

Deep adaptative spectral zoom for improved remote heart rate estimation

Joaquim Comas¹, Adrià Ruiz², Federico Sukno¹

¹ Department of Information and Communication Technologies, Pompeu Fabra University, Barcelona, Spain

² Seedtag, Madrid, Spain

Abstract—Recent advances in remote heart rate measurement, motivated by data-driven approaches, have notably enhanced accuracy. However, these improvements primarily focus on recovering the rPPG signal, overlooking the implicit challenges of estimating the heart rate (HR) from the derived signal. While many methods employ the Fast Fourier Transform (FFT) for HR estimation, the performance of the FFT is inherently affected by a limited frequency resolution. In contrast, the Chirp-Z Transform (CZT), a generalization form of FFT, can refine the spectrum to the narrow-band range of interest for heart rate, providing improved frequential resolution and, consequently, more accurate estimation. This paper presents the advantages of employing the CZT for remote HR estimation and introduces a novel data-driven adaptive CZT estimator. The objective of our proposed model is to tailor the CZT to match the characteristics of each specific dataset sensor, facilitating a more optimal and accurate estimation of HR from the rPPG signal without compromising generalization across diverse datasets. This is achieved through a Sparse Matrix Optimization (SMO). We validate the effectiveness of our model through exhaustive evaluations on three publicly available datasets —UCLA-rPPG, PURE, and UBFC-rPPG—employing both intra- and cross-database performance metrics. The results reveal outstanding heart rate estimation capabilities, establishing the proposed approach as a robust and versatile estimator for any rPPG method.

I. INTRODUCTION

Recently, the research community has increasingly focused on the camera-based measurement of human physiological signals and their potential applications [43], [2], [24], particularly in the extraction and analysis of vital signs such as the heart rate (HR), heart rate variability (HRV), respiration rate (RR), oxygen saturation (SpO₂), and blood volume pulse (BVP). Among these vital signs, HR has been the most extensively studied due to its relevance to health and human-computer interaction applications, e.g. by using HR to complement facial expressions for the analysis of human emotions [17], [8]. Indeed, the majority of metrics used to evaluate the performance of remote photoplethysmography (rPPG) methods are HR-based, including mean absolute error (MAE), root mean square error (RMSE), and Pearson’s correlation (R), all derived from the rPPG signal.

Conventional rPPG methods involve a two-stage process: initial signal extraction based on rPPG principles and subsequent HR calculation using signal processing. While recent progress has been made in the first stage [7], [32], [58], the second stage relies heavily on traditional hand-crafted methods and remains largely unexplored. Three main approaches to derive HR from rPPG are distinguished: in the temporal

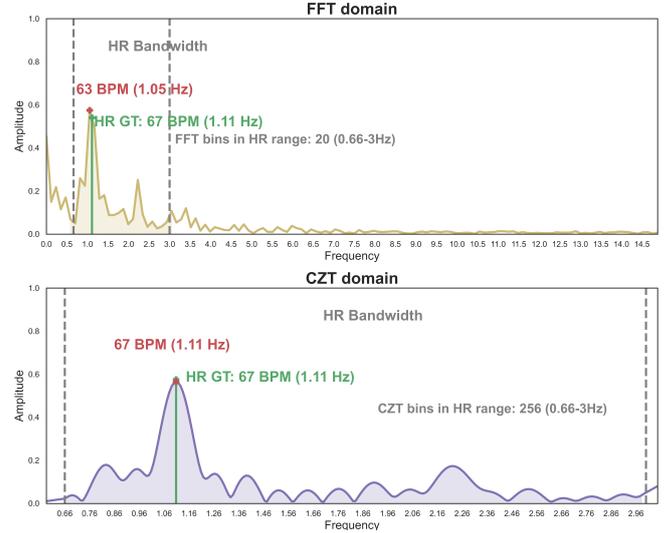


Fig. 1: Frequency spectrum comparison between FFT (above) and CZT (below) for a UCLA-rPPG PPG signal sample (8.5 sec windows). In red is denoted the predicted HR while green represents the HR ground truth.

domain, peak detection algorithms like Pan-Tompkins [34] are commonly used, providing acceptable HR estimates, especially for ECG. However, rPPG signals can often be very noisy, requiring post-processing and manual fine-tuning, making this method challenging. In the frequency domain, the FFT is widely used for HR estimation. Despite its precision, its range is fixed by the Nyquist frequency, which is well above the interest range for HR, under a uniform spectrum sampling. Thus, the frequency resolution of the FFT is also fixed and often sub-optimal for HR estimation, e.g. very few frequency bins fall within the range of interest. Recent data-driven approaches propose end-to-end models, but they face challenges like convergence issues due to realistic rPPG-irrelevant noise [58] or overlook the ordinal nature of HR outputs, as seen in other deep learning problems like facial age estimation [6].

In this paper, we propose the use of the Chirp-Z Transform (CZT) [41] for the task of remote HR estimation. Originating in 1969, this method serves as a generalization of the widely used Discrete Fourier Transform (DFT). With its unique capabilities, it allows for a flexible choice of the spectral resolution, rendering it particularly well-suited for HR estimation—a superiority that we will demonstrate in section IV over traditional FFT estimation. Indeed, the CZT can change

both the resolution and the range of frequencies to be covered, thus allowing to focus on a particular range of interest. This is illustrated in Figure 1. This feature also enables more precise estimations using short temporal windows of the rPPG signal, making it particularly suitable for continuous heart rate (HR) evaluation, as considered in the recent Vision-for-Vitals challenge (V4V) [42]. Most existing video-based physiological recovery methods have traditionally focused on predicting HR values over large window intervals (e.g., 30 seconds or video-level). While this evaluation protocol may provide unbiased performance for these methods, it limits their effectiveness in real-time applications, such as computing Heart Rate Variability (HRV) or diagnosing Atrial fibrillation [36], where the CZT could offer a viable solution for these real-time HR applications.

Finally, inspired by previous works [49], [20] formulating the FFT computation as a neural network, we introduce a novel deep CZT estimator. In particular, we formulate the CZT as a set of fully differentiable network layers, allowing to fine-tune its parameters and thus, capturing better the relationship between the rPPG signal and the HR measurement sensor of the considered dataset. Although the goal is to improve the estimation by learning the characteristics of each HR sensor, we also introduce a Sparse Matrix Optimization (SMO) regularization loss to retain the generalization ability and structure of the standard CZT, balancing the deviations required to adapt to the considered sensor without overfitting. Interestingly, we find that the performance improvement of the proposed CZT is reproducible for different rPPG signal extraction methods, as we will empirically show in section IV-E. This allows it to be integrated within any type of rPPG method, be it handcrafted or data-driven, and accommodates different temporal window sizes, enabling continuous HR estimation using the frequency domain.

