are: i) window size; ii) overlapping window size; iii)  $\alpha$  representing the number of basis in pulse modeling; and iv)  $\hat{c}$  which is required for defining prior in HR tracking (refer Equation (28)). Since minimum possible heart beats can be 42bpm and the proper window should contain at least two cycles of the cardiovascular pulse, we selected the window of 4 second video where slightly more than two cycles can be observed. Similarly, overlap between the successive windows is chosen in an application specific manner. We focus on frequently updating the previous HR from the new HR, typically, twice in a second. Hence, we set the overlap between successive windows at 3.5 seconds. The remaining parameters are set by the value, providing minimum  $MAE$  on the training set. The  $MAE$  for different parameter values are shown in Figure 9. We set  $\alpha$  and  $\hat{c}$  equal to 50 and 0.4 respectively, where minimum  $MAE$  is attained on the training set.

Table 1: Description of subversions of *MOMBAT*

Method	Color channel	Image registration	Pulse modeling <sup>1</sup>	Bayesian tracking
<i>NorSysR</i>	Red	No	NA	No
<i>NorSysI</i>	Green	No	NA	No
<i>NorSys</i>	Green	Yes	NA	No
<i>PulseModP</i>	Green	Yes	Polynomial	No
<i>PulseModL</i>	Green	Yes	Legendre	No
<i>PulseMod</i>	Green	Yes	Fourier	No
<i>BayTrack</i>	Green	Yes	NA	Yes
<b><i>MOMBAT</i></b>	<b>Green</b>	<b>Yes</b>	<b>Fourier</b>	<b>Yes</b>

<sup>1</sup> NA means pulse modeling is avoided. Otherwise, the basis name is mentioned.

#### 4.4. Performance Evaluation

For rigorous performance analysis, we create several other methods from our proposed method, *MOMBAT* by avoiding or replacing its components. The following methods are considered for the performance analysis: a) *NorSysI* which is obtained by avoiding image registration, pulse modeling and HR tracking in *MOMBAT*; b) *NorSysR* which is same as *NorSysI* except that it uses red light instead of green light for extracting the temporal signals; c) *NorSys* which is obtained by avoiding pulse modeling and HR tracking in *MOMBAT*; d) *PulseModP* and *PulseModL* which avoid HR tracking, but utilize polynomial and Legendre basis [19] respectively, instead of Fourier basis for pulse modeling in *MOMBAT*; e) *PulseMod* which is given by considering proposed (or Fourier basis based) pulse modeling, but avoiding HR tracking in *MOMBAT*; and f) *BayTrack* which is given by avoiding pulse modeling, but considering HR tracking from *MOMBAT*. The description of these subversions of *MOMBAT* is provided in Table 1. Furthermore, we compare our method with the following existing well known methods: [1]; CHROM [11]; POS [52]; [22]; [47]; [39]; [35]; and [17]. Pulse signal is extracted from the Lagrangian temporal signals using Principal Component Analysis in [1]. [22] registers the face and utilizes

Eulerian temporal signals. Model based methods are utilized in [11] and [52] where temporal signals are extracted by fusing RGB color channels. Optical and physiological properties of skin reflection are used to perform such a fusion. Methods [1], [22], [11] and [52] provide one HR value. To conform these methods with *MOMBAT*, we extract the window and then analyze each window using these methods for HR monitoring. We are unable to conduct the comparative analysis with [27] for HR monitoring because it requires large window size as described in [47]. In [35], CNN trained on several windows is used for HR monitoring. The training and test sets contain different windows of the same subjects. For more rigorous analysis, we also perform the experimentation with another method *ModCNN* where [35] is used, except that its training and testing sets do not contain windows of the same subjects.

