# Improved YOLOv4-tiny based on attention mechanism for skin detection

**Ping Li** [1,2] , **Taiyu Han** [1] , **Yifei Ren** [1] , **Peng Xu** [1] , **Hongliu Yu** [Corresp. 1]

[1] Institute of Rehabilitation Engineering and Technology, University of Shanghai for Science and Technology, Shanghai, China

[2] Department of Biomedical Engineering, Changzhi Medical College, Changzhi, Shanxi, China

Corresponding Author: Hongliu Yu
Email address: yhl_usst@outlook.com

**Background.** The automatic bathing robot needs to identify the area to be bathed to perform visually guided bathing tasks. The visual perception of the skin area is the first step in the operation of an automatic bathing robot. The deep CNN-based object detection algorithm has excellent robustness to light and environmental changes when performing skin detection. The one-stage object detection algorithm has good real-time performance, which is widely used in practical projects. **Methods.** In our previous work, we perform skin detection using several models and find that YOLOv4 has best comprehensive performance. This study uses the YOLOv4-tiny model for skin detection considering the convenience of practical deployment. In addition, we add three kinds of attention mechanisms to strengthen feature extraction, namely SE, ECA, and CBAM. In particular, we add the attention module to the two feature layers of the backbone output. In the enhanced feature extraction network part, we apply the attention module to the upsampled features. **Results.** The experimental results reveal that the weight file of YOLOv4-tiny without attention mechanisms is reduced to 9.2% of YOLOv4, but the mAP maintains 67.3% of YOLOv4. The performance of the YOLOv4-tiny is improved by combining the CBAM or ECA modules, but the addition of SE deteriorates the performance of YOLOv4-tiny instead. Among them, CBAM is the best, which can enhance YOLOv4-tiny detection accuracy by about 5%.

# Improved YOLOv4-tiny based on attention mechanism for skin detection

Ping Li[1,2], Taiyu Han[1], Yifei Ren[1], Peng Xu[1], Hongliu Yu[1]

[1] Institute of Rehabilitation Engineering and Technology, University of Shanghai for Science and Technology, Shanghai, China

[2] Department of Biomedical Engineering, Changzhi Medical College, Changzhi, Shanxi, China

Corresponding Author:

Hongliu Yu[1]

Email address: yhl_usst@outlook.com

10 # Improved YOLOv4-tiny based on attention mechanism

11 # for skin detection

12 Ping Li[1,2], Taiyu Han[1], Yifei Ren[1], Peng Xu[1], Hongliu Yu[1]

13 [1] Institute of Rehabilitation Engineering and Technology, University of Shanghai for Science and

14 Technology, Shanghai, China

15 [2] Department of Biomedical Engineering, Changzhi Medical College, Changzhi, Shanxi, China

16 Corresponding Author:

17 Hongliu Yu[1]

18 No. 580, Jungong Road, Yangpu District, Shanghai, 200093, China

19 Email address: yhl_usst@outlook.com

20

## 21 Abstract

22 **Background.** The automatic bathing robot needs to identify the area to be bathed to perform

23 visually guided bathing tasks. The visual perception of the skin area is the first step in the

24 operation of an automatic bathing robot. The deep CNN-based object detection algorithm has

25 excellent robustness to light and environmental changes when performing skin detection. The

26 one-stage object detection algorithm has good real-time performance, which is widely used in

27 practical projects.

28 **Methods.** In our previous work, we perform skin detection using several models and find that

29 YOLOv4 has best comprehensive performance. This study uses the YOLOv4-tiny model for skin

30 detection considering the convenience of practical deployment. In addition, we add three kinds

31 of attention mechanisms to strengthen feature extraction, namely SE, ECA, and CBAM. In

32 particular, we add the attention module to the two feature layers of the backbone output. In the

33 enhanced feature extraction network part, we apply the attention module to the upsampled

34 features.

35 **Results.** The experimental results reveal that the weight file of YOLOv4-tiny without attention

36 mechanisms is reduced to 9.2% of YOLOv4, but the mAP maintains 67.3% of YOLOv4. The

37 performance of the YOLOv4-tiny is improved by combining the CBAM or ECA modules, but

38 the addition of SE deteriorates the performance of YOLOv4-tiny instead. Among them, CBAM

39 is the best, which can enhance YOLOv4-tiny detection accuracy by about 5%.

## 40 Introduction

41 CNN (Convolutional Neural Network) is a machine learning model in a supervised learning

42 framework. In 2012, AlexNet first used CNN for image classification (Krizhevsky, Sutskever &

43 Hinton, 2017), winning the ImageNet large scale visual recognition challenge by an

44 overwhelming margin. Since then, CNN has been widely used in computer vision tasks such as

45 image classification (Liu, Soh & Lorang, 2021) and object detection (Zhou et al., 2022). Using

46 massive data as learning samples, we can obtain a CNN model with analysis capability, feature

47 representation capability, and recognition capability to achieve skin detection. In CNN models,

48 the convolutional layer extracts features, the pooling layer performs dimensionality reduction

- -

49   and information integration, and the fully connected layer combines the extracted features and
50   outputs data adapted to a specific problem. In general, an activation function is introduced to
51   give the model a nonlinear representation capability. With the compression of the width and
52   height of feature maps, CNN models acquire robust semantic information and abstract extracted
53   features.
54      Skin perception for automatic bathing robots is a prerequisite for bathing. The intelligent
55   bathing system detects the human skin in the bathing scene based on vision sensors. Skin
56   detection in bathing scenes is a challenging task. From the environmental perspective, the
57   bathing scene is full of water mist and own various lighting and backgrounds. A skin detection
58   algorithm generally extracts skin features and then classifies them using a classifier. Traditional
59   skin detection typically exploits handcrafted features to distinguish between skin and non-skin
60   zones, such as color, texture, and statistical features. Zhang proposes a skin color model based on
61   reference points of the face (Zhang et al., 2022). Shifa conducts skin detection by combined
62   threshold rule-based segmentation in the RGB, HSV, and YCbCr color spaces (Shifa et al., 2020).
63   Sun uses a local skin color model to change a global model performing skin detection for single
64   images (Sun, 2010). Luo performs skin detection by face location method and facial structure
65   estimation (Luo & Guan, 2017). Javadi conducts the skin lesions detection by color properties
66   using a genetic algorithm for selecting the best features and determines the 3D position by
67   Kinect camera (Javadi & Soltanizadeh, 2021). Handcrafted features are not robust to
68   environmental changes and are insufficient for bathing scenarios. Skin detection based on
69   machine learning, which generally uses supervised methods to construct detectors to extract skin
70   features, is less influenced by environmental factors and has gained more applications in recent
71   years. Salah utilizes CNN trained by the skin and non-skin patches to detect skin pixels (Salah,
72   Othmani & Kherallah, 2022). Kim exploits two CNNs for skin detection and compares
73   performance using different training strategies (Kim, Hwang & Cho, 2017). Lin conducts the
74   CNN-based facial skin detection and optimizes the CNN with the Taguchi method (Lin et al.,
75   2021).
76      Different from just identifying skin and non-skin areas, we need to provide the robot with
77   information about specific skin areas (hands, feet, head, etc.) to clean up skins using different
78   modes. We are facing a multi-classification problem rather than a secondary classification
79   problem. For the bathing scenario, traditional algorithms become inadequate. Therefore, deep
80   learning methods are introduced. Skin detection based on deep CNN does not rely on
81   handcrafted features. In application areas, one-stage object detection models based on CNN
82   achieve good real-time performance and are computationally efficient, such as the YOLO
83   models (Redmon et al., 2016; Redmon & Farhadi, 2017; Redmon & Farhadi, 2018; Bochkovskiy,
84   Wang & Liao, 2020). The one-stage framework eliminates the proposals generation and outputs
85   the categories and bounding boxes directly.
86      We perform skin detection based on deep learning methods. Our research is based on
87   previous work by our team (Li et al., 2022), which finds that the YOLOv4 algorithm has a high
88   mAP for skin detection in bath scenes. At the same time, it has extensive computation, leading to