A. Contribution

In this work, we showcase the advantages of incorporating the Chirp-Z Transform for remote HR estimation. Our contribution extends to the introduction of a novel deep-learning CZT estimator module, specially designed to enhance HR estimation in rPPG conditions. The core objective of our proposed HR estimator is to acquire a customized variation of the classical CZT, adapting to the distinctive characteristics of each HR dataset sensor. To train our model effectively, we also introduce a novel Sparse Matrix Optimization (SMO) loss to not only facilitate robust generalization within its original dataset but also promote adaptability across diverse datasets by constraining the adaptation of the data-driven CZT. Another notable feature of the presented CZT estimator module is its flexibility to be integrated into any rPPG approach, providing a more precise HR estimator regardless of the adopted rPPG method. To validate our approach, comprehensive experiments are conducted using 3 publicly available datasets: UCLA-rPPG, PURE and UBFC-rPPG, and including both intra- and cross-dataset evaluation.

The remainder of this paper is organized as follows: firstly, in Section II, we conduct a comprehensive review of the

rPPG approaches and the different methods used to estimate the HR in the existing literature. The proposed approach is presented in Section III, while experimental results are provided in Section IV. Section V summarizes our findings and conclusions.

II. RELATED WORK

A. rPPG measurement

Since Takano et al. [46] and Verkruysse et al. [50] demonstrated the feasibility of remote HR measurement from facial videos, researchers have proposed diverse methods for physiological data recovery. Some focus on regions of interest, utilizing techniques such as Blind Source Separation [40], [39], [21], Normalized Least Mean Squares [23], or self-adaptive matrix completion [48]. Others leverage the skin optical reflection model by projecting RGB skin pixel channels into an optimized subspace, mitigating motion artifacts [10], [52].

In recent years, deep learning-based methods [45], [57], [37], [19], [28], [33] have surpassed conventional techniques, achieving state-of-the-art performance in estimating vital signs from facial videos. Some combine knowledge from traditional methods with Convolutional Neural Networks (CNNs) to exploit sophisticated features [31], [32], [44]. Recent works like [19], [26] explore unsupervised approaches using meta-learning, demonstrating improved generalization in out-of-distribution cases. On the other hand, some researchers focus on fully end-to-end approaches [7], [56], [37]. Unlike previous methods, end-to-end models use facial videos as input to predict the rPPG signal directly. Unsupervised techniques have gained traction in end-to-end approaches, with transformer-based models like Physformer [58] and RADIANT [13] leveraging long-range spatiotemporal features. However, these models currently lack optimization for computational efficiency, rendering them unsuitable for deployment on mobile devices. Additionally, they often depend on a prior fine-tuning stage due to limited generalization capabilities, especially in computer vision tasks [11]. While promising, they may not yet demonstrate a significant performance advantage over CNN-based models [25]. Finally, works like [25] and [9] propose lightweight rPPG frameworks with competitive HR estimation results while controlling computational cost.

B. HR estimation from rPPG

Once the rPPG signal is successfully extracted, a post-processing stage is employed to derive the corresponding HR value. Traditional literature distinguishes between temporal and frequency analysis. In the temporal domain, the peak detection method calculates the inter-beat intervals (IBI) between consecutive beats. Subsequently, HR is determined by averaging all IBIs over a time window and computing the inverse. Approaches [39], [1], [29], [52] have heavily relied on these peak detector algorithms for HR estimation. Temporal HR analysis remains prevalent in state-of-the-art methods, especially in real-time applications like [14], [12], or when intrinsic characteristics like HRV require analysis,

as demonstrated in [30], [57]. However, these methods depend on accurate peak detection for robust instantaneous HR estimation, leading researchers to use customized peak detection functions for each presented database, making standardization challenging.

On the other hand, in frequency-based methods, the conventional approach involves converting the extracted rPPG signal to the frequency domain using the FFT, typically implementing Welch’s method for density estimation. This technique assumes HR periodicity, where the highest power in the spectrum resides within the HR bandwidth. To mitigate the noise present in the rPPG signal, bandpass filtering in the HR range is commonly applied. This method, ranging from early rPPG handcrafted approaches [40], [21], [48] to recent deep learning methods [7], [37], [13], is widely adopted. However, a significant limitation of the FFT lies in assuming signal periodicity, restricting its frequency resolution. Consequently, the FFT may struggle to accurately discern closely spaced frequency components, especially in the narrowed HR band. Some attempts to address the weaknesses of the FFT were introduced by Irani et al. [16] and Yu et al. [55]. The first one, suggested replacing the FFT with the Discrete Cosine Transform (DCT), demonstrating improved accuracy compared to FFT, while the second employed the Short-Fast-Fourier Transform (SFFT), proving more suitable for rapidly changing heart rate trends. However, these solutions only partially mitigate FFT issues. The DCT, being a derivative of the FFT, still suffers from the same FFT issues, while the SFFT involves a trade-off between temporal and frequency resolutions. Instead, in this paper, we propose a solution using the CZT, a generalization of the FFT that offers flexibility in specifying bandwidth without sacrificing good frequency and time resolution simultaneously.

Despite many researchers still relying on traditional signal-processing techniques, the emergence of deep learning methods has led to various data-driven solutions. Among them, several researchers proposed the usage of deep HR regressors. Spetlik et al. [45] pioneered this approach by combining a two-step convolutional neural network, with the second network being a deep HR estimator to extract the HR value. Niu et al. [32] introduced a deep regression model through a Gated Recurrent Unit (GRU), while works like [54], [5], [22] incorporated fully connected layers into their end-to-end frameworks for direct HR extraction without prior rPPG signal recovery. Despite architectural differences, these methods share similar HR optimization using L1 loss or categorical cross-entropy loss. Even so, these optimization functions often force the model to learn periodic features within target frequency bands, which is challenging due to the present noise in the rPPG signal. Furthermore, the usage of categorical cross-entropy and the treatment of the HR band as a multiclass classification problem does not take into account the inter-class relationships between the different classes. To consider the ordinal structure of the HR values, we adopt a Squared Earth Mover’s Distance Loss [15] in the spectral density domain to optimize our deep CZT adaptive estimator. This loss aims to adapt the parameters

of our deep CZT model to fit the characteristics of the HR sensor, which, jointly with the SMO loss, also guarantees a constrained adaptation, promoting generalization in cross-dataset evaluation.

III. METHODOLOGY

In this section, we introduce and define the conventional CZT. Subsequently, we present our proposed deep CZT adaptive estimator and finally present our optimized objective function for remote HR estimation.