#### 4.5. Performance Analysis on Our Dataset

Our experimental results on our dataset are presented in Table 2. It can be inferred from the table that [1] provides the most spurious HR monitoring because it requires the tracking of facial features, which is easily affected by expressions. Likewise, [22] exhibits lower performance than the other methods except [1] because it averages all the temporal signals for pulse extraction. This is error-prone because large noise in few temporal signals due to facial expressions, can tremendously affect the cardiovascular pulse after averaging. Both [1] and [22] employ highly time consuming feature tracking and BSS. In contrast, [39] extracts the pulse signal using only green channel intensity differences of full face and avoiding computationally expensive BSS step and feature tracking. It enables [39] to perform in the most computationally efficient manner, but such a method performs spuriously because it is easily affected by facial deformation, as shown in [17].

*NorSysR* and *NorSysI* are different only in the way that they utilize red and green light respectively for the temporal signal extraction. It can be observed from Table 2 that *NorSysI* performs better than *NorSysR*, which indicates that green light is more effective in photo-plethysmographic imaging than red light. This observation is also mentioned in [48]. Furthermore, *NorSysI* performs better than CHROM [11] and POS [52], which utilize optical and physiological properties of skin reflection to con-

Table 2: Comparative Results of HR Monitoring on our Dataset

Method	$\mu$	$\sigma$	$err_5$	$MAE$	$\rho$	$t_s$
[1]	-18.4120	27.6195	38	22.5972	-0.1793	25.72
[22]	-9.0745	20.2534	75	10.0305	0.2515	30.37
[39]	9.5317	21.0467	70	11.3885	0.3310	1.24
CHROM [11]	-8.1932	19.1917	76	9.843	0.3015	6.81
POS [52]	-8.9827	19.7920	77	10.214	0.2912	6.81
<i>NorSysR</i>	-10.8246	21.4921	72	10.946	0.2847	6.80
<i>NorSysI</i>	-6.9634	18.7418	80	8.7382	0.3106	6.80
[47]	6.8242	18.3521	81	8.1864	0.4256	19.41
<i>NorSys</i>	-6.4405	17.4389	83	7.3813	0.4486	6.84
<i>PulseModP</i>	-6.4079	17.3964	82	7.2503	0.4504	9.63
<i>PulseModL</i>	-6.3865	17.3726	83	7.2057	0.4542	9.63
<i>PulseMod</i>	-6.1783	17.1010	85	7.1853	0.4627	9.63
<i>BayTrack</i>	-0.5864	5.7052	94	2.2154	0.8821	7.06
[35]	0.9275	7.3472	90	3.4588	0.8226	9.92
<i>ModCNN</i>	-10.1539	21.6781	31	18.2169	-0.0337	9.92
[17]	0.4667	4.8230	89	2.4968	0.8601	16.81
<b><i>MOMBAT</i></b>	<b>-0.1041</b>	<b>2.6172</b>	<b>97</b>	<b>1.3293</b>	<b>0.9746</b>	<b>9.78</b>

Unit of: i)  $err_5$  is %; ii)  $\mu$ ,  $\sigma$  and  $MAE$  is bpm; and  $t_s$  is seconds.

solidate RGB color channels for temporal signal extraction. It is because CHROM and POS do not provide correct HR estimation when pulse signal and noise share similar amplitudes [52]. In addition, POS fails when face videos are illuminated by multiple light sources [52].

*NorSysI* is the same as *NorSys* except that *NorSysI* avoids image registration and Table 2 points out that *NorSys* performs better than *NorSysI*. It indicates that performance can be increased by utilizing image registration. Average computational time of *NorSysI* and *NorSys* are 6.8 second and 6.84 second, respectively, out of which, BSS is the most computationally expensive step requiring 5.13 seconds. Method *NorSys* also performs better than [47] due to better ROI selection, image registration and proper BSS technique. [47] utilizes matrix completion to mitigate the noise, which increases the computation time significantly. Also, it can be observed from the table that *PulseModP*, *PulseModL*, *PulseMod* and *BayTrack* perform better HR monitoring than *NorSys*. *PulseModP*, *PulseModL* and *PulseMod* mitigate the problems of out-of-plane movements in *NorSys* by modeling the pulse signal (refer Figure 5 for example) and in return, they incur an additional average time of 0.31 and 2.48 sec for the modeling and quality estimation respectively. *PulseMod* performs better than *PulseModL* and *PulseModP* which demonstrate that Fourier basis is better suited for pulse signal modeling. Similarly, *BayTrack* performs better monitoring than *NorSys* because it rectifies the HR estimates by incorporating the prior knowledge of the HR estimates. Some of its examples are shown in Figure 7. It incurs an additional average time of 0.22 sec than *NorSys* due to PSNR estimation.