89    the slow speed of YOLOv4 after being deployed to embedded devices. In the study, we adopt

90    YOLOv4-tiny (Zhao et al., 2022a) for skin detection, the lightweight model of YOLOv4, and

91    investigate the effect of attention mechanisms on the YOLOv4-tiny.

92          The remaining parts of the paper are arranged as follows: "Materials & Methods" section

93    offers an introduction to data sets acquisition, YOLOv4-tiny, improved YOLOv4-tiny based on

94    attention mechanisms, transfer learning, experimental setup, and evaluation indicators. The

95    "Results" section describes the experimental results. The "Discussion" section discusses the

96    results related to our application scenarios. The "Conclusion" section summarizes our research

97    and looks at future work.

## Materials & Methods

99    **Data sets acquisition**

100   A total of 1500 images containing human skin are collected, considering factors such as position,

101   illumination, resolution, blurring, and the presence of water mist. Finally, our data sets choose

102   1000 images based on the image quality. The image annotation tool LabelImg (Bhatt et al., 2022)

103   is used to generate XML files corresponding to the images. The XML file includes the file name,

104   ground truth box information, and category information in the corresponding image. Example

105   images in the data set are exhibited in Fig. 1.

106   **YOLOv4-tiny**

107   The structure of YOLOv4-tiny is shown in Fig. 2. The backbone is CSPDarknet53-tiny, which is

108   utilized for feature extraction. CSPDarknet53-tiny is composed of DarknetConv2D_BN_Leaky

109   modules and Resblock_body modules. A DarknetConv2D_BN_Leaky module combines a two-

110   dimensional convolutional layer, a normalized processing layer, and an activation function. The

111   Mish activation function (Misra, 2019) in the YOLOv4 is replaced by a Leaky Relu function (He

112   et al., 2015) to improve detection speed. The structure of Resblock_body is illustrated in Fig. 3.

113   The skip connection can better combine the semantic information and let the model converge

114   quickly, preventing both model degradation and gradient disappearance (Furusho & Ikeda, 2020).

115   Feat1 and Feat2 are the output feature layers from the Resblock_body module. The Feat2 output

116   branch of the first two Resblock_body modules is the input of the next module. FPN (Lin et al.,

117   2017) is used for enhancing feature extraction and performing feature fusion to combine feature

118   information at different scales. For the output of the third Resblock_body module, Feat1 is

119   directly used as the first input of the FPN. The second input of the FPN is the result obtained by

120   processing Feat2 by the DarknetConv2D_BN_Leaky module. The output P2 of FPN is obtained

121   by convolution processing on the second input of the FPN. The output P1 of FPN is obtained by

122   stacking Feat1 and the result which is obtained by convolution and up-sampling operations on P2.

123   The structure of FPN is simple, allowing YOLOv4-tiny to have excellent real-time performance.

124   Compared with YOLOv4, YOLOv4-tiny has two detection heads and predicts at two scales. The

125   YOLO head is used to obtain classification and regression prediction results. The structure of the

126   YOLO head is straightforward. The two prediction feature layers for prediction are acquired by a

127   small amount of convolution of P1 and P2. YOLOv4-tiny is still making the detection based on

128   anchors, using fixed-size anchors as a prior for object boxes, tiling many anchors on images, and

129 adjusting anchors to bounding boxes by the prediction results. "13×13" and "26×26" represent
130 the granularity of grids. "33" represents the prediction results adapted to our application, i.e.,
131 3×(4+1+6), where "3" represents the number of anchors, "4" indicates the number of location
132 parameters, "1" denotes the confidence score, and "6" means the number of categories to be
133 identified.

134     The loss function generally includes bounding box location loss $L_{loc}$, classification loss $L_{cls}$,
135 and confidence loss $L_{conf}$. The overall loss L is calculated as Eq. (1).

136
$$L = L_{loc} + L_{cls} + L_{conf} \quad (1)$$

137     $L_{loc}$ measures the position error (height $h$, width $w$, and central coordinates) between the
138 prediction box and the GT box. The evaluation indicators include IOU, GIOU (Rezatofighi et al.,
139 2019), DIOU, and CIOU (Zheng et al., 2019), as summarized in Table 1. We introduce CIOU
140 loss as $L_{loc}$, as indicated in Eq. (2).

141
$$L_{loc} = 1 - IoU + \rho^2\left(b, b^{gt}\right)\big/d^2 + \alpha\upsilon \quad (2)$$

142
$$\upsilon = 4\big/\left(\pi^2\right) * \left(\arctan\left(w^{gt}/h^{gt}\right) - \arctan\left(w/h\right)\right)^2 \quad (3)$$

143
$$\alpha = \upsilon\big/\left(\upsilon + 1 - IoU\right) \quad (4)$$

144 $\rho^2(b, b^{gt})$ represents the European distance between the central points of the prediction box and
145 the GT box, $d$ represents the diagonal distance of the minimum enclosed area, including the
146 prediction box and the GT box, $\alpha$ is weight, and $\upsilon$ expresses the consistency of aspect ratio. $\alpha$
147 and $\upsilon$ are calculated as demonstrated in Eq. (3) and Eq. (4).

148     $L_{cls}$ measures the category error between the prediction box and the GT box, as shown in Eq.
149 (5).