A. Chirp-Z Transform

The CZT [41] computes the z-transform of the finite duration signal $x[n]$ along a general spiral contour in the z-plane. Therefore, the CZT is defined using the following formula:

$$CZT(x[n]) = \sum_{n=0}^{N-1} x[n] \cdot z_k^{-n} \quad (1)$$

Unlike the DFT, which evaluates the Z-transform of $x[n]$ on N equally spaced points on the unit circle in the z-plane, the CZT is not constrained to operate along the unit circle, evaluating the z-transform along spiral contours described as:

$$z_k = A \cdot W^{-k}, \quad k = 0, 1, \dots, M-1 \quad (2)$$

where A is the complex starting point, W is a complex scalar describing the complex ratio between points on the contour, and M is the length of the transform. In addition, the CZT can also be expressed with the following matrix expression:

$$\underbrace{\begin{bmatrix} X_0 \\ X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}}_x = \underbrace{\begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & W^1 & W^2 & \dots & W^{(n-1)} \\ 1 & W^2 & W^4 & \dots & W^{2(n-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & W^{(m-1)} & W^{2(m-1)} & \dots & W^{(m-1)(n-1)} \end{bmatrix}}_W \underbrace{\begin{bmatrix} A^{-0} & 0 & 0 & \dots & 0 \\ 0 & A^{-1} & 0 & \dots & 0 \\ 0 & 0 & A^{-2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & A^{-n} \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}}_x \quad (3)$$

where A is a diagonal matrix N by N and W is an M by N Vandermonde matrix.

During this work, we employ this matrix operation to compute the CZT, which is relevant for the design of our deep CZT adaptive estimator, described in subsection III-B. Although the CZT can be implemented in an efficient form such as FFT, by using the Bluestein algorithm [3], we do not aim to find an efficient computation of CZT in this work. In this paper, we are interested in the CZT as an alternative form to estimate HR from rPPG signals by taking advantage of its spectral zoom property. Therefore, we propose to limit the CZT to the region of the unit circle because we are only interested in a narrow bandwidth of the rPPG frequency spectrum. So, we can define a zoom region that begins at A and ends at $(M-1) \cdot W$. Following the literature and considering the majority of rPPG databases, we decided to limit the HR band within 0.66 and 3 Hz, equivalent to 40 to 180 beats-per-minute (BPM). Apart from defining the zoom spectrum region, the used bin density in the CZT is also configurable. After our preliminary experiments, we set the size of the CZT, M , equal to the size of the input size N ,

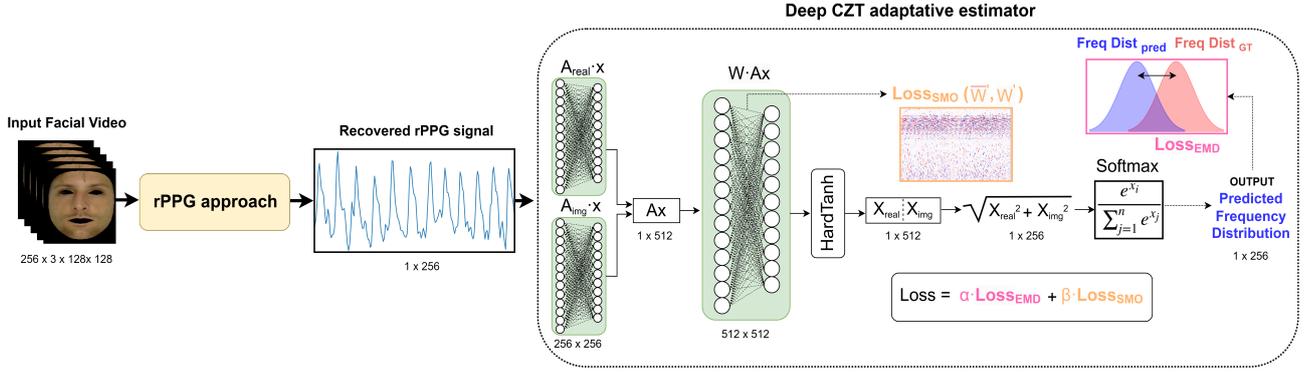


Fig. 2: Overall structure of our proposed deep CZT adaptive estimator model.

as shown in the example of Figure 1. This means, at the typical sampling rate of 30 frames per second, CZT has approximately 13 times higher frequency resolution than FFT in HR bandwidth, i.e. the same number of bins is used to cover 2.33 Hz instead of 30 Hz.

B. Deep CZT estimator

Our presented model, depicted in Figure 2, is designed to adapt the conventional CZT for estimating HR values from the rPPG signal. This adaptation involves making constrained modifications to the general structure of the CZT to learn the relationship between the input signal and the characteristics of the HR sensor. To achieve this, we represent the CZT framework as a neural network. However, addressing the complex nature of the transformation and ensuring proper optimization is challenging due to difficulties in the training dynamics between real and imaginary parts of the network, particularly at the backpropagation stage [47]. We overcome this issue by relying on the structure of the CZT to avoid the direct use of the imaginary component. We duplicate the size of our A and W matrices, with the left half containing the real values and the right half containing the imaginary values, following Euler's formula:

$$e^{\pm i\theta} = \cos \theta \pm i \sin \theta \quad (4)$$

Applying this formulation in equation 3 we can express our deep CZT as the following matrix multiplication operation:

$$\begin{bmatrix} X_{\text{Re}} \\ X_{\text{Im}} \end{bmatrix} = \begin{bmatrix} W_{\text{Re}} & -W_{\text{Im}} \\ W_{\text{Im}} & W_{\text{Re}} \end{bmatrix} \begin{bmatrix} Ax_{\text{Re}} \\ Ax_{\text{Im}} \end{bmatrix} \quad (5)$$

Initially, the input rPPG signal x is fed to two fully connected layers, which contain the weights of the A matrix. One matrix for the real part (the $\cos(\cdot)$ term), and another for the imaginary part (the $\sin(\cdot)$ term). As A is a diagonal matrix, and its diagonal denotes the starting frequency point, each fully connected layer contains a single parameter. Notably, this parameter is kept fixed, emphasizing optimization within a predetermined bandwidth between 0.66 and 3 Hz.

Subsequently, adjustments are made to the W matrix to accommodate complex number matrix multiplication between W and Ax :

$$(W_{\text{Re}} + iW_{\text{Im}})(Ax_{\text{Re}} + iAx_{\text{Im}}) = (W_{\text{Re}} \cdot Ax_{\text{Re}} - (W_{\text{Im}} \cdot Ax_{\text{Im}}) + i((W_{\text{Re}} \cdot Ax_{\text{Im}}) + (W_{\text{Im}} \cdot Ax_{\text{Re}}))) \quad (6)$$

To optimize the W matrix efficiently while preserving the CZT structure, we form the last fully connected layer as a square matrix initialized with the values of W_{Re} and W_{Im} from the classical CZT. To reduce the number of learnable parameters by half for W , we learn W_{Re} and W_{Im} once and then replicate them on the respective sides of the square W matrix in Equation 5. As W contains \cos and \sin values, and its weights fall within the range of -1 to 1, we constrain the learning of new weights to this range using a HardTanh function. Finally, we combine the real and imaginary components of the learned CZT by computing the modulus to obtain the learned frequential distribution. This distribution is then normalized using a Softmax function, and, similar to Welch's method, the heart rate is computed by extracting the frequency with the maximal power response and multiplying it by 60.