Table 2 indicates that [35] exhibits good HR monitoring when the training and test sets contain different windows of the same subjects. But when the training and testing sets do not contain windows of the same subjects (that is, *ModCNN*) then there is significant performance degradation. It points out that CNN employed by [35], leverages the facial texture and skin color for HR monitoring. This is obviously a wrong way of performing HR monitoring. Likewise, [17] relying on only face reconstruction is incompetent to handle out-of-plane movements and hence, provide spurious HR estimates. But our method, *MOMBAT* handles most of the spurious cases by utilizing the pulse modeling and Bayesian tracking. Furthermore, it provides the best HR

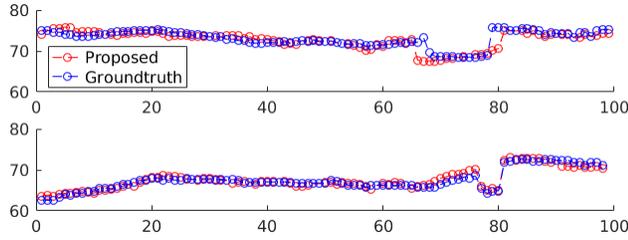


Figure 10: HR Monitoring by Our *MOMBAT* under Large HR Fluctuations. X and Y-axes denote the window number and HR (in bpm) respectively. (Figures best viewed in colors)

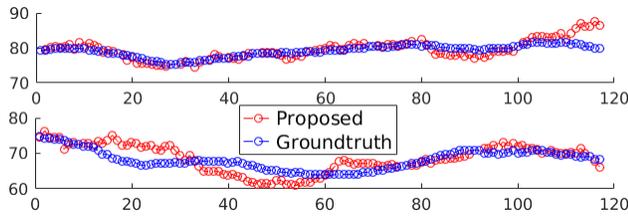


Figure 11: Erroneous HR Monitoring by Our *MOMBAT*. X and Y-axes denote the window number and HR (in bpm) respectively. (Figures best viewed in colors)

monitoring amongst all the methods. But it incurs an additional average time of 2.94 sec when compared with *NorSys* due to modeling and tracking. Such small time differences can be neglected to achieve significantly better HR monitoring for the face videos of 54 sec. A few cases where *MOMBAT* has successfully performed the HR monitoring are shown in Figure 10. Just like other existing methods, *MOMBAT* may perform spuriously when the face video contains noise that persists for long duration. Some such spurious monitoring cases by *MOMBAT* are shown in Figure 11.

#### 4.6. Performance Analysis on COHFACE

One major factor that hampers the progress in this realm of HR analysis using face video is the lack of appropriate datasets [10]. It is stated by [10] that several existing publicly available datasets that are extensively used in the literature, are not appropriate for HR estimation using face videos. One such example is MAHNOB-HCI dataset [44] which involves negligible illumination variations induced by the movie and thus, unable to cater complex real-world scenarios. COHFACE dataset is regarded

