

A Cloud-Edge-Terminal Collaborative System for Temperature Measurement in COVID-19 Prevention

Zheyi Ma*, Hao Li*, Wen Fang*, Qingwen Liu*, Bin Zhou[†] and Zhiyong Bu[†]

*Dept. of Computer Science and Technology, Tongji University, Shanghai, China

[†]Key Laboratory of Wireless Sensor Network and Communications,

Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, Shanghai, China

Email: {zheyima, lihao1101, wen.fang, qliu}@tongji.edu.cn, {bin.zhou, zhiyong.bu}@mail.sim.ac.cn

Abstract—To prevent the spread of coronavirus disease 2019 (COVID-19), preliminary temperature measurement and mask detection in public areas are conducted. However, the existing temperature measurement methods face the problems of safety and deployment. In this paper, to realize safe and accurate temperature measurement even when a person's face is partially obscured, we propose a cloud-edge-terminal collaborative system with a lightweight infrared temperature measurement model. A binocular camera with an RGB lens and a thermal lens is utilized to simultaneously capture image pairs. Then, a mobile detection model based on a multi-task cascaded convolutional network (MTCNN) is proposed to realize face alignment and mask detection on the RGB images. For accurate temperature measurement, we transform the facial landmarks on the RGB images to the thermal images by an affine transformation and select a more accurate temperature measurement area on the forehead. The collected information is uploaded to the cloud in real time for COVID-19 prevention. Experiments show that the detection model is only 6.1M and the average detection speed is 257ms. At a distance of 1m, the error of indoor temperature measurement is about 3%. That is, the proposed system can realize real-time temperature measurement in public areas.

Index Terms—COVID-19 Prevention, Mobile Temperature Measurement, Cloud-Edge-Terminal Collaborative System

I. INTRODUCTION

In December 2019, a novel coronavirus named coronavirus disease 2019 (COVID-19) was identified in Wuhan, China, and spread quickly around the world. A year later in December 2020, more than 70 million people were infected with the virus, causing more than 1.6 million deaths [1]. The research [2] shows that the clinical features of COVID-19 pneumonia are similar to other pneumonia, but liver function damage is more frequent in COVID-19 than in non-COVID-19 patients. That is, COVID-19 is a more infectious and dangerous disease.

Fever is an initial symptom of COVID-19 [3]. Therefore, temperature measurement is an effective way to find patients. Nowadays, pedestrians are forced to measure their body temperature before entering public areas, such as subway stations and train stations. There are two main methods of temperature measurement:

The work was supported by the National Key Research and Development Project under Grant 2020YFB2103900 and Grant 2020YFB2103902. It was also supported by the National Natural Science Foundation of China under Grant 61771344 and Grant 62071334.

- Handheld infrared thermometer for short-distance temperature measurement.
- Thermal camera for remote temperature measurement.

The infrared thermometer is cheap and easy to operate, but it requires close contact between the inspector and pedestrians. It is unsafe when pedestrians carry the virus. To realize safe temperature measurement, thermal cameras which can measure temperature from a long distance have been developed. However, the current design of thermal cameras usually requires the purchase of other hardware, such as displays and computers, to support the operation of the system. In addition, body temperature information cannot be used by the system for COVID-19 prevention effectively and timely. Finally, to prevent the spread of the virus through respiratory droplets, people always wear masks in public areas. Thus, mask detection in public areas needs to be taken into account.

A skin temperature extraction method, proposed by Aryal in [4], uses the binocular camera with an RGB lens and a thermal lens to detect people's skin temperature. However, the occlusion problem and the method capable of accurately positioning temperature measurement areas were not analyzed in the paper. In this paper, to accurately measure body temperature even when a person's face is partially obscured, we propose a cloud-edge-terminal collaborative system for remote temperature measurement and develop a detection model which can select a more accurate temperature measurement area on the forehead. As shown in Fig. 1, the detection model can achieve real-time detection on mobile devices. From the RGB and thermal images, we can get the person's temperature and whether he is wearing a mask. Finally, the cloud-edge-terminal collaborative system uploads the location and body temperature information to the cloud.