150
$$L_{cls} = -\sum_{i=0}^{K\times K} I_{ij}^{obj} \sum_{c\in classes} \left[ \hat{p}_i(c)\log\left(p_i(c)\right) + \left(1 - \hat{p}_i(c)\right)\log\left(1 - p_i(c)\right) \right] \quad (5)$$

151 $K\times K$ represents the number of grids on feature maps of different scales, and $c$ represents the
152 category. If the $j$-th prior box of the $i$-th grid has objects to be predicted, $I_{ij}^{obj}=1$; otherwise,
153 $I_{ij}^{obj}=0$.   $(c)$ and $p_i(c)$ represent the actual value and predicted value of the probability that the $j$-
154 th prior box of the $i$-th grid belongs to category $c$, respectively.

155     $L_{conf}$ adopts a cross-entropy loss function, as shown in Eq. (6). $M$ represents the number of
156 prior boxes.   and $C_i$ represent the actual and predicted values of confidence. If the $j$-th prior box
157 of the $i$-th grid has no object to be predicted, $I_{ij}^{noobj}=1$; otherwise, $I_{ij}^{noobj}=0$.

158
$$L_{conf} = \sum_{i=0}^{K\times K}\sum_{j=0}^{M} I_{ij}^{obj}\left[\hat{C}_i \log\left(C_i\right) + \left(1 - \hat{C}_i\right)\log\left(1 - C_i\right)\right] - \sum_{i=0}^{K\times K}\sum_{j=0}^{M} I_{ij}^{noobj}\left[\hat{C}_i \log\left(C_i\right) + \left(1 - \hat{C}_i\right)\log\left(1 - C_i\right)\right] \quad (6)$$

159 **Improved YOLOv4-tiny based on attention mechanisms**
160 The attention mechanism is a normal tip for deep learning, which has a variety of
161 implementations (Niu, Zhong & Yu, 2022). The core of the attention mechanism is to make the
162 network pay attention to where it needs more attention. In general, attention mechanisms can be

163  divided into the channel attention mechanism, the spatial attention mechanism, and a
164  combination of the two (Tian et al., 2021).
165     In this paper, the following attention mechanisms are used:
166     (1) SE (Squeeze-and-Excitation) (Hu, Shen & Sun, 2018). SE is a typical implementation of
167  the channel attention mechanism, obtaining the weights of each channel in the feature maps. The
168  inter-dependencies among channels are modeled explicitly. Instead of introducing a new-built
169  spatial dimension for the fusion of feature channels, SE uses a feature rescaling strategy.
170  Specifically, the importance of each channel is acquired spontaneously by self-learning.
171  Following the degree of matter, the helpful features are enhanced, and the useless features are
172  suppressed, achieving the adaptive calibration of feature channels. SE includes squeeze and
173  excitation operations. The squeeze operation conducts feature compression across the spatial
174  dimension, converting a two-dimensional feature map into a real number that owns a global
175  receptive field. The output size matches the number of input channels. The excitation operation
176  is equivalent to the mechanics of gates in recurrent neural networks, where weights are created
177  for each channel employing learned parameters, explicitly modeling the correlation between
178  feature channels. Finally, the weights, which are output by excitation operation, represent
179  importance of each channel. The rescaling of features in the channel dimension is accomplished
180  by multiplying the weights by features of each channel (Huang et al., 2019). The specific
181  implementation of SE is shown in Fig. 4.
182     (2) ECA. ECA is an improved version of SE. Wang argues that seizing all channel
183  dependencies is ineffective and unessential for SE block (Wang et al., 2020). Convolution
184  operation owns the cross-channel information capture capability. ECA removes the fully
185  connected layer of SE and learns weights by 1D convolution operation on the globally averaged
186  pooled features. The specific implementation of ECA is shown in Fig. 5.
187     (3) CBAM (Convolutional Block Attention Module). CBAM (Woo et al., 2018) performs
188  channel attention and spatial attention mechanism processing for feature maps, respectively. The
189  specific implementation of CBAM is shown in Fig. 6. The implementation of the channel
190  attention module (CAM) can be divided into two parts. Global average pooling and maximum
191  global pooling are performed separately for the input feature maps. The outputs are processed
192  using a shared, fully connected layer. Sum the two processed results, and then take the sigmoid
193  operation, obtaining the weights (between 0 and 1) of each channel of the input features. The
194  weights are multiplied by the original input features to get the output of CAM. The spatial
195  attention module (SAM) takes the maximum value and the average value on each channel of
196  each feature point. The two results are stacked. Adjust the number of channels using a
197  convolution operation. Get the weights of each feature point of the input features through the
198  following sigmoid function. Obtain the final output by multiplying the weights by the original
199  input features.
200     In this study, the above attention mechanisms are applied to the YOLOV4-tiny. As shown
201  in Fig. 7, we add attention mechanisms on the two feature layers extracted from the backbone
202  network and attention mechanisms on the up-sampled results in FPN.

**Transfer Learning**

Training a network from scratch requires a enormous amount of labeled data. Manual labeling of data sets is time-consuming and labor-intensive, which introduces human error. Small data sets combined with transfer learning techniques can produce a desirable model quickly (Pratondo & Bramantoro, 2022). The ImageNet contains more than 14 million images covering more than 20,000 categories, of which more than one million images have explicit annotations and corresponding labels at objects' locations in the image (Russakovsky et al., 2015). The pre-trained models on ImageNet can learn fundamental features such as textures, lines, etc., which are general in object detection. The pre-trained weights on the ImageNet are the initial weights for all models in this study.

**Experimental setup and evaluation indicators**

For the fairness of model comparison, we use the same data sets as our previous work (Li et al., 2022), with a ratio of 60%:20%:20% for the training, validation, and test sets. All models are trained with the help of the high-performance computing center of the University of Shanghai for Science and Technology. Mosaic data augmentation is used in the training process in which four randomly stitched images are input to the network for training to increase the background diversity (Bin et al., 2022). The learning rate cosine decline strategy is used during the model training. The loss function is optimized using a label smoothing approach to suppress the overfitting problem during training (Zhang et al., 2021). The probability $p_i(c)$ distribution before and after label smoothing is shown in Eq. (7), with $\delta = 0.05$.

$$p_i(c) = \begin{cases} 1, i = y \\ 0, i \neq y \end{cases} \Rightarrow p_i(c) = \begin{cases} 1-\delta, i = y \\ \delta, i \neq y \end{cases} \quad (7)$$

We use the Pytorch framework for model building and training. The initial value of the learning rate is set to 0.001, and the decay rate is set to 0.01. The batch size is set to 16, which indicates the number of images input to the model for training every time. SGD is utilized as the optimizer for model training. When training, the weights of the backbone are frozen first for 50 epochs, and all weights are trained after 50 epochs, which increases the convergence speed and training performance of models.

