

Completely Contactless and Online Finger Knuckle Identification

Zhenyu ZHOU

June 27, 2022

1 Abstract

2 Introduction

3 Matching Contactless Finger Knuckle

One of our contributions is the online finger knuckle identification. In this kind of situation, we choose the RFNet [8] as our feature extraction backbone, because the model not only is lightweight enough, but it achieves state-of-the-art performance on the palmprint dataset. Meanwhile, the paper [8] uses the soft-shifted triplet loss function, called SSTL to train the model and matching two features for dealing with translation problem. However, in generally, feature maps of the same class will not only just shift along two axes, but also will have local deformable transformation. For solving it, we propose a new loss and also a new matching method, called translation and rotation triplet loss function (TRTL). With the TRTL, the feature maps can be translated along the x-axis and y-axis, and can be rotated clockwise and counterclockwise. Then we will get the minimal value after translation and rotation as the similarity scores.

3.1 Translated and Rotated Triplet Loss Function

As for a new loss function, the most important point is whether it can be differentiable. With a differentiable loss, the back propagation process can proceed smoothly, and the learnable parameters can be updated to get the minimal loss. In this section, we will discuss the derivation of the TRTL loss function. Because our neural networks were trained using the architecture of triplet network [13], we used TRTL as loss function to update convolutional kernel of our models.

In generally, the TRTLoss is still a variant of triple loss, so that the TRTLoss can be written as a format of triple loss function as the Equation 1. As for the N , it means the batch size during training iteration, and $T(I^a)$ is the output template of input anchor image I^a through neural network. The hard margin parameter m can determine the distance between different class cluster by pushing them away during training process.

$$TRTL = \frac{1}{N} \sum_i^N [L(T(I_i^a), T(I_i^p)) - L(T(I_i^a), T(I_i^n)) + m]_+ \quad (1)$$

In order to adapt to tasks with different degrees of deformation, and balance performance and speed, we set translation and rotation ranges as a hyperparameter. The $L(T_1, T_2)$ will get the minimal distance of two templates $D_{w,h,\theta}(T_1, T_2)$ after translation and rotation in the range $-W \leq w \leq W$, $-H \leq h \leq H$, $-\Theta \leq \theta \leq \Theta$. Meanwhile, the distance $D_{w,h,\theta}(T_1, T_2)$ calculates the pixel-wise MSE value when template T_1 is translated w pixel along x-axis and h pixel along y-axis and rotated θ angle in the Equation 3.

$$L(T_1, T_2) = \min_{-W \leq w \leq W, -H \leq h \leq H, -\Theta \leq \theta \leq \Theta} D_{w,h,\theta}(T_1, T_2) \quad (2)$$

$$D_{w,h,\theta}(T_1, T_2) = \frac{1}{|C_{w,h,\theta}|} \sum_{(x,y) \in C_{w,h,\theta}} (T_1^{(w,h,\theta)}[x, y] - T_2[x, y])^2 \quad (3)$$

In terms of $C_{w,h,\theta}$, it represents the common region between two templates after one template shifted along x-axis with w , shifted along y-axis with h , and rotated with θ . As for the (T_a, T_p) pair, we can assume when the T_a is rotated angle of θ_{ap} and shifted with (w_{ap}, h_{ap}) pixels can get the minimal $D_{w_{ap}, h_{ap}, \theta_{ap}}(T_a, T_p)$, then $L(T_a, T_p) = D_{w_{ap}, h_{ap}, \theta_{ap}}(T_a, T_p)$. Meanwhile, with the $(w_{an}, h_{an}, \theta_{an})$, the (T_a, T_n) pair can get the minimal $D_{w_{an}, h_{an}, \theta_{an}}(T_a, T_n)$.

$$\frac{\partial Loss}{\partial T_i^p} = \begin{cases} 0, \text{ if } (x, y) \notin C_{w_{ap}, h_{ap}, \theta_{ap}} \text{ or } Loss = 0 \\ \frac{-2(T_i^a[[x_{w_{ap}}, y_{h_{ap}}] * M(\theta_{ap})] - T_i^p[x, y])}{N|C_{w_{ap}, h_{ap}, \theta_{ap}}|}, \text{ otherwise} \end{cases} \quad (4)$$

The $M(\theta_{ap})$ is the rotation matrix.

$$\frac{\partial Loss}{\partial T_i^n} = \begin{cases} 0, \text{ if } (x, y) \notin C_{w_{an}, h_{an}, \theta_{an}} \text{ or } Loss = 0 \\ \frac{-2(T_i^a[[x_{w_{an}}, y_{h_{an}}] * M(\theta_{an})] - T_i^n[x, y])}{N|C_{w_{an}, h_{an}, \theta_{an}}|}, \text{ otherwise} \end{cases} \quad (5)$$