The contributions of this paper include:

- We propose a cloud-edge-terminal collaborative system, which has a lightweight face alignment model with a mask detection branch and is easy to deploy in mobile devices. The system can upload location and body temperature information to the cloud in real time.
- An accurate affine transformation matrix is calculated to align the image pairs taken by the binocular camera.
- Experiments show that our model can achieve real-time

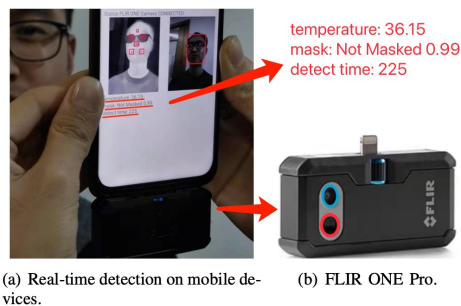


Fig. 1. Real-time detection on the mobile device.

temperature measurement on mobile devices. The detection model is only 6.1M and the average detection speed is 257ms. The error of indoor temperature measurement is about 3% at a distance of 1m.

In the rest of this paper, we show the design of the cloud-edge-terminal collaborative system in section II. In section III, we present the infrared temperature measurement method. Then, a face alignment model with a mask detection branch is proposed and described in section IV. In addition, we depict our experiments in section V. Finally, we make a conclusion in section VI.

II. CLOUD-EDGE-TERMINAL COLLABORATIVE SYSTEM

In this section, a cloud-edge-terminal collaborative system is proposed. With the help of this system, the timeliness of data transmission and the stability of the temperature measurement system can be ensured.

The work process of the cloud-edge-terminal collaborative system is shown in Fig. 2. At first, the input of the system is a pair of RGB and thermal images captured at the same time. Secondly, if faces exist in the RGB image, the facial landmark detection model can locate the landmarks. Otherwise, the camera retakes the valid images. Next, an affine transformation is performed and facial landmarks on the RGB image can be transformed into facial landmarks on the thermal image. Then, these facial landmarks are used for temperature measurement in the thermal image. Finally, the location and temperature information are uploaded to the cloud.

A. Image Capturing

As shown in Fig. 1(b), FLIR ONE Pro is a professional thermal camera, which is a binocular camera with an RGB lens marked by a blue circle and a thermal lens with a red circle [5]. Due to the demand for mobile device deployment, we choose FLIR ONE Pro to capture image pairs. In addition, FLIR ONE Pro has an accuracy of $\pm 3^\circ\text{C}$ or $\pm 5\%$ for temperature measure when the device is between 15°C and 35°C and the scene is between 5°C and 120°C . That is, the recommended working environment temperature meets the needs of daily body temperature detection.

Moreover, the thermal camera can take a pair of RGB and thermal images at the same time. The RGB image resolution

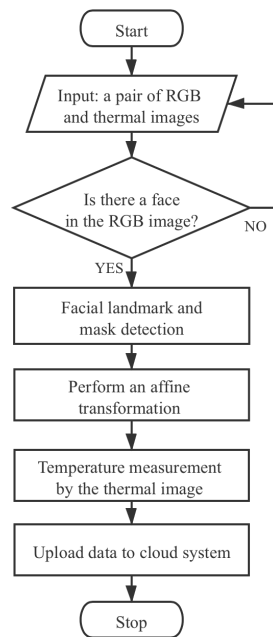


Fig. 2. Flow chart of temperature measurement.

is 1440×1080 pixels and the thermal image resolution is 640×480 pixels. We can get facial landmarks from the RGB images and get body temperature from the thermal images.

B. Cloud-Edge-Terminal Collaborative System

After aligning the image pairs, we can select an accurate area for temperature. However, a small amount of data cannot be mined for additional information, we need to collect data to analyze COVID-19. To improve the efficiency of information transmission, a cloud-edge-terminal collaborative system shown in Fig. 3 is necessary.