Table 3: Comparative Results of HR Monitoring on COHFACE

Method	$\mu$	$\sigma$	$err_5$	$MAE$	$\rho$
[1]	-10.1664	19.2326	70	12.3612	0.0897
[22]	9.3481	18.5818	74	11.2493	0.1264
[39]	6.8215	19.4267	67	16.3568	-0.1129
CHROM [11]	-7.9135	19.1191	75	10.384	0.1315
POS [52]	-7.1263	18.3902	77	9.621	0.1617
<i>NorSysR</i>	-7.8460	18.9212	72	10.696	0.1248
<i>NorSysI</i>	-6.2179	16.3592	80	9.0424	0.1721
[47]	6.9628	15.6330	81	8.9354	0.1813
<i>NorSys</i>	-5.3672	15.8627	83	8.7956	0.2065
<i>PulseModP</i>	-4.6730	14.7966	85	8.2718	0.2108
<i>PulseModL</i>	-4.5203	14.7171	85	8.1629	0.2256
<i>PulseMod</i>	-4.0942	13.9043	86	7.9250	0.2614
<i>BayTrack</i>	-1.0576	9.4852	90	6.4797	0.5352
[35]	2.3874	11.9163	89	6.8128	0.5036
<i>ModCNN</i>	-4.3561	18.2781	56	14.8243	-0.0696
[17]	1.4666	12.6595	88	6.5411	0.5252
<b><i>MOMBAT</i></b>	<b>-0.9832</b>	<b>7.3823</b>	<b>92</b>	<b>5.8923</b>	<b>0.6184</b>

Unit of  $err_5$  is % while unit of  $\mu$ ,  $\sigma$  and  $MAE$  is bpm.

as a more challenging dataset to cater more realistic conditions than MAHNOB-HCI by [21]. Thus, we have conducted our experiments on COHFACE dataset as well. But it lacks significant motion variations and thus, we create and conduct experiments on our dataset for better evaluation of our proposed method, *MOMBAT*.

The COHFACE dataset contains 160 face videos acquired from 40 subjects. Experiments are conducted on this dataset using the performance metrics, parameter selection and methods described in Section 4.2, Section 4.3 and Section 4.4 respectively. The corresponding results are shown in Table 3. It can be observed that these results are similar to the results on our dataset, that is, *MOMBAT* performs best amongst the

considered methods. Furthermore, it can be observed that the efficacy of *MOMBAT* reduces slightly when COHFACE dataset is considered rather than our dataset. It is because the COHFACE dataset contains compressed videos which deteriorate the HR analysis [36].

## 5. Conclusions

This paper has proposed an HR monitoring method, *MOMBAT*, that is, MONitoring using Modeling and BAYesian Tracking. It has utilized the face videos acquired from a low cost camera in contact-less manner, for HR monitoring. HR monitoring using face videos can be error-prone due to facial expressions, out-of-plane movements, camera parameters and environmental factors. Thus, our *MOMBAT* have alleviated these issues to improve the HR monitoring by introducing pulse modeling and Bayesian HR tracking. The proposed Fourier basis based modeling mitigates the out-of-plane movements by successfully reconstructing the poor quality pulse signal estimates using the good quality pulse signal estimates. The noise can result in some spurious HR estimates, but our proposed Bayesian decision theory based HR tracking mitigates such cases to improve the HR monitoring.

Experimental results have demonstrated that HR monitoring can be significantly improved when both pulse modeling and HR tracking are incorporated. Further, they have indicated that our *MOMBAT* perform the monitoring in near real-time, with an average absolute error of 1.3293 bpm and the Pearson correlation of 0.9746 between predicted and actual HR. This indicates that our method, *MOMBAT* can be effectively used for HR monitoring.

Our method *MOMBAT* can perform spuriously when the face video contains motion that persists for long duration. Our future work will investigate the possibilities to handle this issue by fusing it with Lagrangian techniques. Our method requires some parameter selection. Amongst them, the number of basis depends on the sampling rate. We will be collecting a larger database at different sampling rates using different video compression techniques. It will be used to explore the efficacy of convolutional neural networks [35] and mitigate the compression artifacts for better HR analysis.

## Acknowledgement

Our method was tested on a publicly available dataset COHFACE Dataset in Section 4.6. It was provided by the Idiap Research Institute, Martigny, Switzerland.