As for the $T_i^a[x, y]$ derivation, because we shift and rotate the anchor in the above formula, we can inversely shift and rotate the positive and negative input feature.

$$\frac{\partial Loss}{\partial T_i^a[x, y]} = - \frac{\frac{\partial Loss}{\partial T_i^p[[x - w_{ap}, y - h_{ap}] * M(-\theta_{ap})]} + \frac{\partial Loss}{\partial T_i^n[[x - w_{an}, y - h_{an}] * M(-\theta_{an})]}}{\quad} \quad (6)$$

4 Experiments and Results

For proving TRTL loss function performance, we will compare its performance with Soft-Shift Triplet (STTL)[8] loss function on different public finger knuckle database based on the RFNet [8]. With TRTL loss function, the RFNet is represented by RFNet-TRTL, on the country, RFNet-STTL represents with STTL loss function. Compare to convolution layer or dilated convolution [21], the deformable convolution [22] can solve local deformable by sampling different location and different weight. We also replace the RFNet convolution layer with deformable convolution

layer called DeConvRFNet. As for the RFNet and DeConvRFNet, we will firstly pretrain on the HK PolyU Finger Knuckle Images Database (V1.0) [15] as the pretrained weights.

Meanwhile, we will also compare with the FKNet [2] which get the state-of-the-art performance on 3D finger knuckle identification, and EfficientNetV2-S [14]. Both of FKNet and EfficientNetV2-S are classification neural network. As a classification neural network, it commonly has a problem when the number of classes of testing dataset is not as same as the training set classes, result in fine-tuning on the testing set. Therefore, we use the vector before soft-max layer as the feature vector, and then calculate the MSE of two feature vectors as the similarity score during matching finger knuckle. We use the ResNet-50 pretrained weights as the FKNet initial weights, and use the pretrained weights on the ImageNet21K as the initial weights of EfficientNetV2-S.

In generally, public finger knuckle database already offer segmented finger knuckle images, but we use the

In this section, all experiment will use the finger knuckle segmented by YOLOv5-CSL as the input image to train all models and test the matching performance.

EfficientNetV2-S is the original classification model, we keep the same architecture and just change the FC layer of the head part with convolution layer for fitting TRTL and STTL to compose EfficientNetV2-S-STTL and EfficientNetV2-S-TRTL model. When trained these EfficientNetV2-S model, we use the pretrained weight on the ImageNet21K. As for the FKNet, we use the pretrained ResNet-50 weights.

4.1 Model Complexity Analysis

Model	Params (M)	Input Size	FLOPs (B)	Feature Extraction (s)	Matching (s)
DeConvRFNet-STTL	0.36M	128x128	1.29B		
DeConvRFNet-TRTL	0.36M	128x128	1.29B		
EfficientNetV2-S [14]	20.18M	300x300	5.40B		
EfficientNetV2-S-STTL	20.00M	300x300	5.38B		
EfficientNetV2-S-TRTL	20.00M	300x300	5.38B		
FKNet [2]	7.28M	96x64	0.28B		
RFNet-STTL [8]	0.46M	128x128	1.39B	0.0062s	0.049s
RFNet-TRTL	0.46M	128x128	1.39B	0.0062s	

Table 1: Model complexity analysis

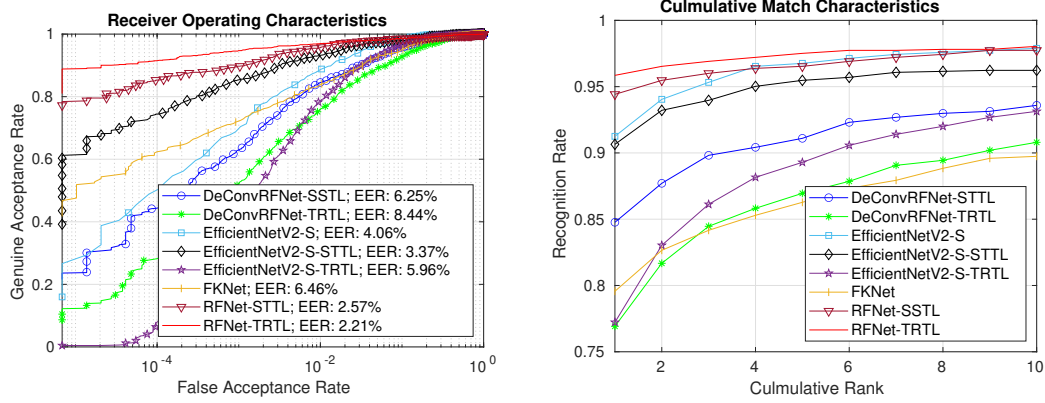
.....