C. Loss function

To optimize the weights of the W matrix in our deep CZT adaptive estimator, we employ a combined loss function. The aim of this combined loss function is two-fold: 1) Learn the closest frequency distribution from the rPPG signal to the distribution of the HR sensor used as ground truth. 2) Constrain the learning process to avoid capturing irrelevant noise features from specific databases, guiding the model to focus on meaningful adaptations of the classic CZT for remote HR estimation aligned with sensor-specific characteristics.

For the first objective, we expect the HR estimator to predict HR distributions where sub-bands closer to the HR ground truth distribution are assigned higher probabilities than those further away. Unlike Categorical Cross-Entropy or MAE loss, which do not consider the inter-class relationship in HR distributions, we propose an alternative approach using Earth Mover's Distance (EMD). EMD is defined as the minimum cost to transport the mass of one distribution onto another one. To implement this, we adopt the Squared Earth

Mover’s Distance-based Loss proposed in [15], expressed as:

$$\mathcal{L}_{EMD} = \frac{1}{N} \sum_{i=1}^N (CDF_i(p) - CDF_i(t))^2 \quad (7)$$

where $CDF()$ denotes the cumulative density function, p and t denote two compared distributions of same size N , which represents the batch dimensionality. In our framework, $CDF(p)$ and $CDF(t)$ represent the frequency density functions of the predicted HR and the ground truth HR, as illustrated in Figure 2.

For the second objective, we seek a trade-off between the model customization of the HR sensor and the preservation of the classical CZT. Although the goal of our deep CZT adaptive estimator is to enhance the accuracy of HR estimation by learning the intrinsic features of the HR sensor, we aim to retain the structure of the classical CZT, which already demonstrates outstanding HR results, as will be shown in Section IV. Therefore, to control the adaptation of the most relevant features, our \mathcal{L}_{SMO} acts as a regularization loss of the \mathcal{L}_{EMD} by constraining the learning of W . Referring to Subsection III-B, we focus exclusively on the W learnable weights. Therefore, we construct a \tilde{W} matrix comprising the W_{Re} and W_{Im} matrices without their replications needed for the full W matrix:

$$\tilde{W} = \underbrace{\left[\begin{array}{c|c} W_{Re} & W_{Im} \end{array} \right]}_N \Bigg\} M \quad (8)$$

Here, N represents the number of columns, which is twice the input signal’s size after separating the real and imaginary components, and M denotes the number of rows, where each row corresponds to a frequency sub-band between 0.66 and 3 Hz. To control the flexibility in learning the optimized W through the \mathcal{L}_{EMD} , we perform HR frequency sub-band optimization. This involves comparing each row value between \tilde{W}' and \tilde{W} , subsequently averaging all the sub-band errors, expressing \mathcal{L}_{SMO} as:

$$\mathcal{L}_{SMO} = \frac{1}{L} \sum_{i=1}^M \sum_{j=1}^N |\tilde{W}'(i, j) - \tilde{W}(i, j)| \quad (9)$$

where \tilde{W}' represents the learnable version of \tilde{W} from our CZT estimator and L denotes the total number of frequency bins. In summary, the overall loss function \mathcal{L}_{HR} can be formulated as:

$$\mathcal{L}_{HR} = \alpha \cdot \mathcal{L}_{EMD} + \beta \cdot \mathcal{L}_{SMO}, \quad (10)$$

where α and β are balancing parameters. In our experiments, we set $\alpha = 100$ and $\beta = 0.01$ empirically based on our preliminary experiments. These settings yield a balanced contribution of both losses, taking into account that the magnitude of \mathcal{L}_{SMO} is much larger than that of \mathcal{L}_{EMD} .

IV. EXPERIMENTS

In this section, we introduce three benchmark datasets and outline the implementation of our methodology. Our analysis begins by comparing HR estimations derived from PPG ground truth and sensor HR data. Specifically, in the initial experiment, we compare the classical CZT with baseline methods for HR estimation, concurrently assessing the influence of signal length. Subsequently, we evaluate remote HR estimation performance, referencing the ground truth data. Our exploration then extends to examining the impact of the loss function in our deep CZT model. We conduct a comprehensive evaluation through both intra-database and cross-database assessments, showcasing its robustness in the context of remote HR estimation. Finally, we test our deep CZT model with various rPPG models, highlighting its versatile capability with different rPPG approaches.

A. Datasets

We assessed our approach on the following rPPG datasets, since each of them contains PPG and HR ground truth data.

The UCLA-rPPG dataset [53] comprises 489 videos from 98 subjects with diverse characteristics, including skin tones, ages, genders, and ethnicities. Each subject underwent five trials, with each trial lasting approximately 1 minute. The recordings were captured at a resolution of 640×480 pixels and 30 frames per second (FPS), in an uncompressed format. Synchronous gold-standard PPG and HR measurements were collected alongside the facial videos. Due to the lack of predefined folds in this dataset, we split the data into training (80%), validation (10%), and testing (10%) sets.

The UBFC-rPPG [4] includes 42 RGB videos from 42 subjects. The subjects were asked to play a time-sensitive mathematical game, emulating a standard human-computer interaction scenario, to obtain varied HR during the experiment. The recorded facial videos were acquired indoors with varying sunlight and indoor illumination at 30 FPS with a webcam (Logitech C920 HD Pro) at a resolution of 640×480 in uncompressed 8-bit RGB format. The bio-signals ground truth was acquired using a CMS50E transmissive pulse oximeter to record the PPG signal and heart rate. In our experiments, we used UBFC-rPPG in a cross-dataset evaluation, where all 42 videos were used only for testing.

The Pulse Rate Detection Dataset (PURE) contains 60 videos from 10 subjects (8 male, 2 female) performing six different head motion tasks: steady, talking, slow translation, fast translation, small rotation, and medium rotation. The facial videos were recorded using an ECO274CVGE camera with a resolution of 640×480 pixels and 30 FPS. Each video is about 1 minute long and stored in uncompressed PNG format. The gold-standard measures of BVP and SpO2 were collected using a finger pulse oximeter. Similar to the UBFC-rPPG dataset, we considered the PURE dataset only for cross-dataset evaluation.