## References

- [1] Guha Balakrishnan, Fredo Durand, and John Gutttag. Detecting pulse from head motions in video. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3430–3437, 2013.
- [2] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. 3d constrained local model for rigid and non-rigid facial tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2610–2617. IEEE, 2012.
- [3] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. Constrained local neural fields for robust facial landmark detection in the wild. In *IEEE International Conference on Computer Vision Workshops (ICCV-W)*, pages 354–361, 2013.
- [4] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. Openface: an open source facial behavior analysis toolkit. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–10. IEEE, 2016.
- [5] Gary G Berntson, J Thomas Bigger, Dwain L Eckberg, Paul Grossman, Peter G Kaufmann, Marek Malik, Haikady N Nagaraja, Stephen W Porges, J Philip Saul, Peter H Stone, et al. Heart rate variability: origins, methods, and interpretive caveats. *Psychophysiology*, 34(6):623–648, 1997.
- [6] Brojeshwar Bhowmick, Apurba Mallik, and Arindam Saha. Mobiscan3d: A low cost framework for real time dense 3d reconstruction on mobile devices. In *Intl Conf on Ubiquitous Intelligence and Computing*, pages 783–788, 2014.
- [7] Guido Borghi, Matteo Fabbri, Roberto Vezzani, Rita Cucchiara, et al. Face-from-depth for head pose estimation on depth images. *IEEE transactions on pattern analysis and machine intelligence*, 2018.

- [8] A John Camm, Marek Malik, J Thomas Bigger, Günter Breithardt, Sergio Cerutti, Richard J Cohen, Philippe Coumel, Ernest L Fallen, Harold L Kennedy, RE Kleiger, et al. Heart rate variability: standards of measurement, physiological interpretation and clinical use. task force of the european society of cardiology and the north american society of pacing and electrophysiology. 1996.
- [9] Kingshuk Chakravarty, Suraj Suman, Brojeshwar Bhowmick, Aniruddha Sinha, and Abhijit Das. Quantification of balance in single limb stance using kinect. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 854–858, 2016.
- [10] Xun Chen, Juan Cheng, Rencheng Song, Yu Liu, Rabab Ward, and Z Jane Wang. Video-based heart rate measurement: Recent advances and future prospects. *IEEE Transactions on Instrumentation and Measurement*, 68(10):3600–3615, 2019.
- [11] Gerard De Haan and Vincent Jeanne. Robust pulse rate from chrominance-based rppg. *IEEE Transactions on Biomedical Engineering*, 60(10):2878–2886, 2013.
- [12] Lawrence T DeCarlo. On the meaning and use of kurtosis. *Psychological methods*, 2(3):292, 1997.
- [13] Richard O Duda, Peter E Hart, and David G Stork. *Pattern classification*. John Wiley & Sons, 2012.
- [14] Puneet Gupta, Brojeshwar Bhowmick, and Arpan Pal. Accurate heart-rate estimation from face videos using quality-based fusion. In *IEEE International Conference on Image Processing, (ICIP)*, pages 4132–4136. IEEE, 2017.
- [15] Puneet Gupta, Brojeshwar Bhowmick, and Arpan Pal. Serial fusion of eulerian and lagrangian approaches for accurate heart-rate estimation using face videos. In *IEEE International Conference of the Engineering in Medicine and Biology Society (EMBC)*, pages 2834–2837. IEEE, 2017.
- [16] Puneet Gupta, Brojeshwar Bhowmick, and Arpan Pal. Exploring the feasibility of face video based instantaneous heart-rate for micro-expression spotting. In *IEEE*

- Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1316–1323, 2018.
- [17] Puneet Gupta, Brojeshwar Bhowmik, and Arpan Pal. Robust adaptive heart-rate monitoring using face videos. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 530–538. IEEE, 2018.
- [18] Puneet Gupta and Phalguni Gupta. An accurate finger vein based verification system. *Digital Signal Processing*, 38:43–52, 2015.
- [19] Puneet Gupta and Phalguni Gupta. Fingerprint orientation modeling using symmetric filters. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 663–669. IEEE, 2015.
- [20] Puneet Gupta and Phalguni Gupta. An accurate fingerprint orientation modeling algorithm. *Applied Mathematical Modelling*, 40(15):7182–7194, 2016.
- [21] Guillaume Heusch, André Anjos, and Sébastien Marcel. A reproducible study on remote heart rate measurement. *arXiv preprint arXiv:1709.00962*, 2017.
- [22] Chong Huang, Xin Yang, and Kwang-Ting Tim Cheng. Accurate and efficient pulse measurement from facial videos on smartphones. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–8. IEEE, 2016.
- [23] Ming-Chun Huang, Jason J Liu, Wenyao Xu, Changzhan Gu, Changzhi Li, and Majid Sarrafzadeh. A self-calibrating radar sensor system for measuring vital signs. *IEEE transactions on biomedical circuits and systems*, 10(2):352–363, 2016.
- [24] Peter J Huber. Projection pursuit. *The annals of Statistics*, pages 435–475, 1985.
- [25] Antony Lam and Yoshinori Kuno. Robust heart rate measurement from video using select random patches. In *International Conference on Computer Vision (ICCV)*, pages 3640–3648, 2015.
- [26] Dongseok Lee, Jeehoon Kim, Sungjun Kwon, and Kwangsuk Park. Heart rate estimation from facial photoplethysmography during dynamic illuminance changes.

- In *IEEE International Conference of the Engineering in Medicine and Biology Society (EMBC)*, pages 2758–2761. IEEE, 2015.
- [27] Xiaobai Li, Jie Chen, Guoying Zhao, and Matti Pietikainen. Remote heart rate measurement from face videos under realistic situations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4264–4271, 2014.
- [28] Apurba Mallik, Brojeshwar Bhowmick, and Shah Nawaz Alam. A multi-sensor information fusion approach for efficient 3d reconstruction in smart phone. In *International Conference on Image Processing, Computer Vision, and Pattern Recognition*, 2015.
- [29] A Jonathan McLeod, Dante PI Capaldi, John SH Baxter, Grace Parraga, Xiong-biao Luo, and Terry M Peters. Analysis of periodicity in video sequences through dynamic linear modeling. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 386–393. Springer, 2017.
- [30] Mariana Nogueira, Mathieu De Craene, Sergio Sanchez-Martinez, Devyani Chowdhury, Bart Bijmens, and Gemma Piella. Analysis of nonstandardized stress echocardiography sequences using multiview dimensionality reduction. *Medical Image Analysis*, 60:101594, 2020.
- [31] Michel Owayjan, Ahmad Kashour, Nancy Al Haddad, Mohamad Fadel, and Ghinwa Al Souki. The design and development of a lie detection system using facial micro-expressions. In *International Conference on Advances in Computational Tools for Engineering Applications (ACTEA)*, pages 33–38. IEEE, 2012.
- [32] Constantinos B Papadias. Globally convergent blind source separation based on a multiuser kurtosis maximization criterion. *IEEE Transactions on Signal Processing*, 48(12):3508–3519, 2000.
- [33] Son Lam Phung, Abdesselam Bouzerdoum, and Douglas Chai. A novel skin color model in ycbcr color space and its application to human face detection. In

- International Conference on Image Processing (ICIP)*, volume 1, pages I–289. IEEE, 2002.
- [34] 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, 2011.
- [35] Ying Qiu, Yang Liu, Juan Arteaga-Falconi, Haiwei Dong, and Abdulmotaleb El Saddik. EVM-CNN: Real-time contactless heart rate estimation from facial video. *IEEE Transactions on Multimedia*, 21(7):1778–1787, 2019.
- [36] Michal Rapczynski, Philipp Werner, and Ayoub Al-Hamadi. Effects of video encoding on camera based heart rate estimation. *IEEE Transactions on Biomedical Engineering*, 66(12):3360–3370, 2019.
- [37] Michal Rapczynski, Philipp Werner, Frerk Saxen, and Ayoub Al-Hamadi. How the region of interest impacts contact free heart rate estimation algorithms. In *International Conference on Image Processing (ICIP)*, pages 2027–2031. IEEE, 2018.
- [38] Juergen Reichert. Automatic classification of communication signals using higher order statistics. In *International Conference on Acoustics, Speech, and Signal Processing*, volume 5, pages 221–224. IEEE, 1992.
- [39] Angel Melchor Rodríguez and J Ramos-Castro. Video pulse rate variability analysis in stationary and motion conditions. *Biomedical engineering online*, 17(1):11, 2018.
- [40] Arun A Ross, Karthik Nandakumar, and Anil Jain. *Handbook of multibiometrics*, volume 6. Springer Science & Business Media, 2006.
- [41] Lizawati Salahuddin, Jaegeol Cho, Myeong Gi Jeong, and Desok Kim. Ultra short term analysis of heart rate variability for monitoring mental stress in mobile settings. In *International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 4656–4659. IEEE, 2007.