4.2 Within-Database Experiments

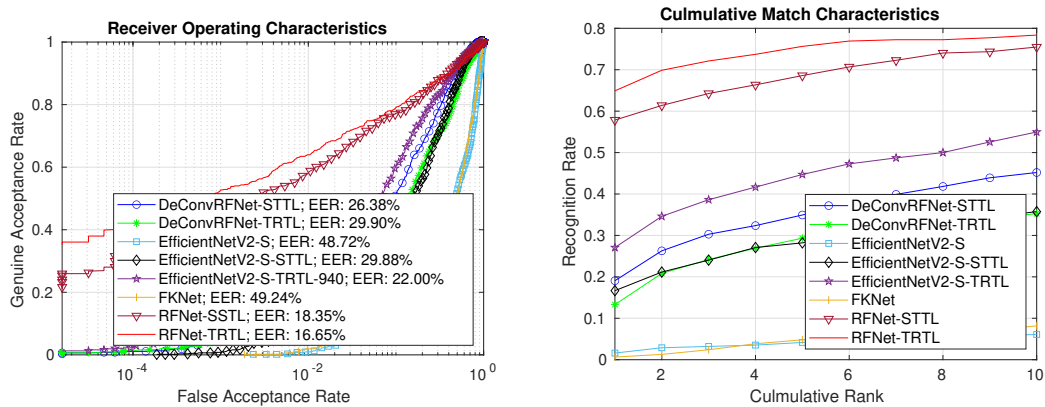
4.2.1 Finger Knuckle V3 Database with Deformation

The Finger Knuckle V3 Database have 1-104 subjects that have two session samples, and the rest subjects of first session 105-221 just offer one session samples. So as the first experiment, I

firstly fine-tuned my model on the second session of 1-104 subjects, and test on the first session 1-104 subjects. So it will have $104 * 6 = 624$ genuine matching scores, and have $104 * 103 * 6 = 64272$ imposter matching scores. From the below figure, if the false accept rate is below 10^{-2} , the RFNet-128-WRS is better than the RFNet-128-WS. I also use the FKNet to train on this database, and the performance of FKNet is not better than the RFNet depend on the ROC figure. From the CMC and ROC, each model with WRS is better than WS on this dataset. For the ROC curve, I add EfficientNetV2-S model performance.



As for the two-session protocol on the database. I should fine-tune my model on the 105-221 subjects, and use two-session protocol to evaluate my model performance on the 1-104 subjects dataset. In totally, it will generate $104 * 6 = 624$ genuine scores, and $104 * 103 * 6$ imposter scores. However, the FKNet is a classification task, and the output number classes should be same when training and testing. So the two session protocol experiment is not fit for FKNet. If the FKNet train on the 105-221 subjects and test on the 1-104 subjects with two sessions, the classes is different.

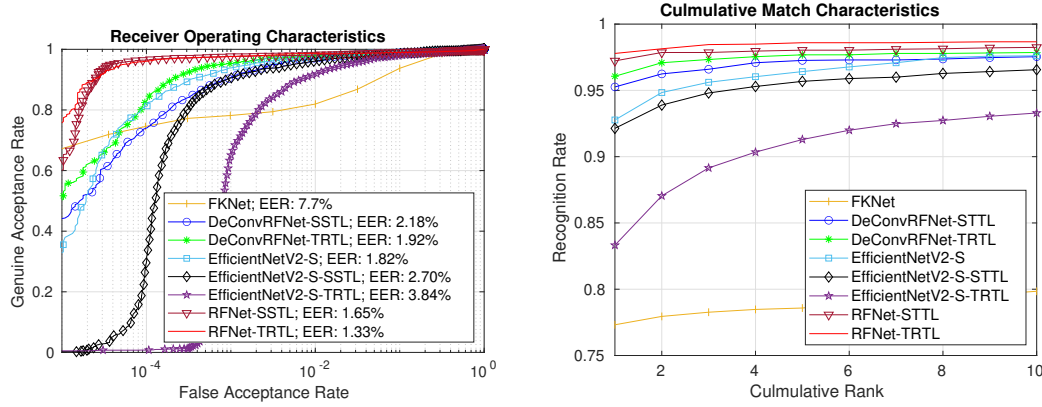


The two-session protocol will use the session1 as the probe and use the session2 as the enrollment. As for the genuine matching scores, each sample of a subject will choose the minimal matching score when compare to the rest samples. In this kind of situation, it will have 104×6 genuine matching scores. Meanwhile, as for the imposter matching scores, it will also choose the minimal value result in $104 * 103 * 6$ imposter matching scores on the confusion matrix.

With Yolov5-CSL segmented finger knuckle, the RFNet performance is slightly higher than the local feature descriptors based on key points matching [6], and performance higher than the paper [7]. If we want to compare different method performance, I think we should use same

dataset. In this kind of situation, the method of [6] maybe will get higher performance on the segmented finger knuckle by YOLOV5.