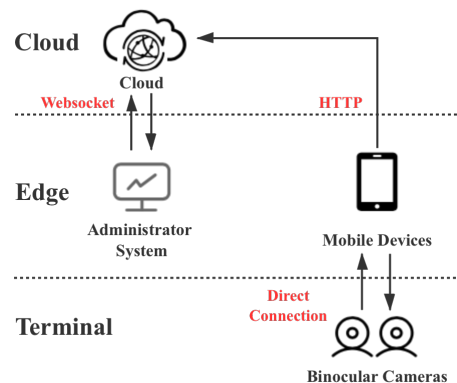


Fig. 3. Cloud-edge-terminal collaborative system design.

1) *Cloud and Mobile Devices Communication:* The binocular camera is directly connected to the mobile device through

Type-C port, and the computing power of mobile devices is utilized to process the captured images. Once a pedestrian is detected, we record the position of the mobile device and the temperature of the pedestrian. The position and temperature of pedestrians can be uploaded to the cloud in real time by hypertext transfer protocol (HTTP) [6], if the pedestrian's body temperature is abnormal.

2) *Cloud and Administrator System Communication*: HTTP is a fundamental web protocol, and information can be sent to the cloud in one-way via HTTP. However, the cloud cannot actively send information to the administrator. The WebSocket Protocol, proposed in [7], enables two-way communication between a client and a host. We utilize WebSocket to realize full-duplex communication between the cloud and the administrator system. When the temperature exceeds the threshold, the cloud will immediately send an alert signal to the administrator system.

III. INFRARED TEMPERATURE MEASUREMENT METHOD

In this section, we present a mobile framework design for infrared temperature measurement. To detect on mobile devices by a specialized thermal camera, a lightweight detection model is proposed and deployed below.

A. Facial Landmark Detection

Facial landmark detection is an important part of infrared temperature measurement. With the help of facial landmarks, we can obtain body temperature from a more accurate area.

Due to the influence of COVID-19, more and more people are used to wearing face masks. As depicted in Fig. 4, The temperature of these areas covered by face masks or glasses is lower than other areas. Therefore, the temperature of a random area or the entire face area cannot be recognized as the human body temperature.

There are 5 most critical points on the face, including the left and right corners of the mouth, the center of the two eyes, and the nose. These points are the internal key points of the face. Since the forehead area is the least likely to be covered in the current scene, we prefer to use the temperature of this area as the human body temperature. The localization method of these 5 points will be shown in section IV.

B. Image Alignment

As shown in Fig. 4, a pair of RGB and thermal images have different image resolutions. Besides, binocular disparity refers to the difference in image location of an object seen by the left and right eyes. Similarly, there is a disparity between a pair of images taken by a binocular camera because the two lenses cannot overlap in physical space. We cannot directly apply the landmark coordinate (x, y) on the RGB image to the thermal image because of the resolution and disparity.

However, the spatial positions of the two lenses are fixed. That is, the straightness and parallelism of the RGB image will not change in the thermal image. Therefore, we can use affine transformation to align a pair of RGB and thermal images.

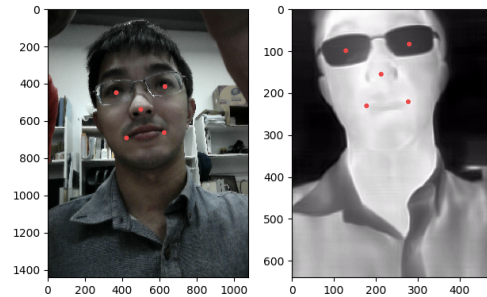


Fig. 4. A pair of images with landmarks taken by FLIR ONE Pro.

An affine transformation includes scaling, translation, rotation, reflection, and shearing. It can be depicted as:

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = A \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & t_x \\ a_3 & a_4 & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}, \quad (1)$$

where the coordinate (x, y) represents the landmark on the RGB image and (x', y') represents the landmark on the thermal image. An affine transformation matrix A transforms (x, y) to (x', y') . The translation of landmarks is controlled by parameters t_x and t_y . Other transformations are controlled by parameter a_{1-4} . If these six unknown parameters are determined, the affine transformation matrix A can be determined. We can use the matrix to locate the facial landmarks on the thermal image according to the facial landmarks detected on the RGB image. Afterwards, the temperature can be accurately measured.