B. Implementation details

1) *Preprocessing and training procedure:* In all our experiments, we adopt the preprocessing stage from [9] for

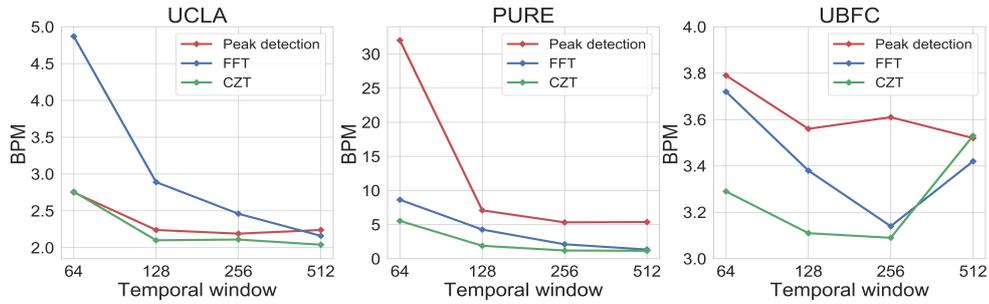


Fig. 3: MAE comparison between handcrafted HR estimation methods using different temporal windows with ground truth data.

each dataset. However, a minor modification is made to facial region estimation compared to [9]. Instead of using the MTCNN [59] algorithm, we choose for facial preprocessing segmentation, utilizing the Mediapipe Face Mesh¹ model [18]. This adjustment simplifies rPPG estimation by focusing on facial skin without introducing additional complexity in terms of parameters or optimization tasks. After masking the facial video, each frame is resized to 128×128 pixels.

Our deep CZT estimator module undergoes a two-stage training process. Initially, we pre-train all rPPG approaches used in subsections IV-D and IV-E, adhering to the specifications of each approach following [27] and [9]. Subsequently, we train our CZT estimator by initializing all weights similar to the traditional CZT, while keeping the weights of the rPPG method frozen. The implementation uses PyTorch 1.8.0 [35] and is trained on a single NVIDIA GTX1080Ti. Sequences of 256 frames without overlap are used during training. We employ the AdamW optimizer with a learning rate of 0.0001 and a ReduceLROnPlateau learning rate scheduler with a patience of 5 and a factor of 0.9. Regarding the experiment reproducibility, the details of our experimental setup, as presented in Section IV, will be made publicly available in a GitHub repository.

2) *Metrics and evaluation*: To evaluate the HR estimation performance of the proposed model, we adopt the same metrics used in the literature, such as the mean absolute HR error (MAE), the root mean squared HR error (RMSE), the mean absolute percentage error (MAPE) and Pearson’s correlation coefficients R [23], [26], [25]. Ablation experiments are performed using 256-frame sequences (8.5 seconds approx) with no overlap to compute HR estimation for all reported metrics, which is more challenging and informative than HR estimation based on the whole video sequence at once. Importantly, our goal is to approximate the HR sensor’s value, and thus, all reported results are compared against the HR sensor ground truth data of each dataset without deriving HR from the PPG signal, as is commonly done in other works more focused on assessing the quality of the rPPG signal than on the estimated HR.

C. Preliminary experiment: Why CZT for remote HR estimation?

In our preliminary experiment, we aim to highlight the advantages of employing the CZT for remote HR estimation. This initial study has two main objectives: 1) to quantify the improvements in accuracy with respect to standard practice, and 2) to show the possibility of shorter-window estimates.

Firstly, we estimate and compare HR values using directly the PPG signal provided by the evaluated datasets (i.e. simulating a perfect rPPG estimation), benchmarking against two widely used techniques for HR estimation: peak detection and FFT. Secondly, we evaluate the performance of each technique across varying temporal windows for each dataset, investigating the impact of temporal window size on HR estimation.

In Figure 3, the MAE between the HR extracted from the PPG ground truth signal and HR ground truth data is depicted for UCLA-rPPG, PURE, and UBFC-rPPG datasets. The behaviour of each HR estimation method is explored across four temporal windows (2.5 to 17 seconds, or 64 to 512 samples), incremented in powers of two, for each dataset. A consistent trend emerges across datasets, showing HR error decreasing with larger window sizes, as shorter temporal windows pose greater challenges. Comparing the three methods, the CZT (in green) achieves the lowest error in the three datasets across most window sizes, maintaining a similar error across different temporal windows. We observe two important deviations from these general observations:

- 1) In the UBFC-rPPG dataset for a window size of 512 samples, the HR error for both frequential methods increases instead of reducing with respect to the error using shorter windows. This is due to the high variability of HR over time for this dataset, as subjects were playing a mathematical game while recorded. In such case, estimating HR based on a unique maximum of the spectrum, as done in FFT and CZT, inevitably leads to inaccuracies.
- 2) In the PURE dataset for a window size of 64 samples, we observe that the temporal method (in red) produces comparatively large errors. This happens because the majority of recordings in the PURE dataset are for subjects with HR between 40 and 60 BPM; therefore there are very few peaks in such a short temporal window and the resulting estimates are poor.

¹https://github.com/google/mediapipe/blob/master/docs/solutions/face_mesh.md

TABLE I: Comparison of different loss functions (beats per minute).

Loss	Intra-database evaluation			
	MAE↓	RMSE↓	MAPE↓	R↑
\mathcal{L}_{CE}	2.22±0.28	2.98±2.18	2.91±0.35	0.97±0.03
\mathcal{L}_{EMD}	1.88±0.32	2.94±2.90	2.45±0.38	0.96±0.04
$\mathcal{L}_{EMD} + \mathcal{L}_{SMO}$	1.74±0.26	2.53±1.67	2.26±0.32	0.97±0.03

A notable observation from this preliminary experiment is that none of the employed methods yields an absolute error of zero in any dataset. While in the ideal scenario, we anticipate a perfect match between the PPG signal and the HR values acquired from the HR sensor, this initial experiment reveals an error gap between both ground truth signals. Even with a highly precise rPPG approach, an error persists in the recovered signal and the HR values [12]. Determining the cause of this difference is challenging, given that most oximeter devices do not provide information on how HR is estimated from the PPG signal, and this process may vary across devices. This underscores the importance of establishing a mapping between the recovered PPG signal and the HR value, which is the objective of our deep CZT estimator, studied in the next section.

D. Ablation studies

In this subsection, we will evaluate our proposed deep CZT estimator module. For our ablation studies we choose as rPPG approach the TDM model [9], because it consists of a lightweight method that obtains competitive rPPG results controlling the amount of parameters.

1) *Loss function impact:* In Table I, we perform a loss function experiment for our deep CZT estimator, testing different losses in intra-database evaluation. As detailed in subsection III-B, our proposed model outputs the estimated frequency distribution. Therefore, we explore three distinct losses: cross-entropy, squared Earth Mover’s Distance-based loss, and our proposed loss, which combines the Squared Earth Mover’s Distance-based Loss with SMO regularization loss.

Comparing the cross entropy and the squared EMD loss, we can appreciate a significant difference, especially in MAE and MAPE metrics where the error is considerably lower. As explained in section II-B, this discrepancy can be attributed to the cross-entropy loss not considering inter-class relationships among different HR classes, whereas the EMD loss incorporates the ordinal behavior of the HR frequency distribution. Finally, we find the performance of our combined loss, $\mathcal{L}_{EMD} + \mathcal{L}_{SMO}$, which denotes the best HR results for all the metrics surpassing the previous EMD and cross-entropy losses. This superior performance can be attributed to the regularization imposed over the parameters of the matrix W , compelling our estimator to modify only relevant information. This preserves the traditional CZT structure, preventing the learning of inappropriate modifications or overfitting.