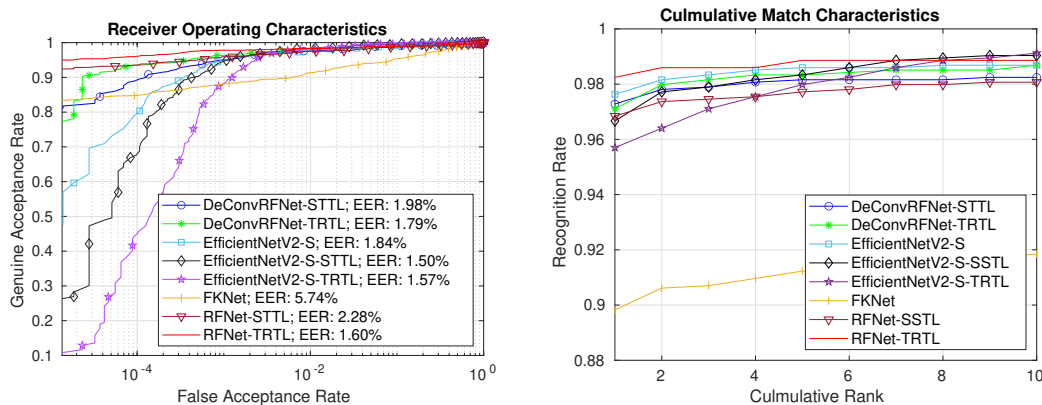
4.2.2 Index Finger Knuckle of Hand Dorsal Image Database



As for the experiment, the dataset totally contains 712 subjects, and I use the segmented Index finger knuckle as my dataset. And I fine-tuned my model on the first sample of each subject, and then use the rest four sample as the testing dataset. At the testing process, it has $712 * 4 = 2848$ genuine matching scores, and has $712 * 711 * 4 = 2024928$ imposter matching scores. The performance of RFN-128-WRS and RFN-128-WS is similar, but the RFN-128-WS is slightly better than RFN-128-WRS depend on the EER value. And we can get an information that the RFNet is better than the rest network in the ROC figure, including the FKNet.

4.2.3 2D Samples of 3D Finger Knuckle Database

First experiment on the database is to use the one session 190 subjects image to fine-tune models and then to test on the another session 190 subjects. It has $190 * 6$ genuine matching scores and $190 * 189 * 6$ imposter matching scores. From the result, we can see that these RFN-128-WRS, RFN-128-WS, EfficientNetV2 can get very high matching accuracy. Meanwhile, the RFNet-TRTL has the minimal EER value among these models. As for the FKNet performance, it gets a very bad result on the 2D images of 3D finger knuckle. I think I have fully trained the FKNet. Maybe the model is overfitting on the training dataset.

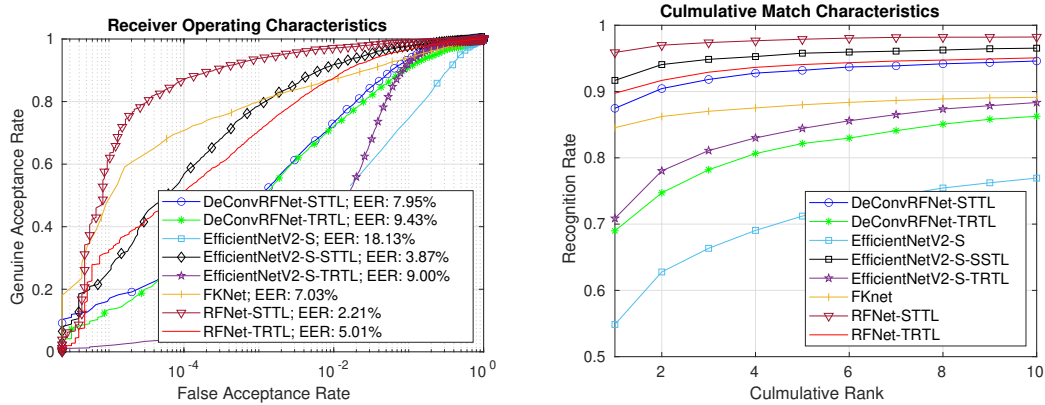


And then use the two session protocol. I use the rest samples of session1, and it has 191-228 subjects. In this kind of situation, the training dataset is too small. The two session protocol will test on the 190 subjects, these subjects can offer two session samples. Due to the training set is too small, so the matching performance is not very good. As for the FKNet, it cannot fit on two session protocol due to classification task.

