

1 We would like to thank the reviewers for their extensive and constructive feedback and are glad they found our work a
2 *strong contribution*. We have made an effort to answer all their comments and will update our paper accordingly.

3 **R2 & R5: Why is a high frame rate important?** The standard sampling rate for contact pulse and respiration sensors
4 is between 100-256Hz. This is because derivative information is often computed from these waveforms. For example,
5 for heart and respiration rate variability measurement 20 Hz is not sufficient because peak timings are too imprecise.

6 **R2 & R3: How does the illumination and surface reflection help in understanding the resp or pulse waveforms?**

7 The volumetric changes of blood in the surface of the skin cause changes in light absorption (\uparrow volume of hemoglobin =
8 \uparrow light absorption). This in turns affects the amount of visible light reflected from the skin, this is the source of the
9 photoplethysmogram. The mechanical force of blood pumping around the body also causes subtle motions and these
10 are the source of the ballistocardiogram. These color and motion changes in the video help us extract the pulse signal
11 and heart rate frequency. Respiration is captured mostly via motions of the torso as the subjects breathes in and out. But,
12 the principle of respiratory sinus arrhythmia also means that the color changes of the skin pixels due to blood flow also
13 contain some respiration information. Separating these motions and color changes of interest from specular reflections
14 and other motion noise is the primary task of camera-based physiological measurement.

15 **R1:** 1) We believe the computational benefits offered by our approach are a helpful contribution to the research literature.
16 Applying knowledge distillation and model pruning to Hybrid/3D-CAN are definitely interesting lines of investigation
17 but are beyond the scope here. 2) The attention module learns higher weights for skin regions in the pulse task, and
18 torso regions in the resp. task. This boosts the SNR in the motion branch because changes from other sources can be
19 more effectively ignored. Fig. 3-B shows that the temporal shifting in the motion branch can introduce more signal but
20 also more noise, therefore the attention module is even more important. 3) Hybrid-CAN uses a 2D attention branch
21 with an averaged frame whereas 3D-CAN uses a 3D attention branch with all the frames. We have shown in Table 1
22 that Hybrid-CAN was able to achieve a MAE of 1.12 while 3D-CAN had a MAE of 1.18. We experimented using the
23 middle frame and last frame as the input to the 2D attention branch. We found that the mean frame provided the best
24 results. Since the 2D branch only requires a single frame as input (the average frame), the generated attention mask can
25 be shared with all the frames (10 in our case) in the motion branch, helping significantly reduce computation with little
26 impact on accuracy. 4) TS-CAN sometimes performs better than MTTs-CAN because we train two separate “dedicated”
27 TS-CAN models, one for pulse and one for resp. However, we only train one multi-task (MTTs-CAN) model which
28 has to share the weights for two tasks and is much more computationally efficient but is not always as accurate.

29 **R2:** 1) We agree with that our approach could have applications beyond telehealth monitoring during a pandemic. We
30 wrote our introduction to resonate with the current situation, but we agree it would be good to reduce this emphasis and
31 highlight more applications in other domains. 2) We will add a breakdown of the performance by participant in the
32 supp. material. The st.dev. across subjects was (MTTs-CAN): HR MAE: 1.15 BPM (AFRL) and 4.82 BPM (MMSE)
33 and respiration MAE: 1.62 breaths/min (AFRL), these combined with the violin and Bland-Altman plots will be helpful.
34 We agree about evaluation on patient groups of interest (e.g., COVID and A.Fib.) but a full clinical validation was
35 beyond the scope of this work. Generally speaking we expect that improvements in accuracy on healthy people will also
36 help with performance on patients. 3) Vital sign measurement is still imperative even if the cardiopulmonary symptoms
37 are severe. Scalable and accessible health monitoring has many applications and could help democratize care.

38 **R3:** 1) The spatial resolution of the videos does not need to be very high (658x492 for AFRL) and we actually
39 downsample the face regions to 36x36 pixels, which is quite effective (Chen & McDuff, 2018). Studies have shown that
40 off-the-shelf webcams are sufficient for physiological measurement. We would argue that high resolution cameras are
41 not necessary. 2) The reason we proposed the use of multi-task learning and tensor shifting is to enable our network to
42 run on edge devices in real-time while maintaining a sufficient frame rate and state-of-the-art accuracy. Our goal in Sec.
43 3.2 was to provide an explanation of how to realize on-device real-time temporal and spatial modeling. Table 1 and
44 Figure 3-A capture these results. 3) The Firefly RK-3399 has a CPU with 2 Cortex A72 (1.8Ghz) and 4 Cortex-A53
45 cores. This is equivalent to the CPU in the Snapdragon (Snap) 650 Mobile Platform¹. Snap 650 is a mid-to-low-end
46 mobile platform that Qualcomm launched in 2015. Snap 650 achieves a score of 275 in single-core evaluation and a
47 score of 830 in multi-score evaluation². The modern Snap 865 achieves a score of 903 in single-core evaluation and a
48 score of 3304 in multi-score evaluation. To contextualize this, Snap 650 has a little worse performance than the Apple
49 A9 chip used in the iPhone 6s. Thus, we believe that the computing platform we tested on is a fair benchmark for the
50 class of mobile devices we are targeting. 4) Clothing and facial hair do obscure the skin and can affect the measurement
51 of the pulse signal, indeed a limitation of optical measurement via PPG. Skin type can also have an effect, we coded the
52 skin type and gender of subjects in AFRL using the Fitzpatrick scale. The results (MTTs-CAN) are as follows (Skin
53 Type-MAE): I-2.37 BPM, II-2.62 BPM, 1.77-BPM, IV-0.06 BPM. (Gender-MAE): Men-1.34 BPM, Women-3.92 BPM.

54 **R5:** We will emphasis why our method makes the deployment of camera-based vital measurement more tractable given
55 that the lower computation budget requires less power consumption per inference, we appreciate that suggestion.

¹ <https://www.qualcomm.com/products/snapdragon-650-mobile-platform> ² <https://browser.geekbench.com/>