

# Orientation-conditioned Facial Texture Mapping for Video-based Facial Remote Photoplethysmography Estimation

## Supplementary Material

### 7. Supplementary Material

#### 7.1. Data Visualization

We provide select video frames approved for publication from the PURE [4] dataset after different video processing steps but excluding the normalized frame-difference, pixel outlier clipping and standardization steps ( $F_D$ ) for visualization and comparison.

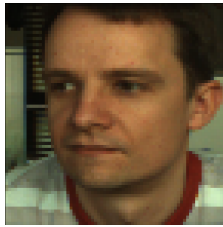


Figure 6. Frame of *Subject 1* in scenario *Medium Rotation* from PURE [4] with static cropping ( $\times 1.5$ -scale Box) applied. Subsequent frames will use the same bounding box, hence the in-plane position of the face will vary due to subject motion.

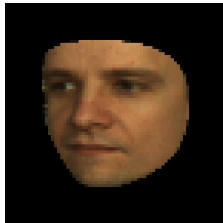


Figure 7. Frame of *Subject 1* in scenario *Medium Rotation* from PURE [4] with static cropping ( $\times 1.5$ -scale Box) and facial segmentation applied. Subsequent frames will use the same bounding box, hence the in-plane position of the face will vary due to subject motion.

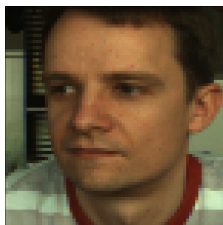


Figure 8. Frame of *Subject 1* in scenario *Medium Rotation* from PURE [4] with dynamic cropping ( $\times 1.5$ -scale Box) and square padding applied. However, subsequent frames will remain centered on the square padded and scaled facial region allowing for larger in-plane subject motion within a video sequence.



Figure 9. Frame of *Subject 1* in scenario *Medium Rotation* from PURE [4] with UV transformation and masking ( $\Theta \geq 45^\circ$ ) applied - the UV transformation process inherently dynamically localizes and segments the facial region. Subsequent frames will have the same structure with varying texture.

#### 7.2. Intra-dataset Testing

In Table 5 we report the full set of performance metrics referenced in Section 4.3 using the evaluation pipeline and metric implementations provided in [15] for intra-dataset testing on the PURE [33] dataset using subject-independent cross-validation. We obtain these results using the protocol described in Section 4.4. We report both the results including and excluding samples from *Subject 7 - Talking (S7-T)* to provide insight into the evaluation variability. We denote the sequence first-order normalized frame difference, pixel outlier clipping, and standardization operations as  $F_D$  for brevity. We also denote the UV transformation operation as  $T_{UV}$ .

#### 7.3. Cross-dataset Testing

In Table 6 we report the full set of performance metrics referenced in Section 4.3 using the evaluation pipeline and metric implementations provided in [15]. We obtained these results using the protocol described in Section 4.5, we perform cross-dataset testing on the MMPD [34] dataset using PhysNet models trained on the PURE [4] dataset. We report additional ablations for the operations applied after  $T_{UV}$  to demonstrate the impact of the sequence of operations, and provide additional internally consistent comparisons. We denote the sequence first-order normalized frame difference, pixel outlier clipping, and standardization operations as  $F_D$  for brevity. We also denote the UV transformation operation as  $T_{UV}$ .

Video Processing Pipeline	MAE $\pm$ SE (BPM)	RMSE $\pm$ SE (BPM)	$r \pm$ SE	SNR $\pm$ SE (dB)
$F_D + \text{Crop}_{Static} (\times 1.5\text{-Box}) + \text{Resize}$	$0.492 \pm 0.172$	$1.408 \pm 0.946$	$0.998 \pm 0.008$	$10.721 \pm 1.044$
$\text{Crop}_{Static} (\times 1.5\text{-Box}) + \text{Resize} + F_D$ ( <b>PhysNet-XY</b> )	$1.318 \pm 0.979$	$7.632 \pm 56.531$	$0.945 \pm 0.043$	$11.061 \pm 1.025$
$\text{Crop}_{Static} (\times 1.5\text{-Box}) + \text{Resize} + F_D$ (Excl. S7-T)	$0.341 \pm 0.138$	$1.108 \pm 0.727$	$0.999 \pm 0.007$	$11.457 \pm 0.963$
$T_{UV} + F_D + \text{Resize}$	$2.734 \pm 1.510$	$11.918 \pm 98.699$	$0.862 \pm 0.067$	$11.546 \pm 1.135$
$T_{UV} + F_D + \text{Resize}$ (Excl. S7-T)	$0.594 \pm 0.217$	$1.739 \pm 1.435$	$0.996 \pm 0.011$	$12.228 \pm 1.063$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 90^\circ) + F_D + \text{Resize}$	$1.393 \pm 0.938$	$7.338 \pm 51.500$	$0.949 \pm 0.042$	$12.011 \pm 1.140$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 90^\circ) + F_D + \text{Resize}$ (Excl. S7-T)	$0.462 \pm 0.171$	$1.381 \pm 0.959$	$0.998 \pm 0.008$	$12.470 \pm 1.063$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 60^\circ) + F_D + \text{Resize}$	$1.676 \pm 1.316$	$10.243 \pm 102.807$	$0.899 \pm 0.058$	$12.211 \pm 1.084$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 60^\circ) + F_D + \text{Resize}$ (Excl. S7-T)	$0.356 \pm 0.139$	$1.114 \pm 0.727$	$0.999 \pm 0.007$	$12.617 \pm 1.023$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 45^\circ) + F_D + \text{Resize}$ ( <b>PhysNet-XY</b> )	$1.639 \pm 1.141$	$8.919 \pm 76.940$	$0.924 \pm 0.051$	$11.842 \pm 1.106$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 45^\circ) + F_D + \text{Resize}$ (Excl. S7-T)	$0.500 \pm 0.171$	$1.397 \pm 0.958$	$0.998 \pm 0.008$	$12.159 \pm 1.079$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 30^\circ) + F_D + \text{Resize}$	$1.594 \pm 1.113$	$8.693 \pm 72.994$	$0.928 \pm 0.049$	$11.486 \pm 1.108$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 30^\circ) + F_D + \text{Resize}$ (Excl. S7-T)	$0.485 \pm 0.172$	$1.397 \pm 0.958$	$0.998 \pm 0.008$	$11.817 \pm 1.077$