4.3 Cross-Database Experiments

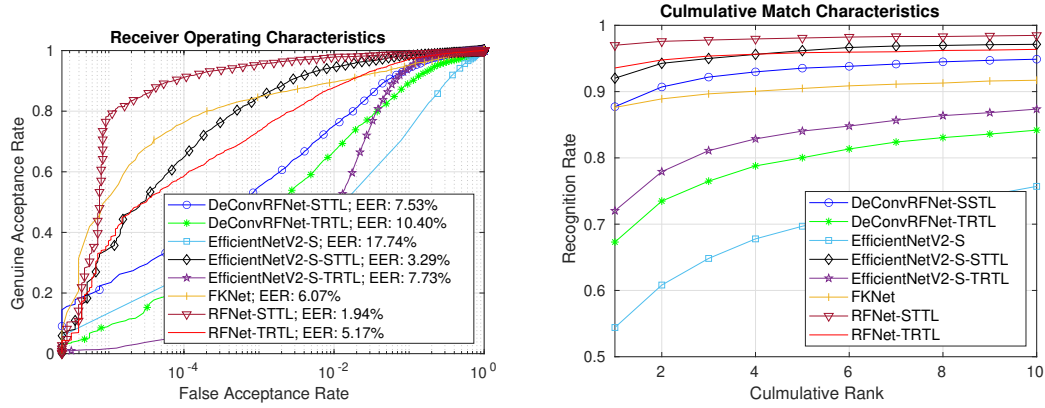
I firstly pre-trained my models on the Finger Knuckle V1 Database, and then fine-tuned models on the Finger Knuckle V3 Database (with deformable). I use these kind training method, and use these models to test performance on the Index Finger Knuckle of Hand Dorsal Image and Tsinghua Finger Knuckle Database as a cross database experiment. The label in the finger curve, the content in parentheses indicates the training samples. Such as RFN-WS(1-104), it uses 1-104 subjects of Finger Knuckle V3 Database to train models.

4.3.1 Index Finger Knuckle of Hand Dorsal Image

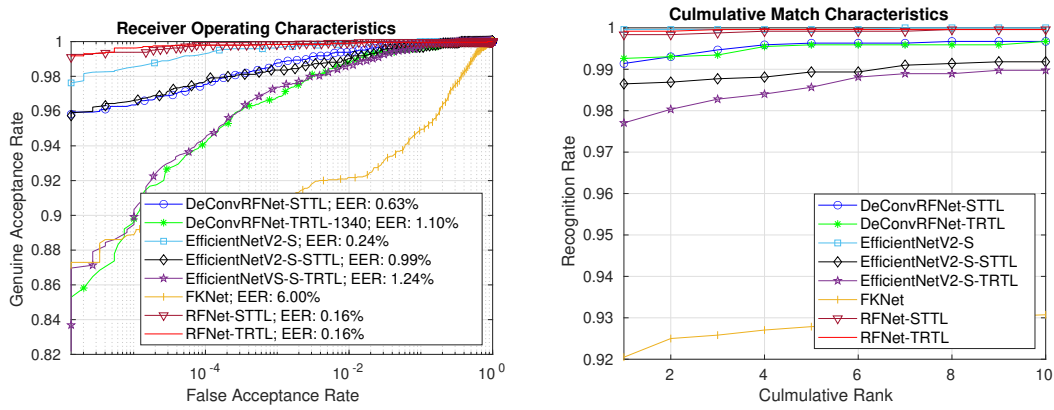


The database totally has 712 subjects, and each subject has 5 samples. Therefore, it will have $712 * 5$ genuine matching scores and $712 * 711 * 5$ imposter matching scores. From the curve, the performance of RFNet-SSTL and RFNet-TRTL is similar, and the RFNet-SSTL is slightly better than RFNet-TRTL while using the same training samples.

4.3.2 Middle Finger Knuckle of Hand Dorsal Image



4.3.3 Tsinghua Finger Knuckle Database



The database has 610 subjects, and each subject can offer 4 samples. Then as the cross database experiment, it will have $610 * 4$ genuine matching scores and $610 * 609 * 4$ imposter matching scores. In this database, all models can get very high matching performance from the table and figure.