- [42] Jason M Saragih, Simon Lucey, and Jeffrey F Cohn. Deformable model fitting by regularized landmark mean-shift. *International Journal of Computer Vision*, 91(2):200–215, 2011.
- [43] Sanjana Sinha, Brojeshwar Bhowmick, Kingshuk Chakravarty, Aniruddha Sinha, and Abhijit Das. Accurate upper body rehabilitation system using kinect. In *International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 4605–4609, 2016.
- [44] Mohammad Soleymani, Jeroen Lichtenauer, Thierry Pun, and Maja Pantic. A multimodal database for affect recognition and implicit tagging. *IEEE Transactions on Affective Computing*, 3(1):42–55, 2012.
- [45] Emily J Lam Po Tang, Amir HajiRassouliha, Martyn P Nash, Andrew J Taberner, Poul MF Nielsen, and Yusuf O Cakmak. Removing drift from carotid arterial pulse waveforms: A comparison of motion correction and high-pass filtering. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 111–119. Springer, 2019.
- [46] Mika P Tarvainen, Perttu O Ranta-Aho, Pasi A Karjalainen, et al. An advanced detrending method with application to HRV analysis. *IEEE Transactions on Biomedical Engineering*, 49(2):172–175, 2002.
- [47] Sergey Tulyakov, Xavier Alameda-Pineda, Elisa Ricci, Lijun Yin, Jeffrey F Cohn, and Nicu Sebe. Self-adaptive matrix completion for heart rate estimation from face videos under realistic conditions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2396–2404, 2016.
- [48] Wim Verkrusse, Lars O Svaasand, and J Stuart Nelson. Remote plethysmographic imaging using ambient light. *Optics express*, 16(26):21434–21445, 2008.
- [49] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 511–518. IEEE, 2001.

- [50] Nannan Wang, Xinbo Gao, Dacheng Tao, and Xuelong Li. Facial feature point detection: A comprehensive survey. *arXiv preprint arXiv:1410.1037*, 2014.
- [51] Peng Wang, Matthew B Green, Qiang Ji, and James Wayman. Automatic eye detection and its validation. In *IEEE Conference on Computer Vision and Pattern Recognition-Workshops (CVPRW)*, pages 164–164. IEEE, 2005.
- [52] 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, 2016.
- [53] Hao-Yu Wu, Michael Rubinstein, Eugene Shih, John Guttag, Frédo Durand, and William T. Freeman. Eulerian video magnification for revealing subtle changes in the world. *ACM Transactions on Graphics*, 31(4), 2012.
- [54] Chenggang Yu, Zhenqiu Liu, Thomas McKenna, Andrew T Reisner, and Jaques Reifman. A method for automatic identification of reliable heart rates calculated from ECG and PPG waveforms. *Journal of the American Medical Informatics Association*, 13(3):309–320, 2006.
- [55] Zhilin Zhang, Zhouyue Pi, and Benyuan Liu. TROIKA: A general framework for heart rate monitoring using wrist-type photoplethysmographic signals during intensive physical exercise. *IEEE Transactions on biomedical engineering*, 62(2):522–531, 2015.