Table 5. Intra-dataset subject-independent performance of PhysNet across different video processing pipelines on the PURE [33] dataset using averaged results across all folds from subject-independent cross-validation training on the PURE [4] dataset.

Video Processing Pipeline	MAE $\pm$ SE (BPM)	RMSE $\pm$ SE (BPM)	$r \pm$ SE	SNR $\pm$ SE (dB)
$F_D + \text{Crop}_{Static} (\times 1.5\text{-Box}) + \text{Resize}$	$17.492 \pm 0.307$	$24.827 \pm 16.908$	$0.047 \pm 0.017$	$-6.225 \pm 0.074$
$\text{Crop}_{Static} (\times 1.5\text{-Box}) + \text{Resize} + F_D$ ( <b>PhysNet-XY</b> )	$14.905 \pm 0.295$	$22.542 \pm 15.837$	$0.155 \pm 0.017$	$-6.882 \pm 0.080$
$\text{Crop}_{Static} (\times 1.5\text{-Box}) + \text{Segment} + \text{Resize} + F_D$	$15.237 \pm 0.312$	$23.524 \pm 17.217$	$0.120 \pm 0.017$	$-6.053 \pm 0.088$
$\text{Crop}_{Dynamic} (\times 1.5\text{-Box}) + \text{Pad}_{Square} + \text{Resize} + F_D$	$17.988 \pm 0.307$	$25.183 \pm 16.488$	$0.033 \pm 0.017$	$-6.263 \pm 0.072$
$\text{Crop}_{Dynamic} (\times 1.5\text{-Box}) + \text{Pad}_{Square} + \text{Segment} + \text{Resize} + F_D$	$14.683 \pm 0.298$	$22.563 \pm 15.526$	$0.138 \pm 0.017$	$-6.553 \pm 0.082$
$T_{UV} + \text{Resize} + F_D$	$13.168 \pm 0.285$	$21.014 \pm 14.394$	$0.227 \pm 0.017$	$-6.606 \pm 0.084$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 90^\circ) + \text{Resize} + F_D$	$13.547 \pm 0.288$	$21.391 \pm 14.610$	$0.210 \pm 0.017$	$-6.644 \pm 0.085$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 60^\circ) + \text{Resize} + F_D$	$12.949 \pm 0.284$	$20.840 \pm 14.129$	$0.243 \pm 0.017$	$-6.305 \pm 0.087$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 45^\circ) + \text{Resize} + F_D$	$15.222 \pm 0.302$	$23.072 \pm 15.474$	$0.156 \pm 0.017$	$-6.105 \pm 0.083$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 30^\circ) + \text{Resize} + F_D$	$15.771 \pm 0.298$	$23.285 \pm 15.670$	$0.133 \pm 0.017$	$-6.643 \pm 0.082$
$T_{UV} + F_D + \text{Resize}$	$12.687 \pm 0.280$	$20.454 \pm 13.843$	$0.248 \pm 0.017$	$-6.679 \pm 0.088$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 90^\circ) + F_D + \text{Resize}$	$13.038 \pm 0.285$	$20.900 \pm 14.249$	$0.216 \pm 0.017$	$-6.473 \pm 0.086$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 60^\circ) + F_D + \text{Resize}$	$12.890 \pm 0.280$	$20.629 \pm 13.794$	$0.256 \pm 0.017$	$-6.284 \pm 0.088$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 45^\circ) + F_D + \text{Resize}$ ( <b>PhysNet-UV</b> )	$12.187 \pm 0.273$	$19.849 \pm 13.102$	$0.294 \pm 0.017$	$-6.265 \pm 0.092$
$T_{UV} + \text{Mask} (\Theta_{UV} \geq 30^\circ) + F_D + \text{Resize}$	$13.300 \pm 0.279$	$20.834 \pm 13.611$	$0.277 \pm 0.017$	$-6.496 \pm 0.087$

Table 6. Cross-dataset performance of PhysNet across different video processing pipelines on the MMPD [33] dataset using averaged results across all folds from subject-independent cross-validation training on the PURE [33] dataset.