4.4 Discussion

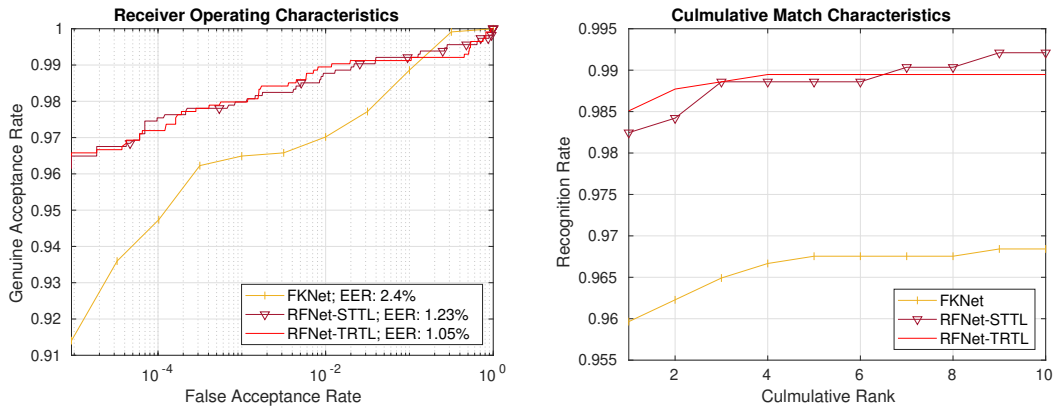
From these experiment results, we can see that EfficientV2-S model is better than FKNet in some dataset. Because EfficientNetV2 model use MBConv as a block unit for replacing residual block. As for MBConv block, it uses depth-wise convolutions to decrease training weights and use Squeeze-Excited block as channel transformer. Meanwhile, the depth of EfficientNetV2-S is deeper than the FKNet.

There is another conclusion is that TRTL generalization ability is lower than STTL loss from the cross database experiment. But in the within database experiment, these model with TRTL loss is better than STTL loss.

.....

5 3D Finger Knuckle of 3D Finger Knuckle Database

I have use the matlab code that offered by the FKNet to generate the 3D finger knuckle images for getting the depth information. But it is different that the input image size. The FKNet will resize the original image size $148 * 212$ to $70 * 100$ as the testing dataset, and crop from the $70 * 100$ to $48 * 80$ as the training dataset. As for RFNet, I just use the original image as the input data. Then the experiment protocol will generate $190 * 6$ genuine matching scores, and $190 * 189 * 6$ imposter matching scores. From the experiment result, we can get that the RFNet is the best model for the 3D Finger Knuckle Database.



6 Ablation Study

7 Online Contactless Finger Knuckle Identification

With TRTL loss, the RFNet [8] can outperform state-of-the-art methods. In the previous section, we have estimated its verification and identification performance on different public finger knuckle database, including within-db and cross-db experiments. As for a completely contactless and online finger knuckle identification, the finger knuckle detector is a very important module for automatically detect and segment finger knuckle region. However, as for traditional segmentation algorithm, they cannot correctly segment the finger knuckles in the presence of complex background interference, multiple finger knuckles in the same field of view, obscured finger knuckles or bent finger knuckles. Meanwhile, as for neural network, the current based on YOLO [11], [9], [10], [1], [17] and R-CNN [4], [3], [12], [5] series object detection and segmentation approaches cannot simultaneously obtain the angle of finger knuckle and the segmentation with high precision. Especially, the angle of the finger knuckle is a vital factor for identification. If we can get the angle of finger knuckle, we can use angle information to align two feature maps for increasing matching accuracy and efficiency. For solving above problems, we propose rotated bounding box detection based on YOLOv5 model for segmenting and getting angle information.

7.1 Contactless Finger Knuckle Detection

7.1.1 Rotated Bounding Box Based on YOLOv5

In order to solve the problem of finger knuckle detection in the real world, we choose to use YOLOv5 model because the YOLO series is famous for its fast detection speed and high accuracy. Especially, the YOLOv5's [17] speed can meet our online detection requirements.

Rotated Bounding Box

However, the YOLOv5 just detect horizontal bounding boxes which cannot offer angle information and will segment a lot of background information. In order to solve these above problem, a rotated bounding box will be predicted instead of horizontal bounding box. As analyzed in this paper [20], the rotated bounding boxes loss will mainly come from angular periodicity and the exchangeability of edges. When use the long side definition of rotated bounding box, it can deal with the exchangeability of edges problem. Meanwhile, using classification task to predict angle can make model easier to train. A periodic coding method called Circular Smooth Label (CSL) [20] soft coding can also solve the problem that One-Hot cannot distinguish class relationship. Formula 7 $g(x)$ is the window function to smooth One-Hot label, and r is a window function of the radius.

$$CSL(x) = \begin{cases} g(x), & \theta - r < x < r + \theta \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

Furthermore, in this paper, we used the Gaussian function for the Equation 7 window function, a commonly available function, and used a window radius of 6 to smooth the labels.