We propose a loss calculation method to evaluate the accuracy of the affine transformation matrix. As shown below, Eq. (2) represents the loss L_x on the x-axis, Eq. (3) represents the loss L_y on the y-axis and Eq. (4) represents the Euclidean distance loss L_{euc} .

$$L_x = \sum_i^n |x'_i - x_i^{pred}| / (width_t \times n) \quad (2)$$

$$L_y = \sum_i^n |y'_i - y_i^{pred}| / (height_t \times n) \quad (3)$$

$$L_{euc} = \sum_i^n \sqrt{(x'_i - x_i^{pred})^2 + (y'_i - y_i^{pred})^2} / (diag_t \times n) \quad (4)$$

where n denotes the number of point pairs, (x'_i, y'_i) is the coordinate of a marked point, and (x_i^{pred}, y_i^{pred}) is the coordinate of the prediction result. Parameters $width_t$, $height_t$, and $diag_t$ are the width, height and diagonal length of the thermal image, respectively.

IV. MTCNN WITH A MASK DETECTION BRANCH

In this section, we propose a two-stage model based on multi-task cascaded convolutional networks (MTCNN) [8] and a single shot multibox detector (SSD) [9] detection model.

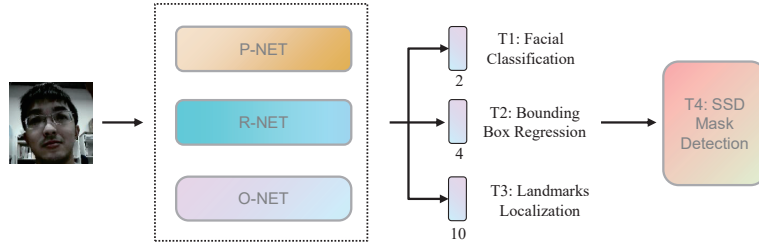


Fig. 5. Overview of MTCNN with a mask detection branch.

The task of MTCNN is face detection and face alignment, while mask detection is the goal of SSD. Since both MTCNN and SSD are lightweight architectures, they are suitable for running on mobile devices. With this model, we can find facial areas with a smaller coverage area to improve temperature measurement results.

A. MTCNN Backbone

MTCNN has a cascaded architecture, which has face detection and face alignment results. The output of MTCNN includes 3 tasks: face classification, bounding box regression, and facial landmark localization. The cascaded architecture contains three nets: P-Net, R-Net, and O-Net. The layers of each net are shown in TABLE I.

TABLE I
LAYERS OF MTCNN PROPOSED IN [8]

Layer Number	P-Net	R-Net	O-Net
Input Size	12×12×3	24×24×3	48×48×3
Layer 1	Conv:3×3 MP:2×2	Conv:3×3 MP:3×3	Conv:3×3 MP:3×3
Layer 2	Conv:3×3	Conv:3×3 MP:3×3	Conv:3×3 MP:3×3
Layer 3	Conv:3×3	Conv:2×2	Conv:3×3 MP:2×2
Layer 4	-	FC128	Conv:2×2
Layer 5	-	-	FC256

MTCNN detects faces from coarse to fine. A cascaded architecture means that the output of the previous network is the input of the next network. The previous network gives a rough judgment firstly, and quickly deletes the areas that do not contain faces. The results are filtered by the next complex network to obtain accurate areas which contain faces.

The output of all three subnets is the same, including facial classification, bounding box regression, and facial landmark localization. As illustrated in Fig. 5, for each sample x_i , MTCNN has three tasks:

1) *Facial classification*: This task contains 2 units as output, representing the probability of whether the image is a face. The loss function of facial classification L_i^{cls} is [8]:

$$L_i^{\text{cls}} = -(y_i^{\text{cls}} \log(p_i) + (1 - y_i^{\text{cls}}) (1 - \log(p_i))), \quad (5)$$

where p_i is the probability that sample x_i is a face and $y_i^{\text{cls}} \in \{0, 1\}$ means the ground-truth label.