Recall and precision can be used to measure performance but are not fully representative of the detector quality. Many sets of recall and precision values are obtained by taking different thresholds. Then, plot a P-R curve (Naing et al., 2022). AP characterizes the area enclosed by the P-R curve and the coordinate axes. The sum of AP values of all classes is then divided by the total number of classes to get mAP, which is the crucial evaluation metric of detectors for multiple categories detection.

# Results

After the training is completed, models are selected based on the results of the validation sets, and the performance is tested using the test sets. The mAPs and weight file information of models are exhibited in Table 2.

In our previous work, the mAP of YOLOv4 reaches 78%, but it has a weight file of 244 MB. After the light-weighting process, the mAP of YOLOv4-tiny is 67.3% of YOLOv4, but the

242  weight file is reduced to 9.2% of YOLOv4. Based on the YOLOv4-tiny, we add attention
243  mechanisms as shown in Fig. 7. As can be seen from Table 2, the detection result decreases
244  instead, and mAP is reduced by 0.9% after adding the SE. There is a 1.1% improvement in mAP
245  after adding the ECA. The performance improvement is the highest with the addition of CBAM,
246  in which mAP is increased by nearly 5%. After adding ECA, the weight file has hardly increased.
247  After adding SE, the weight file has increased by 0.2 M. After adding CBAM, the weight file has
248  increased by 0.4 M. We have established a comprehensive indicator $W$, as shown in Eq. (8). $A$
249  indicates the change of weight file, and $B$ indicates the shift in mAP. When mAP is less than the
250  mAP of the original YOLOv4-tiny model, $B$ takes a negative value. Otherwise, $B$ takes a positive
251  value. The smaller the $A$, the better the effect. The higher the $B$, the better the performance.
252  Overall, the higher the $W$, the better the outcome. After calculation, $A$, $B$, and $W$ are indicated in
253  Table 3. CBAM_YOLOv4-tiny, which introduces the CBAM modules, achieves the best
254  outcome.

$$W = \frac{B}{e^A}$$

255                                                                                       (8)

256        The AP values for the six categories are shown in Table 4, and P-R curves are exhibited in
257  Fig. 8. Table 4 displays that CBAM_YOLOv4-tiny achieves the highest AP values in all
258  categories except for the upper limb. For the upper limb, YOLOv4-tiny combined with ECA
259  reaches the highest AP value.

260  **Discussion**

261  To perform the bathing tasks, we need to recognize the area to be bathed in the bathing scenario
262  and send the recognition information to the bathing robot arm for bathing behavior planning, as
263  shown in Fig. 9. By combining the skin detection results of 2D images with the depth
264  information obtained from the depth camera, we can model the localization of targets in 3D
265  space. To facilitate the robot to implement distinct bathing patterns for areas of the body, we
266  need to identify the skin located at diverse parts of the body. Therefore, we build small data sets
267  in the bathing scenarios to be used as learning samples for object detection models. And the
268  manual annotation is performed with a labelImg tool to classify skin regions into six categories
269  according to different parts.
270        Among object detection algorithms, one-stage detection algorithms are faster than two-stage
271  and are suitable for application in our scenario where real-time performance is required. In our
272  previous work, we explore the effectiveness of object detection models for skin detection with
273  multiple classifications and find the best YOLOv4 model from five models. We lightweight
274  YOLOv4 and impose three kinds of attention mechanisms on the YOLOv4-tiny. We find that
275  both CBAM and ECA improve the detection effect. Yet, SE makes the detection effect worse
276  instead, which implies that we need carefully choose the attention mechanism during practice.
277  Compared with Salah's work, we input data sets including images with six types of labels for
278  network training instead of skin and non-skin patches. We do not only identify skin or non-skin,
279  but also we want to know to which part of the body the skin belongs. To the best of our

280 knowledge, this is the first time we have investigated skin detection that can identify different
281 body parts.
282     There is rare research on object detection-based skin detection combined with robotic arms
283 for bathing tasks. Our study lightweights the YOLOv4 model and explores which attention
284 mechanism works best by imposing attention mechanisms on the YOLOv4-tiny model. However,
285 the YOLOv4-tiny possesses a reduction in mAP compared with the YOLOv4, creating some
286 challenges for high detection accuracy (Zhao et al., 2022b). The relatively small number of
287 trunks in the data sets results in poor detection of trunks because of individual privacy issues.
288 The foot occupies a small area in the whole body range. Foot features tend to disappear with
289 repeated down-sampling operations, resulting in poor detection of the foot.

## Conclusions

291 When using robots for autonomous bathing tasks, the perception of skin in bathing scenarios
292 needs to be accomplished first. To facilitate the embedded deployment, we use YOLOv4-tiny, a
293 lightweight model of YOLOv4, for skin recognition research based on our previous work. Three
294 kinds of attention mechanisms are overlaid in the YOLOv4-tiny. Use the test sets to test the
295 performance of the four models. Compared to the original YOLOv4-tiny, the YOLOv4-tiny
296 combined with the CBAM or ECA attention modules gives a certain increase in mAP, while the
297 addition of SE produces some degree of decrease. It is feasible to use attention mechanisms for
298 performance improvement of YOLOv4-tiny, but not every attention mechanism is suitable. In
299 addition, the best YOLOv4-tiny based on CBAM with 57.2% mAP is insufficient in practice. In
300 future work, we improve the detection for trunk and foot by expanding the trunk and foot
301 samples in the self-built data sets, aiming to guarantee deployment performance while achieving
302 high detection accuracy. Then, using the model with good performance, we convert the model
303 trained by Pytorch into an open neural network exchange(ONNX) model for easy deployment.

## Disclosures

305 The authors declare that they have no conflicts of interest.

## Acknowledgments

309

## References

311 Bhatt S, Soni H, Pawar T, Kher H. 2022. Diagnosis of pulmonary nodules on CT images using
312 YOLOv4. *International Journal of Online and Biomedical Engineering* 18(5):131-146 doi:
313 10.3991/ijoe.v18i05.29529.
314 Bin Z, Sun CF, Fang SQ, Zhao YH, Su S. 2022. Workshop safety helmet wearing detection
315 model based on SCM-YOLO. *Sensors* 22(17):6702 doi: 10.3390/s22176702.
316 Bochkovskiy A, Wang CY, Liao HYM. 2020. YOLOv4: optimal speed and accuracy of object
317 detection. *arXiv preprint* arXiv: 2004.10934.