Loss function

The original YOLOv5 loss function can have three components. The formula can be simply written as $Loss = CIOU_Loss + Loss_{conf} + Loss_{class}$. Since the rotated bounding box is based on the modification of YOLOv5, only the angle classification loss is added more. So the total loss function is as expressed in Equation 8, with the addition of $Loss_{angle}$ to YOLOv5 loss function.

$$Loss = CIOU_Loss + Loss_{conf} + Loss_{class} + Loss_{angle} \quad (8)$$

$$Loss_{angle} = \sum_{i=0}^{S^2} I_{ij}^{obj} \sum_{a \in [0, 180)} [P_i(\hat{a}) \log(P_i(a)) + (1 - P_i(\hat{a})) \log(1 - P_i(a))] \quad (9)$$

7.1.2 Contactless Finger Knuckle Dataset

Our task is to detect finger knuckles in the contactless and online scenario, but by understanding current public finger knuckle database, their data are collected at specific conditions such as certain angle, certain light. In this kind of situation, this kind of data cannot represent real images of finger knuckle in real world. In order to address the shortcomings of current public finger knuckle dataset for contactless detection, we use a web crawler to get images from the

Unsplash [18] where the keywords are finger knuckles. The Unsplash is an image site that offers uploads and downloads, and uses a copyright license that allows users to download and use them for free or even for commercial use [19]. We have downloaded 2347 images, there are 738 images without knuckles, and these images can be used as background training, and the rest 1609 images that contain at least one finger knuckle are the positive samples for the network model. In the network training process, we use crawled images, 169 finger knuckle images from the HKPolyU Finger Knuckle Database (V1.0) [15], and 64 finger knuckles images from the HKPolyU Hand Dorsal Database [16] as for the training set. And we use the rest data as testing set to evaluate performance. The most important part is the data augmentation which contains flip, rotation, resize, translate and mosaic.

7.1.3 Contactless Finger Knuckle Detection

Detection Performance

The YOLOv4 model predict horizontal bounding box, while the remaining YOLOv5 model predict rotated bounding box with CSL classification, called YOLOv5-CSL. We can see the performance difference between these variations of the YOLOv5 model from the Table 2. Among the downloaded 2580 images, 100 images were randomly selected as the testing set.

Model	Inference Time/ms (1024x1024)	Number of Layers	mAP^{val} 0.5	AP of Major FK	AP of Minor FK
YOLOv5x-CSL	41.395	407	89.9	89.6	90.1
YOLOv5m-CSL	36.252	263	85.7	88.9	80.4
YOLOv4	25.992	161	70.7	83.6	57.7

Table 2: Comparison of the accuracy of the different models of the YOLO series for the detection of the finger knuckle. The calculated values of mAP were measured at a detection threshold of 0.4 as well as an IOU threshold of 0.5.

Segmentation Performance

This section aim to compare quality of finger knuckle between YOLOv5-CSL segmented and dataset offered. Because the segmented finger knuckle on the 3D Finger Knuckle Dataset already have high quality, I mainly test on the Index Finger Knuckle of Hand Dorsal Dataset and the Finger Knuckle Dataset V3 (with deformable).

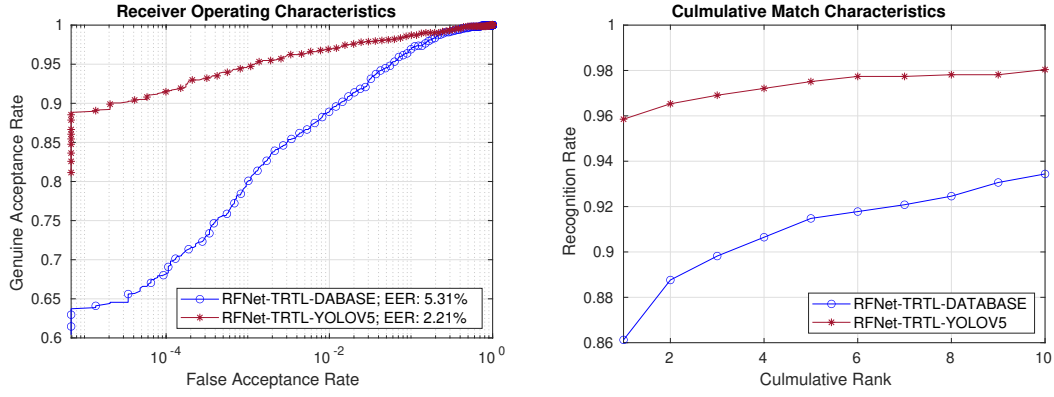


Figure 9: Compare performance on the Finger Knuckle V3 Dataset (with deformable)