2) *Bounding box regression*: This task contains 4 units as output, representing the top-left coordinate, width, and height of the bounding box. The loss function of bounding box regression L_i^{box} is [8]:

$$L_i^{\text{box}} = \|\hat{y}_i^{\text{box}} - y_i^{\text{box}}\|_2^2, \quad (6)$$

where \hat{y}_i^{box} is the prediction coordinate and y_i^{box} is the ground-truth coordinate. There are four parameters, including left, top, height, and width.

3) *Facial landmark localization*: This task contains 10 units as output, representing the coordinates of 5 facial landmarks. The loss function of facial landmark localization L_i^{landmark} is [8]:

$$L_i^{\text{landmark}} = \|\hat{y}_i^{\text{landmark}} - y_i^{\text{landmark}}\|_2^2, \quad (7)$$

where $\hat{y}_i^{\text{landmark}}$ is the prediction coordinate and y_i^{landmark} is the ground-truth coordinate. There are ten parameters, including the coordinates of the left and right corners of the mouth, the center of the two eyes, and the nose.

B. Mask Detection Branch

To avoid the spread of the virus, wearing masks in public areas is recommended. We add a mask detection branch after the MTCNN backbone. For making the model run in mobile devices, the detection branch uses an SSD architecture network to detect whether pedestrians wear masks on their faces.

SSD is a lightweight one-stage detection model, which can quickly and accurately detect masks. The location and classification layers are designed as TABLE II.

TABLE II
SSD ANCHOR CONFIGURATION

Multibox Layers	Feature Map Size	Anchor Size	Aspect Ratio
Layer 1	33×33	0.04, 0.056	1, 0.62, 0.42
Layer 2	17×17	0.08, 0.11	1, 0.62, 0.42
Layer 3	9×9	0.16, 0.22	1, 0.62, 0.42
Layer 4	5×5	0.32, 0.45	1, 0.62, 0.42
Layer 5	3×3	0.64, 0.72	1, 0.62, 0.42

In this branch, we add another new task mask detection, the loss function is similar to facial classification. We use a cross entropy loss function, which is often used in classification problems. The mask detection loss function L_i^m is:

$$L_i^m = -(y_i^m \log(p_i) + (1 - y_i^m)(1 - \log(p_i))), \quad (8)$$

where p_i is the probability that sample x_i wearing a mask and $y_i^m \in \{0, 1\}$ means the ground-truth label.

V. EXPERIMENTS AND ANALYSIS

In this section, we first estimate the affine transformation matrix of the binocular thermal camera. The random sample consensus (RANSAC) estimation method, proposed in [10], is applied to estimate the transformation matrix from our self-made image pairs. Then, the face alignment model with a mask branch is trained. Finally, we develop a TensorFlow Lite model that can be run on Android devices and detection results are uploaded to the cloud-edge-terminal collaborative system in real time.

A. Affine Transformation Matrix

The RANSAC-based robust method is applied to calculate the affine transformation matrix of image pairs captured by the binocular camera. RANSAC is a robust estimation with a two-stage process: i) Classify data points as outliers or inliers. ii) Fit model to inliers while ignoring outliers. In this experiment, we only use a few datasets to calculate the affine transformation matrix and get a good result.

We utilize RANSAC to fit the affine transformation by two sets of point matches: *source* and *target*. The set *source* is the coordinate set of points marked on the RGB image and the set *target* is the coordinate set of points marked on the thermal image. Equation (1) shows that affine transformation has 6 degrees of freedom, so a pair of images and 3 point matches are adequate to estimate the affine transformation matrix. To reduce manual error produced in the marking process, we marked more point matches.

As shown in Fig. 6, we marked 20 pairs of RGB and thermal images for image alignment, and about 10 point matches will be marked on each pair of images. 10 pairs of RGB and thermal images are used for estimation and the other 10 pairs are for testing.