318  Furusho Y, Ikeda K. 2020. Theoretical analysis of skip connections and batch normalization
319  from generalization and optimization perspectives. *Apsipa Transactions on Signal and*
320  *Information Processing* 9:e9 doi: 10.1017/ATSIP.2020.7.
321  He KM, Zhang XY, Ren SQ, Sun J. 2015. Delving deep into rectifiers: surpassing human-level
322  performance on ImageNet classification. In: *2015 IEEE International Conference on Computer*
323  *Vision (ICCV)*. Santiago, CHILE, 1026-1034 doi: 10.1109/ICCV.2015.123.
324  Hu J, Shen L, Sun G. 2018. Squeeze-and-excitation networks. In: *2018 IEEE/CVF Conference*
325  *on Computer Vision and Pattern Recognition(CVPR)*. Salt Lake City, UT, 7132-7141 doi:
326  10.1109/CVPR.2018.00745.
327  Huang GQ, Wan ZN, Liu XG, Hui JP, Wang Z, Zhang ZY. 2019. Ship detection based on
328  squeeze excitation skip-connection path networks for optical remote sensing images.
329  *Neurocomputing* 332:215-223 doi: 10.1016/j.neucom.2018.12.050.
330  Javadi N, Soltanizadeh H. 2021. Automated detection, 3D position of facial skin lesions using
331  genetic algorithm and Kinect camera. *Computer Methods in Biomechanics and Biomedical*
332  *Engineering-Imaging and Visualization* 10(1):48-54 doi: 10.1080/21681163.2021.1972342.
333  Kim Y, Hwang I, Cho NI. 2017. Convolutional neural networks and training strategies for skin
334  detection. In: *24th IEEE International Conference on Image Processing (ICIP)*. Beijing,
335  PEOPLES R CHINA, 3919-3923.
336  Krizhevsky A, Sutskever I, Hinton GE. 2017. ImageNet classification with deep convolutional
337  neural networks. *Communications of the ACM* 60(6), 84-90 doi: 10.1145/3065386.
338  Li P, Yu HL, Li SJ, Xu P. 2022. Comparative study of human skin detection using object
339  detection based on transfer learning. *Applied Artificial Intelligence* 35(15):2370-2388 doi:
340  10.1080/08839514.2021.1997215.
341  Lin HY, Lin CJ, Jeng SY, Yu CY. 2021. Integrated image sensor and hyperparameter
342  optimization of convolutional neural network for facial skin detection. *Sensors and Materials*
343  33(8):2911-2924 doi: 10.18494/SAM.2021.3301.
344  Lin TY, Dollar P, Girshick R, He KM, Hariharan B, Belongie S. 2017. Feature pyramid
345  networks for object detection. In: *30th IEEE/CVF Conference on Computer Vision and Pattern*
346  *Recognition (CVPR)*. Honolulu, HI, 936-944 doi: 10.1109/CVPR.2017.106.
347  Liu Y, Soh LK, and Lorang E. 2021. Investigating coupling preprocessing with shallow and deep
348  convolutional neural networks in document image classification. *Journal of Electronic Imaging*
349  30(4):043024 doi: 10.1117/1.JEI.30.4.043024.
350  Luo Y, Guan YP. 2017. Adaptive skin detection using face location and facial structure
351  estimation. *IET Computer Vision* 11(7):550-559 doi: 10.1049/iet-cvi.2016.0295.
352  Misra D. 2019. Mish: a self regularized non-monotonic neural activation function. *arXiv*
353  *Preprint* arXiv: 1908.08681.
354  Naing KM, Boonsang S, Chuwongin S, Kittichai V, Tongloy T, Prommongkol S, Dekumyoy P,
355  Watthanakulpanich D. 2022. Automatic recognition of parasitic products in stool examination
356  using object detection approach. *Peerj Computer Science* 8:e1065 doi: 10.7717/peerj-cs.1065.

357    Niu ZY, Zhong GQ, Yu H. 2022. A review on the attention mechanism of deep learning.
358    *Neurocomputing* 452:48-62 doi: 10.1016/j.neucom.2021.03.091.
359    Pratondo A, Bramantoro A. 2022. Classification of zophobas morio and tenebrio molitor using
360    transfer learning. *Peer Computer Science* 8:e884 doi: 10.7717/peerj-cs.884.
361    Redmon J, Divvala S, Girshick R, Farhadi A. 2016. You only look once: unified, real-time object
362    detection. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
363    Seattle, WA, 779-788 doi: 10.1109/CVPR.2016.91.
364    Redmon J, Farhadi A. 2017. YOLO9000: better, faster, stronger. In: *30th IEEE/CVF Conference*
365    *on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI, 6517-6525 doi:
366    10.1109/CVPR.2017.690.
367    Redmon J, Farhadi A. 2018. YOLOv3: an incremental improvement. *arXiv preprint* arXiv:
368    1804.02767.
369    Rezatofighi H, Tsoi N, Gwak JY, Sadeghian A, Reid I, Savarese S. 2019. Generalized
370    intersection over union: a metric and a loss for bounding box regression. In: *2019 IEEE/CVF*
371    *Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA, 658-
372    666 doi: 10.1109/CVPR.2019.00075.
373    Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang ZH, Karpathy A, Khosla A,
374    Bernstein M, Berg AC, Li FF. 2015. ImageNet large scale visual recognition challenge.
375    *International Journal of Computer Vision* 115(3):211-252 doi: 10.1007/s11263-015-0816-y.
376    Salah KB, Othmani M, Kherallah M. 2022. A novel approach for human skin detection using
377    convolutional neural network. *Visual Computer* 38(5):1833-1843 doi: 10.1007/s00371-021-
378    02108-3.
379    Shifa A, Imtiaz MB, Asghar MN, Fleury M. 2020. Skin detection and lightweight encryption for
380    privacy protection in real-time surveillance applications. *Image and Vision Computing*
381    94:103859 doi: 10.1016/j.imavis.2019.103859.
382    Sun HM. 2010. Skin detection for single images using dynamic skin color modeling. *Pattern*
383    *Recognition* 43(4):1413-1420 doi: 10.1016/j.patcog.2009.09.022.
384    Tian SS, Chen ZX, Chen BL, Zou WB, Li X. 2021. Channel and spatial attention-based Siamese
385    network for visual object tracking. *Journal of Electronic Imaging* 30(3):033008 doi:
386    10.1117/1.JEI.30.3.033008.
387    Wang QL, Wu BG, Zhu PF, Li PH, Zuo WM, Hu QH. 2020. ECA-Net: efficient channel
388    attention for deep convolutional neural networks. In: *2020 IEEE/CVF Conference on Computer*
389    *Vision and Pattern Recognition (CVPR)*. Seattle, WA, USA, 11531-11539 doi:
390    10.1109/CVPR42600.2020.01155.
391    Woo SH, Park J, Lee JY, Kweon IS. 2018. CBAM: convolutional block attention module.
392    *Computer Vision - ECCV 2018* 11211:3-19 doi: 10.1007/978-3-030-01234-2_1.
393    Zhang CB, Jiang PT, Hou QB, Wei YC, Han Q, Li Z, Cheng MM. 2021. Delving deep into label
394    smoothing. *IEEE Transactions on Image Processing* 30:5984-5996 doi:
395    10.1109/TIP.2021.3089942.