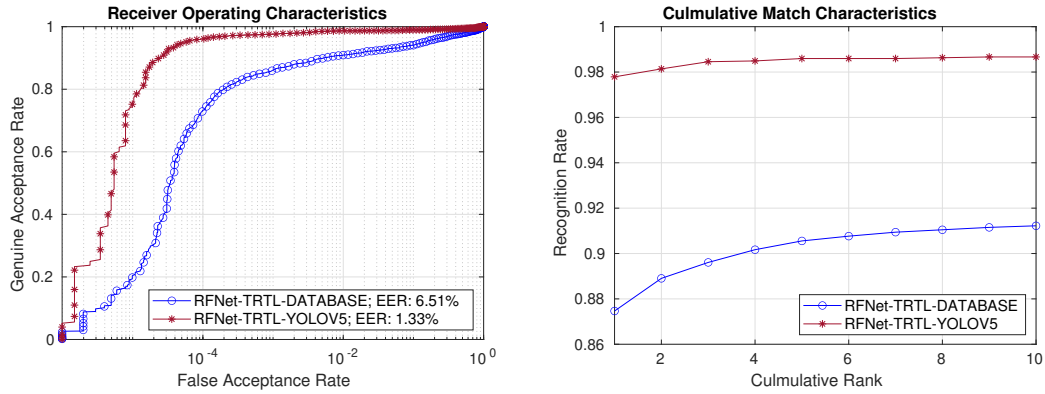


Figure 10: Compare performance on the Index Finger of the Hand Dorsal Image Database.

From the above figures, we can clearly get the conclusion that quality of segmented finger knuckle of YOLOv5 is better than the segmented finger knuckle of dataset through the ROC curve and CMC curve. Especially on the Hand Dorsal Image Database, the EER value can drop from 6.51% to 1.33%.

7.2 Online Contactless Finger Knuckle Identification Performance

.....

8 References

- [1] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. "Yolov4: Optimal speed and accuracy of object detection". In: *arXiv preprint arXiv:2004.10934* (2020).
- [2] Kevin HM Cheng and Ajay Kumar. "Deep feature collaboration for challenging 3D finger knuckle identification". In: *IEEE Transactions on Information Forensics and Security* 16 (2020), pp. 1158–1173.
- [3] Ross Girshick. "Fast r-cnn". In: *Proceedings of the IEEE international conference on computer vision*. 2015, pp. 1440–1448.

- [4] Ross Girshick et al. “Rich feature hierarchies for accurate object detection and semantic segmentation”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014, pp. 580–587.
- [5] Kaiming He et al. “Mask r-cnn”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2961–2969.
- [6] Ajay Kumar. “Contactless finger knuckle authentication under severe pose deformations”. In: *2020 8th International Workshop on Biometrics and Forensics (IWBF)*. IEEE. 2020, pp. 1–6.
- [7] Ajay Kumar. “Toward pose invariant and completely contactless finger knuckle recognition”. In: *IEEE Transactions on Biometrics, Behavior, and Identity Science* 1.3 (2019), pp. 201–209.
- [8] Yang Liu and Ajay Kumar. “Contactless palmprint identification using deeply learned residual features”. In: *IEEE Transactions on Biometrics, Behavior, and Identity Science* 2.2 (2020), pp. 172–181.
- [9] Joseph Redmon and Ali Farhadi. “YOLO9000: better, faster, stronger”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 7263–7271.
- [10] Joseph Redmon and Ali Farhadi. “Yolov3: An incremental improvement”. In: *arXiv preprint arXiv:1804.02767* (2018).
- [11] Joseph Redmon et al. “You only look once: Unified, real-time object detection”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 779–788.
- [12] Shaoqing Ren et al. “Faster r-cnn: Towards real-time object detection with region proposal networks”. In: *Advances in neural information processing systems* 28 (2015), pp. 91–99.
- [13] Florian Schroff, Dmitry Kalenichenko, and James Philbin. “Facenet: A unified embedding for face recognition and clustering”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 815–823.
- [14] Mingxing Tan and Quoc Le. “Efficientnetv2: Smaller models and faster training”. In: *International Conference on Machine Learning*. PMLR. 2021, pp. 10096–10106.
- [15] *The HKPolyU Contactless Finger Knuckle Images Database (V-1.0)*: <http://www4.comp.polyu.edu.hk/~csajaykr/fn1.htm>.
- [16] *The HKPolyU Contactless Hand Dorsal Images Database*: <http://www4.comp.polyu.edu.hk/~csajaykr/knuckleV2.htm>.
- [17] Ultralytics. *YOLOv5*. <https://github.com/ultralytics/yolov5>. 18 May 2020.
- [18] Unsplash. *Unsplash*. <https://unsplash.com/>.
- [19] Unsplash. *Unsplash.com*. ”Unsplash License”. Retrieved 11 January 2017.
- [20] Xue Yang and Junchi Yan. “Arbitrary-oriented object detection with circular smooth label”. In: *European Conference on Computer Vision*. Springer. 2020, pp. 677–694.
- [21] Fisher Yu and Vladlen Koltun. “Multi-scale context aggregation by dilated convolutions”. In: *arXiv preprint arXiv:1511.07122* (2015).
- [22] Xizhou Zhu et al. “Deformable convnets v2: More deformable, better results”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019, pp. 9308–9316.