The marked point matches of 10 pairs of RGB and thermal images, *source* and *target*, are the input of the RANSAC method. And we get the affine transformation matrix A of FLIR ONE Pro:

$$A = \begin{bmatrix} 0.5584 & -0.0062 & -65.9722 \\ -0.0014 & 0.5770 & -156.8899 \\ 0 & 0 & 1 \end{bmatrix}$$

A detection result can be seen in Fig. 6, we apply the affine transformation matrix to the test dataset. On the left is an RGB image with marked points, and on the right is a thermal image with marked points and prediction results. As we can see, blue points cover the marked points on the thermal image,

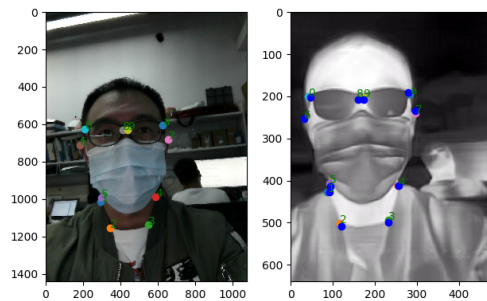


Fig. 6. Transformation result (blue points on the thermal image are prediction result).

which means that the prediction result estimated by the affine transformation matrix is close to the correct result.

Matrix A is applied to estimate the landmarks on the thermal images in our test dataset, and the transformation losses of the test thermal images are presented in TABLE III. L_x , L_y and L_{euc} are all less than 1%, which means the transformation result is close to the actual result.

TABLE III
TRANSFORMATION LOSS

Figure Number	L_x	L_y	L_{euc}
Figure 1	4.9‰	5.8‰	6.2‰
Figure 2	4.4‰	9.3‰	8.5‰
Figure 3	3.9‰	3.7‰	4.2‰
Figure 4	1.8‰	6.3‰	5.2‰
...
Average	3.9‰	7.1‰	6.7‰

B. Facial Landmarks and Temperature Measurement

1) *MTCNN Backbone Training*: Official WIDER FACE training set in [11] and LFW training set in [12] are used for training facial classification, bounding box regression, and facial landmark localization. WIDER FACE contains over 12000 training images and LFW contains 5590 training images.

After training, the loss of classification is 0.1144, the loss of bounding box is 0.05194 and the loss of landmarks is 0.01991. And in the WIDER FACE testing dataset, the accuracy of face detection is 85.1%.

2) *Mask Detection Branch Training*: In the mask detection branch training, 7959 training images with mask annotations are used. These training images are from the WIDER FACE dataset and added mask annotations. We use open-source mask datasets LFW, AgeDB-30 and CFP-FP from [13] as the test datasets. There are approximately 10,000 test images in each of the three test datasets. The accuracy, precision, and recall of the mask detection branch are calculated as follows:

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \quad (9)$$

$$Precision = \frac{TP}{TP + FP} \quad (10)$$

$$Recall = \frac{TP}{TP + FN} \quad (11)$$

In Eqs. (9), (10), and (11), TP represents the number of true positive samples, TN represents the number of true negative samples, FP represents the number of false positive samples, and FN represents the number of false negative samples. The test result is presented in Fig. 7, the accuracy of the model on the three test datasets is higher than 95%.

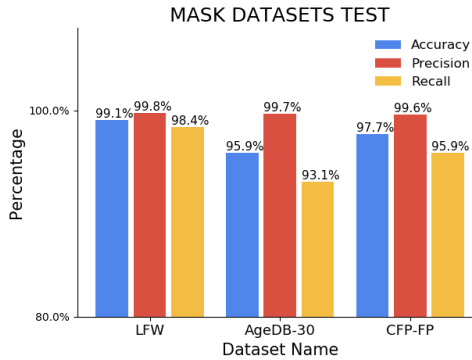


Fig. 7. Test results on three mask datasets(LFW, AgeDB-30, and CFP-PP).