In Figure 4, the SMO results are illustrated, depicting the difference map between the initial \tilde{W} matrix of the standard CZT and our learned \tilde{W}' matrix after training

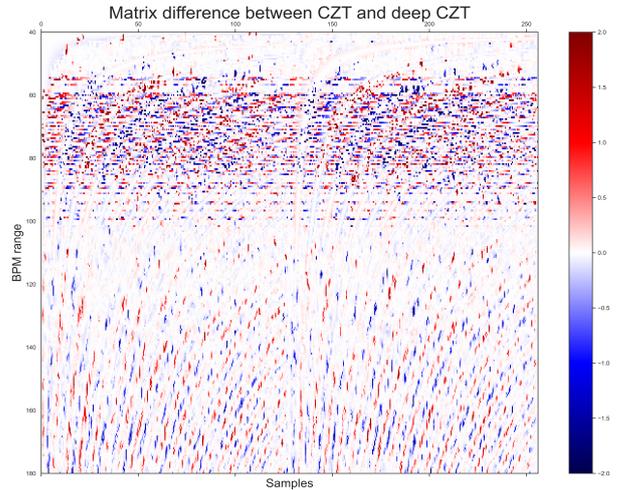


Fig. 4: Sparse Matrix optimization result. Difference between \tilde{W} matrix of standard CZT and \tilde{W}' matrix of trained deep CZT estimator.

our CZT estimator. The majority of changes, fluctuating between -2 and $+2$, mainly occur in the frequency bands between 50 and 100 BPM (on the y-axis), aligning with the central distribution of the UCLA-rPPG dataset around 60-70 BPM, which is the training dataset. Despite observing more changes across the entire matrix, the color is more attenuated, indicating that the difference between the original and learned weights is nearly negligible.

2) *Cross-database evaluation:* In Table II, we evaluate our deep CZT estimator’s performance with the proposed combined loss in cross-dataset scenarios, specifically in the PURE and UBFC-rPPG datasets. The comparison includes the FFT and the standard CZT for a comprehensive analysis.

In the PURE dataset, we can see a notable difference between the traditional FFT and the CZT, where the MAE and RMSE errors are 1 BPM and 2.5 BPM lower, respectively. Regarding deep CZT, we obtain competitive results even in cross-dataset evaluation, improving the HR error from the traditional methods, highlighting an RMSE of 6.11 BPM or a MAPE of 5.23 %. On the other hand, in the UBFC-rPPG dataset, CZT and FFT exhibit more similar errors. CZT secures better results in MAE and MAPE, while FFT shows slightly better RMSE and R. In contrast, the deep CZT estimator significantly enhances performance across all evaluated metrics, with an RMSE of 6.58 and a MAPE of 4.36%, indicating a more robust HR estimation.

TABLE II: Cross-database heart rate evaluation on PURE and UBFC-rPPG (beats per minute).

Dataset	HR Method	Cross-database evaluation			
		MAE↓	RMSE↓	MAPE↓	R↑
PURE	FFT	4.34±1.23	10.37±69.21	7.32±2.36	0.90±0.06
	CZT	3.41±0.93	7.94±34.08	5.84±1.79	0.94±0.05
	Deep CZT	3.37±0.66	6.11±11.54	5.23±1.11	0.97±0.03
UBFC	FFT	4.83±0.94	7.76±20.97	5.40±1.12	0.89±0.07
	CZT	4.66±0.97	7.84±21.21	5.07±1.09	0.88±0.07
	Deep CZT	4.19±0.78	6.58±12.99	4.36±0.82	0.92±0.06

TABLE III: Performance evaluation of deep CZT estimator in SOTA approaches (beats per minute).

rPPG Approach	HR Method	Intra-database evaluation			
		UCLA-rPPG			
		MAE↓	RMSE↓	MAPE↓	R↑
LGI [38]	FFT	3.28±0.47	4.64±6.21	4.24±0.54	0.93±0.05
	CZT	2.88±0.38	3.92±3.67	3.71±0.44	0.96±0.04
	Deep CZT	2.40±0.32	3.29±2.51	3.12±0.38	0.96±0.04
POS [51]	FFT	3.17±0.50	4.71±8.18	4.40±0.76	0.87±0.07
	CZT	2.62±0.34	3.56±2.67	3.60±0.50	0.93±0.05
	Deep CZT	2.16±0.30	3.00±2.10	2.87±0.38	0.97±0.04
Physnet [56]	FFT	2.27±0.28	2.99±1.76	2.96±0.34	0.98±0.03
	CZT	2.26±0.30	3.10±2.21	2.91±0.37	0.98±0.03
	Deep CZT	1.78±0.24	2.47±1.59	2.32±0.30	0.98±0.03
Efficient-Phys [25]	FFT	2.73±0.45	4.17±6.21	3.73±0.67	0.91±0.06
	CZT	2.13±0.34	3.17±3.23	2.76±0.41	0.95±0.04
	Deep CZT	1.68±0.26	2.45±1.61	2.22±0.32	0.97±0.03
TDM [9]	FFT	2.50±0.34	3.43±3.02	3.25±0.41	0.95±0.04
	CZT	2.27±0.30	3.08±2.23	2.93±0.36	0.97±0.03
	Deep CZT	1.74±0.26	2.53±1.67	2.26±0.32	0.97±0.03

E. Adapting deep CZT estimator to SOTA rPPG approaches

Despite SMO regularization aiding intra-dataset optimization, the primary goal is to preserve the CZT structure. As indicated in subsection IV-C, the standard CZT outperforms commonly used techniques in the literature. Therefore, while SMO regularization adjusts weights for intra-dataset characteristics, it also preserves generalization ability across datasets.

In this section, we assess the performance of our HR estimator module using five selected rPPG methods: POS [51], LGI [38], PhysNet [56], EfficientPhys [25] (implemented from the rPPG-toolbox [27]), and the TDM model [9], previously employed in subsection IV-D. Similarly to the previous section, we compare our deep CZT estimator against FFT and CZT, used as benchmarks for each method’s results, as summarized in Table III.

Results consistently align with the findings presented in our preliminary experiment (Section IV.C), showing that CZT tends to accurately extract the HR value with lower error compared to FFT. Among the five methods, CZT exhibits lower HR error in all metrics except in the Physnet model, where the performance is nearly the same with very small differences. On the other hand, our adaptable deep CZT estimator shows a significant reduction in HR error not only compared to FFT but also compared to the standard CZT.

In terms of rPPG estimation, the use of our deep CZT in traditional methods such as LGI or POS, improves their HR performance by reducing the MAE and the RMSE between 0.5 and 1.71 BPM. A similar trend is observed for data-driven methods: using our deep CZT estimator outperform previous results, reducing the MAE and the RMSE below 2 and 3 BPMs, respectively. These promising results indicate that our deep CZT estimator successfully narrows the gap between the HR extracted from the rPPG signal and the HR value from the dataset sensor. Moreover, it proves to be independent of the adopted rPPG method, whether traditional or data-driven, highlighting its potential to be incorporated into any rPPG model designed for HR estimation.