396    Zhang K, Wang YD, Li WY, Li CL, Lei ZC. 2022. Real-time adaptive skin detection using skin
397    color model updating unit in videos. *Journal of Real-time Image Processing* 19(2):303-315 doi:
398    10.1007/s11554-021-01186-9.
399    Zhao SJ, Zheng JC, Sun SD, Zhang L. 2022a. An improved YOLO algorithm for fast and
400    accurate underwater object detection. *Symmetry-Basel* 14(8):1669 doi: 10.3390/sym14081669.
401    Zhao L, Zhang Q, Peng B, Yang L. 2022b. Real-time object detector for low-end devices.
402    *Journal of Electronic Imaging* 31(1):013016 doi: 10.1117/1.JEI.31.1.013016.
403    Zheng ZH, Wang P, Liu W, Li JZ, Ye RG, Ren DW. 2019. Distance-IoU loss: faster and better
404    learning for bounding box regression. *arXiv preprint* arXiv:1911.08287.
405    Zhou QK, Zhang W, Li RZ, Wang J, Zhen SH, Niu F. 2022. Improved YOLOv5-S object
406    detection method for optical remote sensing images based on contextual transformer. *Journal of
407    Electronic Imaging* 31(4):043049 doi: 10.1117/1.JEI.31.4.043049.
408