3) *Temperature Measurement*: Figure 6 shows that if you wear a face mask and glasses, most of the face will be covered. Thus, we cannot measure the temperature of these covered areas. For selecting the forehead which is least covered, we need the bounding box and facial landmark results to locate the best temperature measurement area. We take the coordinates of 4 points, including the top-left and top-right corners of the bounding box, the left eye and the right eye. We connect the top-left corner and the right eye, and connect the top-right corner and the left eye. The intersection of the two straight lines is the center of the temperature measurement area.

C. System Deployment

In order to easily detect the body temperature of pedestrians, we need a model that can be run on a mobile platform. Thus, we use TensorFlow Lite to convert the detection model to a mobile detection model. Our experiment platform is Honor V30 with Kirin 980 SoC. A real-time detection experiment result is presented in Fig. 8.

The final size of the entire model is 6.1M. An Android application is developed to apply this mobile model by Android Studio, and the average detection speed is 257ms. At a distance of 1m, the error of indoor temperature measurement is about 3%. Once a face is detected, the location of the mobile device and the temperature are uploaded to the cloud.

VI. CONCLUSIONS

An accurate mobile body temperature measurement system is proposed in this paper. The cloud-edge-terminal system ensures the timeliness of data transmission and the stability of the temperature measurement system. In addition, we calculate the affine transformation matrix of the image pairs taken by the

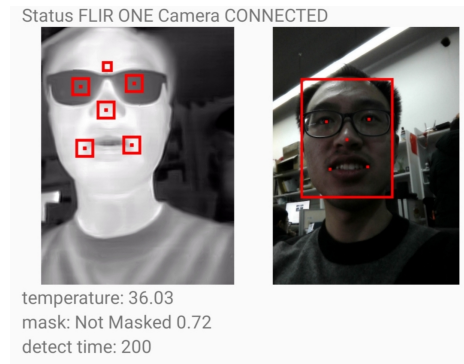


Fig. 8. A real-time detection experiment.

binocular camera. The transformation result, which has only a 6.7% Euclidean error, is close to the actual result. Then, a lightweight face alignment model with a mask detection branch is proposed to realize detection on mobile devices. The average detection speed of our Android model is 257ms, and the error of indoor temperature measurement at a distance of 1m is about 3%, which means that our model can realize real-time and accurate measurement in public areas. Finally, the information can be uploaded to the cloud in real time.

In the future, we plan to prune the network and get a faster and smaller mobile model which will make more contributions to COVID-19 prevention and control.

REFERENCES

- [1] "Who coronavirus disease (covid-19) dashboard," Dec. 2020. [Online]. Available: <https://covid19.who.int/>
- [2] D. Zhao, F. Yao, L. Wang, L. Zheng, Y. Gao, J. Ye, F. Guo, H. Zhao, and R. Gao, "A comparative study on the clinical features of covid-19 pneumonia to other pneumonias," *Clinical Infectious Diseases*, vol. 71, no. 15, pp. 756–761, 2020.
- [3] T. P. Velavan and C. G. Meyer, "The covid-19 epidemic," *Tropical medicine & international health*, vol. 25, no. 3, p. 278, 2020.
- [4] A. Aryal and B. Becerik-Gerber, "Skin temperature extraction using facial landmark detection and thermal imaging for comfort assessment," in *Proceedings of the 6th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation*, 2019, pp. 71–80.
- [5] "Flir one pro," Dec. 2020. [Online]. Available: <https://www.flir.com/products/flir-one-pro/>
- [6] B. A. Forouzan, *TCP/IP protocol suite*. McGraw-Hill, Inc., 2002.
- [7] I. Fette and A. Melnikov, "The websocket protocol," 2011.
- [8] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [10] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [11] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "Wider face: A face detection benchmark," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5525–5533.
- [12] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3476–3483.
- [13] Z. Wang, G. Wang, B. Huang, Z. Xiong, Q. Hong, H. Wu, P. Yi, K. Jiang, N. Wang, Y. Pei, H. Chen, Y. Miao, Z. Huang, and J. Liang, "Masked face recognition dataset and application," *arXiv preprint arXiv:2003.09093*, 2020.