V. CONCLUSIONS

This paper proposes the use of the CZT for remote heart rate measurement. Its inherent flexibility in adjusting frequency resolution overcomes the limitations of the FFT, enabling more robust and precise HR estimation across various temporal window sizes. Additionally, we introduce a novel data-driven CZT estimator tailored to adapt the classical CZT to the unique characteristics of each HR sensor and the recovered PPG signal.

To guide the adaptation of our model, we propose a frequency distribution loss regularized with sparse matrix optimization, showcasing outstanding results in both intra-database and cross-dataset evaluations. Furthermore, we validate the capability of incorporating our deep CZT estimator into several existing rPPG methods, highlighting its adaptability and improved performance compared with current frequency handcrafted methods.

In conclusion, we have presented an alternative to the FFT for remote heart rate estimation, leveraging the CZT. Our results indicate promising performance, even when employing temporal windows as short as approximately 2 seconds. This implies a comparatively lower delay in the estimate of heart rates and suggests the possibility of exploring applications that may require near-instantaneous HR estimation.

ACKNOWLEDGMENTS

This work is partly supported by the eSCANFace project (PID2020-114083GB-I00) funded by the Spanish Ministry of Science and Innovation.

REFERENCES

- [1] G. Balakrishnan, F. Durand, and J. Guttag. Detecting pulse from head motions in video. In *CVPR*, pages 3430–3437, 2013.
- [2] Y. Benezeth, P. Li, R. Macwan, K. Nakamura, R. Gomez, and F. Yang. Remote heart rate variability for emotional state monitoring. In *EMBS Int. Conf. Biomed. Health Inform. BHI*, pages 153–156. IEEE, 2018.
- [3] L. Bluestein. A linear filtering approach to the computation of discrete fourier transform. *IEEE Transactions on Audio and Electroacoustics*, 18(4):451–455, 1970.
- [4] S. Bobbia, R. Macwan, Y. Benezeth, A. Mansouri, and J. Dubois. Un-supervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognit. Lett.*, 124:82–90, 2019.