# Figure 1

Example images in the data sets

# Figure 2

The structure of YOLOv4-tiny

# Figure 3

The structure of Resblock_body

# Figure 4

The specific implementation of SE

Input

↓

AdaptiveAvgPool2d

↓

Linear_Relu

↓

Linear_Sigmoid

↓

Reweight

↓

Output
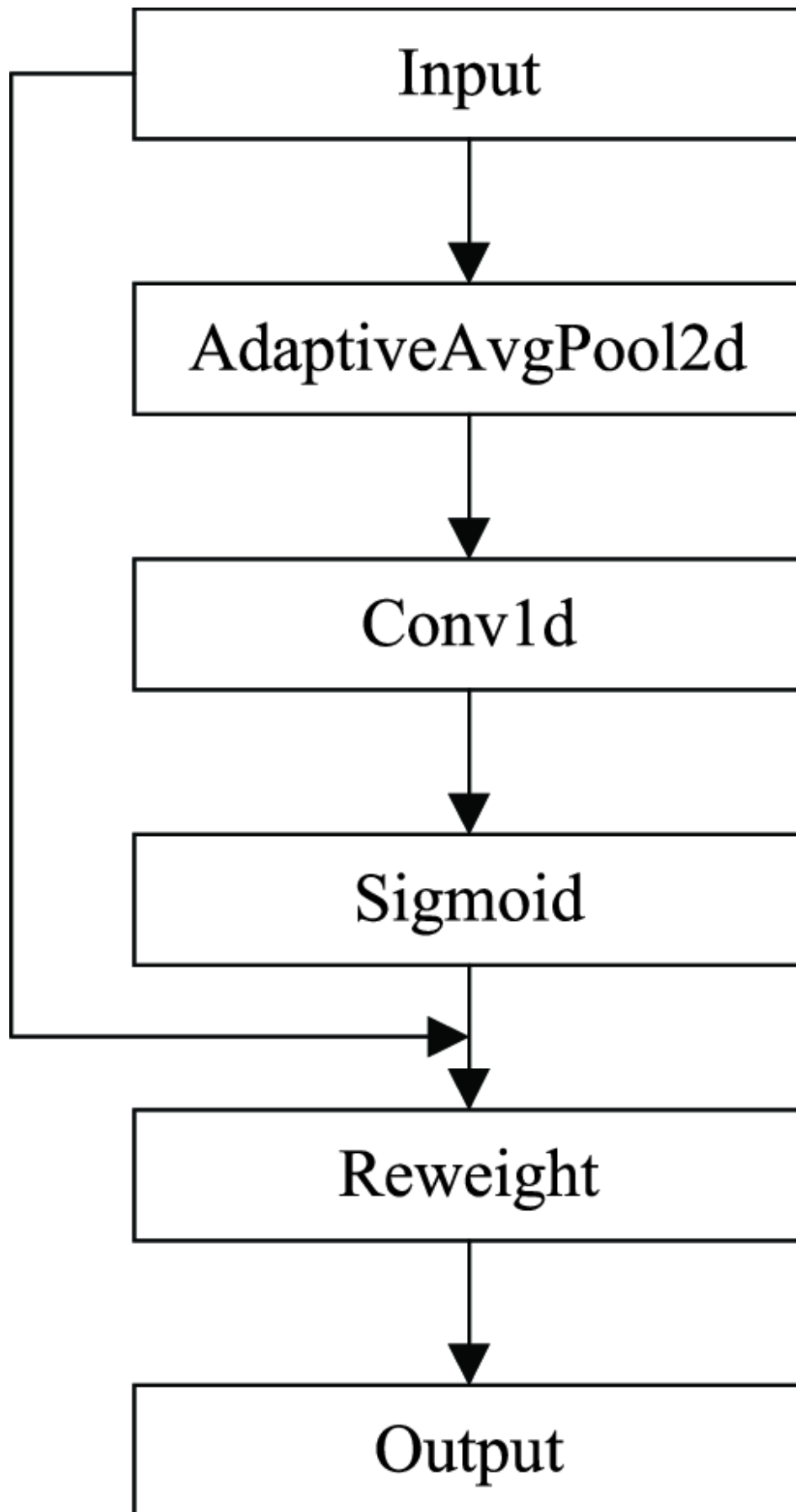
# Figure 5

The specific implementation of ECA
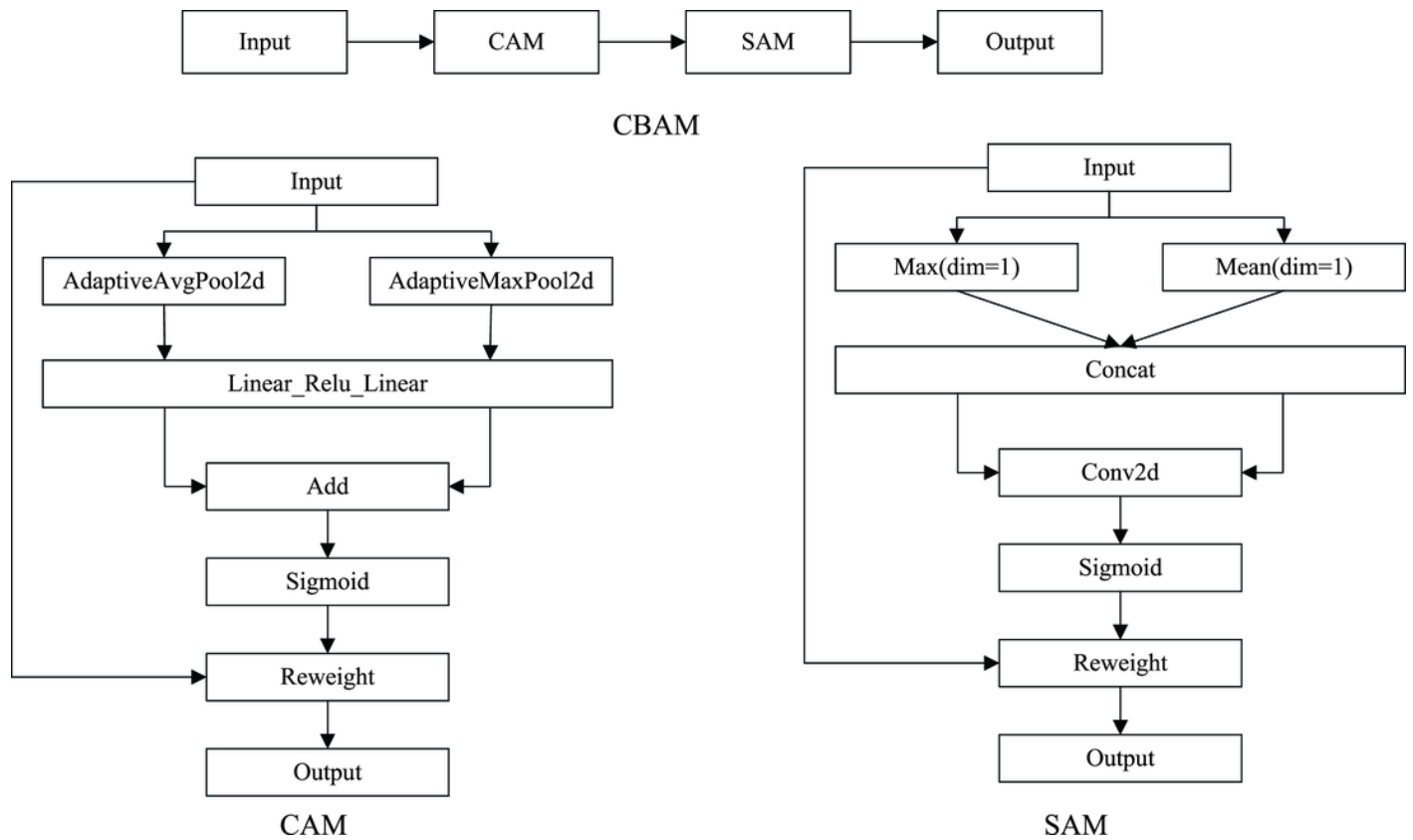
# Figure 6