- [5] F. Bousefsaf, A. Pruski, and C. Maaoui. 3d convolutional neural networks for remote pulse rate measurement and mapping from facial video. *Applied Sciences*, 9(20):4364, 2019.
- [6] W. Cao, V. Mirjalili, and S. Raschka. Rank consistent ordinal regression for neural networks with application to age estimation. *Pattern Recognition Letters*, 140:325–331, 2020.
- [7] W. Chen and D. McDuff. Deepphys: Video-based physiological measurement using convolutional attention networks. In *ECCV*, pages 349–365, 2018.
- [8] J. Comas, D. Aspandi, and X. Binefa. End-to-end facial and physiological model for affective computing and applications. In *2020 15th IEEE international conference on Automatic Face and Gesture Recognition (FG 2020)*, pages 93–100. IEEE, 2020.
- [9] J. Comas, A. Ruiz, and F. Sukno. Efficient remote photoplethysmography with temporal derivative modules and time-shift invariant loss. In *CVPR*, pages 2182–2191, 2022.
- [10] G. De Haan and V. Jeanne. Robust pulse rate from chrominance-based rppg. *IEEE Trans. Biomed. Eng.*, 60(10):2878–2886, 2013.
- [11] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [12] A. Gudi, M. Bittner, and J. Van Gemert. Real-time webcam heart-rate and variability estimation with clean ground truth for evaluation. *Applied Sciences*, 10(23):8630, 2020.
- [13] A. K. Gupta, R. Kumar, L. Birla, and P. Gupta. Radiant: Better rppg estimation using signal embeddings and transformer. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4976–4986, 2023.
- [14] P. Hansen, M. G. Lozano, F. Kamrani, and J. Brynielsson. Real-time estimation of heart rate in situations characterized by dynamic illumination using remote photoplethysmography. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6093–6102, 2023.
- [15] L. Hou, C.-P. Yu, and D. Samaras. Squared earth movers distance loss for training deep neural networks on ordered-classes. In *NIPS Workshop*, volume 5, 2017.
- [16] R. Irani, K. Nasrollahi, and T. B. Moeslund. Improved pulse detection from head motions using dct. In *2014 international conference on computer vision theory and applications (VISAPP)*, volume 3, pages 118–124. IEEE, 2014.
- [17] T.-P. Jung, T. J. Sejnowski, et al. Utilizing deep learning towards multi-modal bio-sensing and vision-based affective computing. *IEEE Transactions on Affective Computing*, 13(1):96–107, 2019.
- [18] Y. Kartynnik, A. Ablavatski, I. Grishchenko, and M. Grundmann. Real-time facial surface geometry from monocular video on mobile gpus. *arXiv preprint arXiv:1907.06724*, 2019.
- [19] E. Lee, E. Chen, and C.-Y. Lee. Meta-rppg: Remote heart rate estimation using a transductive meta-learner. In *ECCV*, pages 392–409. Springer, 2020.
- [20] J. Lee-Thorp, J. Ainslie, I. Eckstein, and S. Ontanon. Fnet: Mixing tokens with fourier transforms. *arXiv preprint arXiv:2105.03824*, 2021.
- [21] M. Lewandowska, J. Rumiński, T. Kocejko, and J. Nowak. Measuring pulse rate with a webcam—a non-contact method for evaluating cardiac activity. In *FedCSIS*, pages 405–410. IEEE, 2011.
- [22] B. Li, P. Zhang, J. Peng, and H. Fu. Non-contact ppg signal and heart rate estimation with multi-hierarchical convolutional network. *Pattern Recognition*, 139:109421, 2023.
- [23] X. Li, J. Chen, G. Zhao, and M. Pietikainen. Remote heart rate measurement from face videos under realistic situations. In *CVPR*, pages 4264–4271, 2014.
- [24] S. Liu, P. C. Yuen, S. Zhang, and G. Zhao. 3d mask face anti-spoofing with remote photoplethysmography. In *ECCV*, pages 85–100. Springer, 2016.
- [25] X. Liu, B. Hill, Z. Jiang, S. Patel, and D. McDuff. Efficientphys: Enabling simple, fast and accurate camera-based cardiac measurement. In *IEEE Winter Conf. Appl. Comput. Vis.*, pages 5008–5017, 2023.
- [26] X. Liu, Z. Jiang, J. Fromm, X. Xu, S. Patel, and D. McDuff. Metaphys: few-shot adaptation for non-contact physiological measurement. In *CHIL*, pages 154–163, 2021.
- [27] X. Liu, X. Zhang, G. Narayanswamy, Y. Zhang, Y. Wang, S. Patel, and D. McDuff. Deep physiological sensing toolbox. *arXiv preprint arXiv:2210.00716*, 2022.
- [28] H. Lu, H. Han, and S. K. Zhou. Dual-gan: Joint bvp and noise modeling for remote physiological measurement. In *CVPR*, pages 12404–12413, 2021.
- [29] D. McDuff, S. Gontarek, and R. W. Picard. Improvements in remote cardiopulmonary measurement using a five band digital camera. *IEEE Transactions on Biomedical Engineering*, 61(10):2593–2601, 2014.
- [30] D. McDuff, S. Gontarek, and R. W. Picard. Remote detection of photoplethysmographic systolic and diastolic peaks using a digital camera. *IEEE Transactions on Biomedical Engineering*, 61(12):2948–2954, 2014.
- [31] X. Niu, H. Han, S. Shan, and X. Chen. Synrhythm: Learning a deep heart rate estimator from general to specific. In *ICPR*, pages 3580–3585. IEEE, 2018.
- [32] X. Niu, S. Shan, H. Han, and X. Chen. Rhythmnet: End-to-end heart rate estimation from face via spatial-temporal representation. *IEEE TIP*, 29:2409–2423, 2019.
- [33] E. M. Nowara, D. McDuff, and A. Veeraraghavan. The benefit of distraction: Denoising camera-based physiological measurements using inverse attention. In *ICCV*, pages 4955–4964, 2021.
- [34] J. Pan and W. J. Tompkins. A real-time qrs detection algorithm. *IEEE transactions on biomedical engineering*, (3):230–236, 1985.
- [35] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Adv Neural Inf Process Syst*, 32, 2019.
- [36] T. Pereira, N. Tran, K. Gadhoumi, M. M. Pelter, D. H. Do, R. J. Lee, R. Colorado, K. Meisel, and X. Hu. Photoplethysmography based atrial fibrillation detection: a review. *NPJ digital medicine*, 3(1):3, 2020.
- [37] O. Perepelkina, M. Artemyev, M. Churikova, and M. Grinenko. Hearttrack: Convolutional neural network for remote video-based heart rate monitoring. In *CVPRW*, pages 288–289, 2020.
- [38] C. S. Pilz, S. Zaunseder, J. Krajewski, and V. Blazek. Local group invariance for heart rate estimation from face videos in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 1254–1262, 2018.
- [39] M.-Z. Poh, D. J. McDuff, and R. W. Picard. Advancements in non-contact, multiparameter physiological measurements using a webcam. *IEEE Trans. Biomed. Eng.*, 58(1):7–11, 2010.
- [40] M.-Z. Poh, D. J. McDuff, and R. W. Picard. Non-contact, automated cardiac pulse measurements using video imaging and blind source separation. *Optics express*, 18(10):10762–10774, 2010.
- [41] L. Rabiner, R. W. Schafer, and C. Rader. The chirp z-transform algorithm. *IEEE transactions on audio and electroacoustics*, 17(2):86–92, 1969.
- [42] A. Revanur, Z. Li, U. A. Ciftci, L. Yin, and L. A. Jeni. The first vision for vitals (v4v) challenge for non-contact video-based physiological estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2760–2767, 2021.
- [43] V. Ronca, A. Giorgi, D. Rossi, A. Di Florio, G. Di Flumeri, P. Aricò, N. Sciaraffa, A. Vozzi, L. Tamborra, I. Simonetti, et al. A video-based technique for heart rate and eye blinks rate estimation: A potential solution for telemonitoring and remote healthcare. *Sensors*, 21(5):1607, 2021.
- [44] R. Song, H. Chen, J. Cheng, C. Li, Y. Liu, and X. Chen. PulseGAN: Learning to generate realistic pulse waveforms in remote photoplethysmography. *IEEE J.Biomed.Health Inform.*, 25(5):1373–1384, 2021.
- [45] R. Špetlík, V. Franc, and J. Matas. Visual heart rate estimation with convolutional neural network. In *BMVC*, 2018.
- [46] C. Takano and Y. Ohta. Heart rate measurement based on a time-lapse image. *Medical engineering & physics*, 29(8):853–857, 2007.
- [47] Z.-H. Tan, Y. Xie, Y. Jiang, and Z.-H. Zhou. Real-valued backpropagation is unsuitable for complex-valued neural networks. *Advances in Neural Information Processing Systems*, 35:34052–34063, 2022.
- [48] S. Tulyakov, X. Alameda-Pineda, E. Ricci, L. Yin, J. F. Cohn, and N. Sebe. Self-adaptive matrix completion for heart rate estimation from face videos under realistic conditions. In *CVPR*, pages 2396–2404, 2016.
- [49] R. Velik. Discrete fourier transform computation using neural networks. In *2008 International Conference on Computational Intelligence and Security*, volume 1, pages 120–123. IEEE, 2008.
- [50] W. Verkruysse, L. O. Svaasand, and J. S. Nelson. Remote plethysmographic imaging using ambient light. *Optics express*, 16(26):21434–21445, 2008.

- [51] W. Wang, A. C. den Brinker, S. Stuijk, and G. De Haan. Algorithmic principles of remote ppg. *IEEE Trans. Biomed. Eng.*, 64(7):1479–1491, 2016.
- [52] W. Wang, S. Stuijk, and G. De Haan. A novel algorithm for remote photoplethysmography: Spatial subspace rotation. *IEEE Trans. Biomed. Eng.*, 63(9):1974–1984, 2015.
- [53] Z. Wang, Y. Ba, P. Chari, O. D. Bozkurt, G. Brown, P. Patwa, N. Vaddi, L. Jalilian, and A. Kadambi. Synthetic generation of face videos with plethysmograph physiology. In *CVPR*, pages 20587–20596, 2022.
- [54] Z.-K. Wang, Y. Kao, and C.-T. Hsu. Vision-based heart rate estimation via a two-stream cnn. In *ICIP*, pages 3327–3331. IEEE, 2019.
- [55] Y.-P. Yu, B.-H. Kwan, C.-L. Lim, S.-L. Wong, and P. Raveendran. Video-based heart rate measurement using short-time fourier transform. In *2013 International Symposium on Intelligent Signal Processing and Communication Systems*, pages 704–707. IEEE, 2013.
- [56] Z. Yu, X.-B. Li, and G. Zhao. Remote photoplethysmograph signal measurement from facial videos using spatio-temporal networks. In *BMVC*, 2019.
- [57] Z. Yu, W. Peng, X. Li, X. Hong, and G. Zhao. Remote heart rate measurement from highly compressed facial videos: an end-to-end deep learning solution with video enhancement. In *ICCV*, pages 151–160, 2019.
- [58] Z. Yu, Y. Shen, J. Shi, H. Zhao, Y. Cui, J. Zhang, P. Torr, and G. Zhao. Physformer++: Facial video-based physiological measurement with slowfast temporal difference transformer. *IJCV*, 131(6):1307–1330, 2023.
- [59] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Sign. Process. Letters*, 23(10):1499–1503, 2016.