The specific implementation of CBAM

# Figure 7

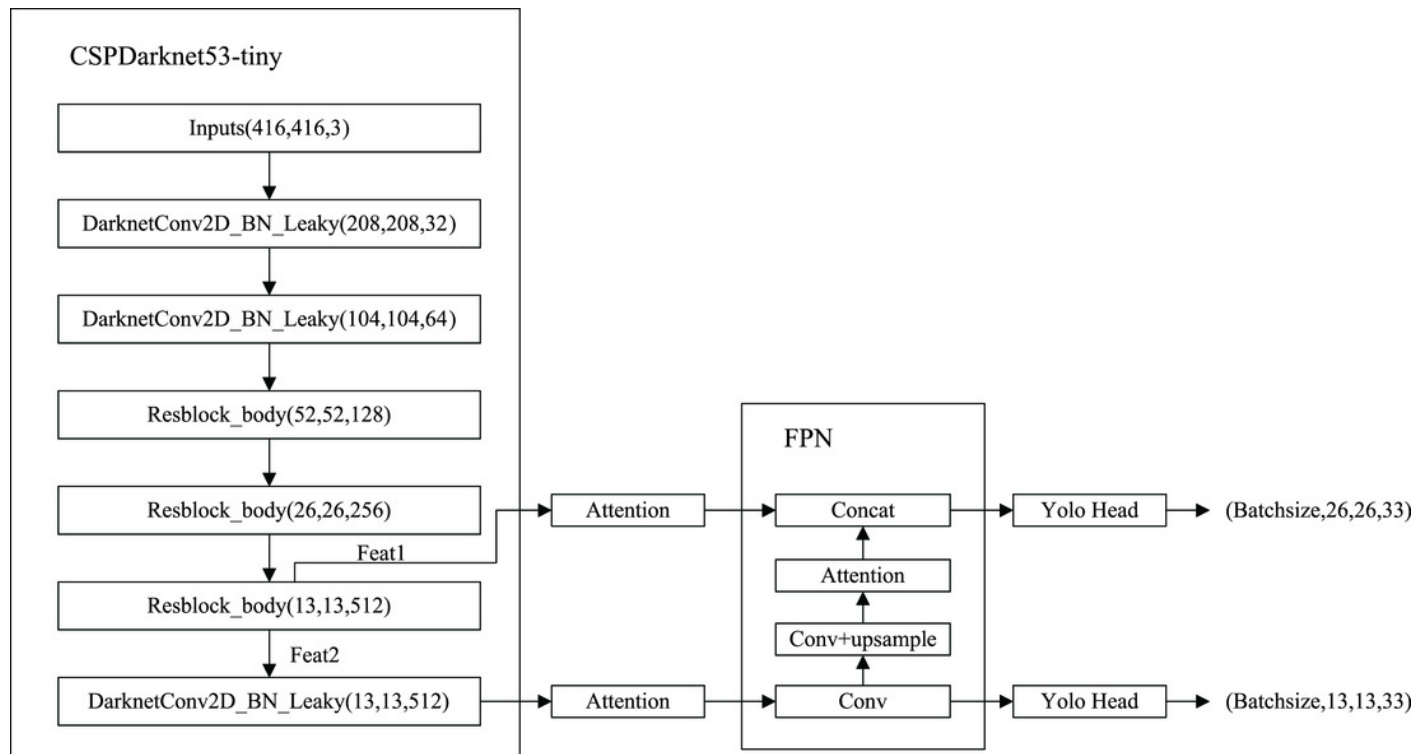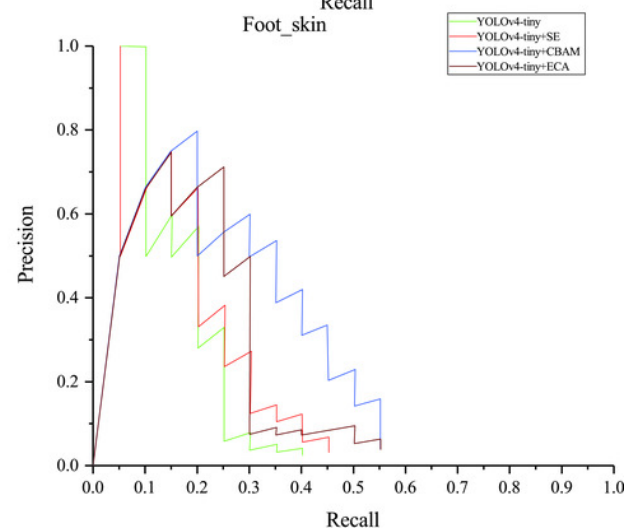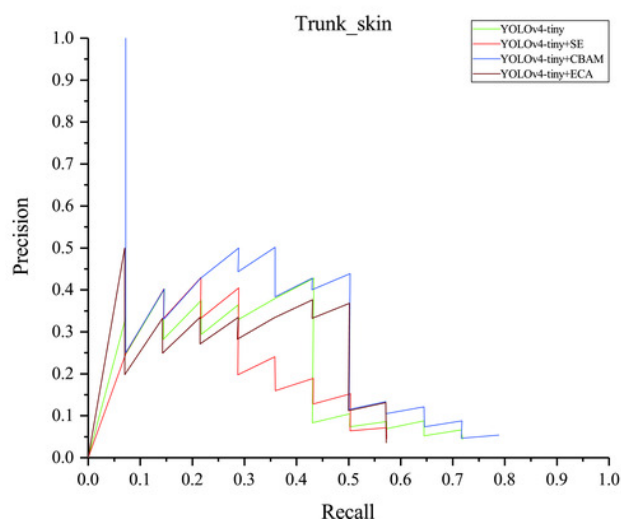Improved YOLOv4-tiny based on attention mechanisms
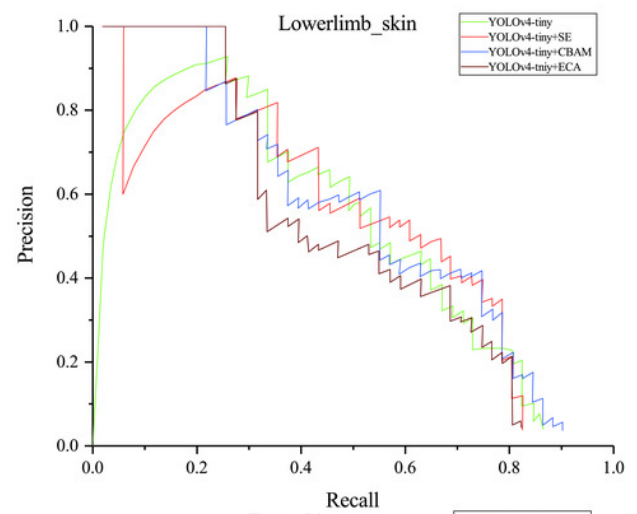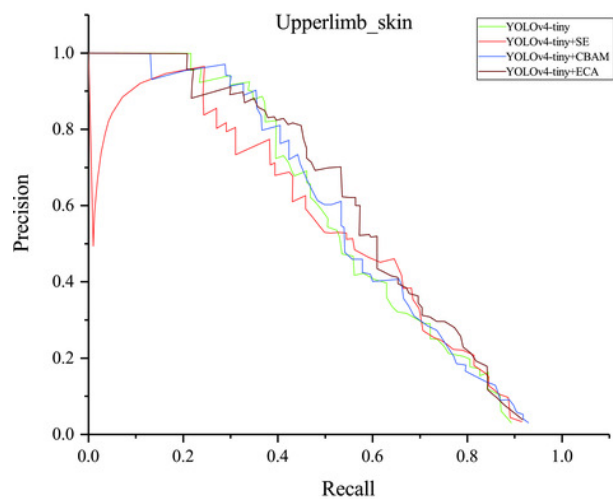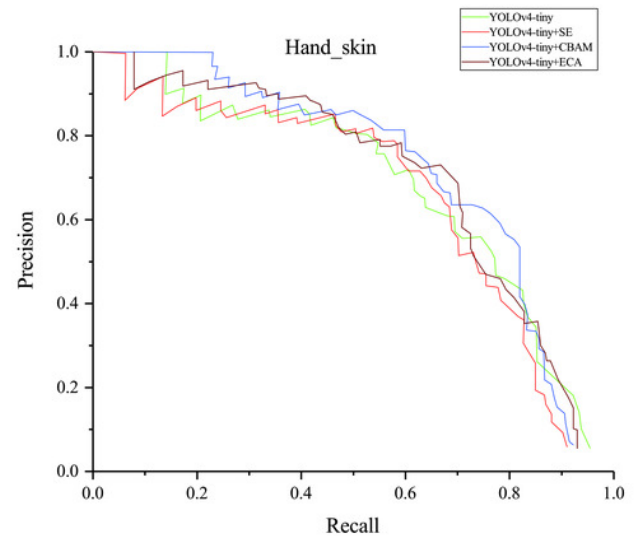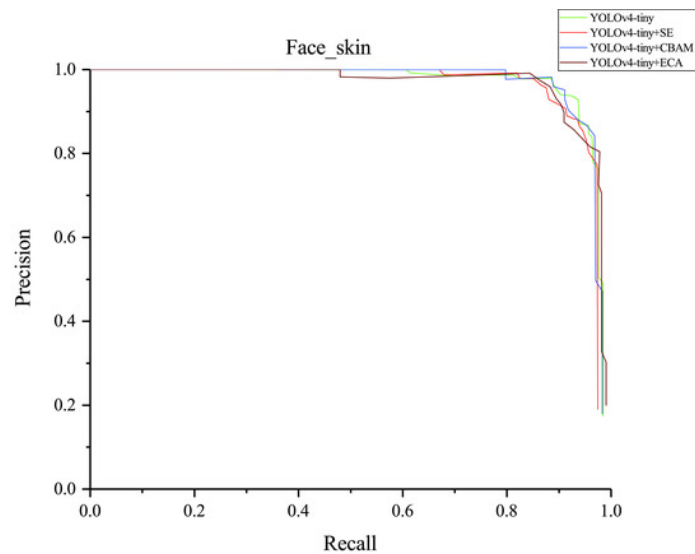
# Figure 8

P-R curves

# Figure 9

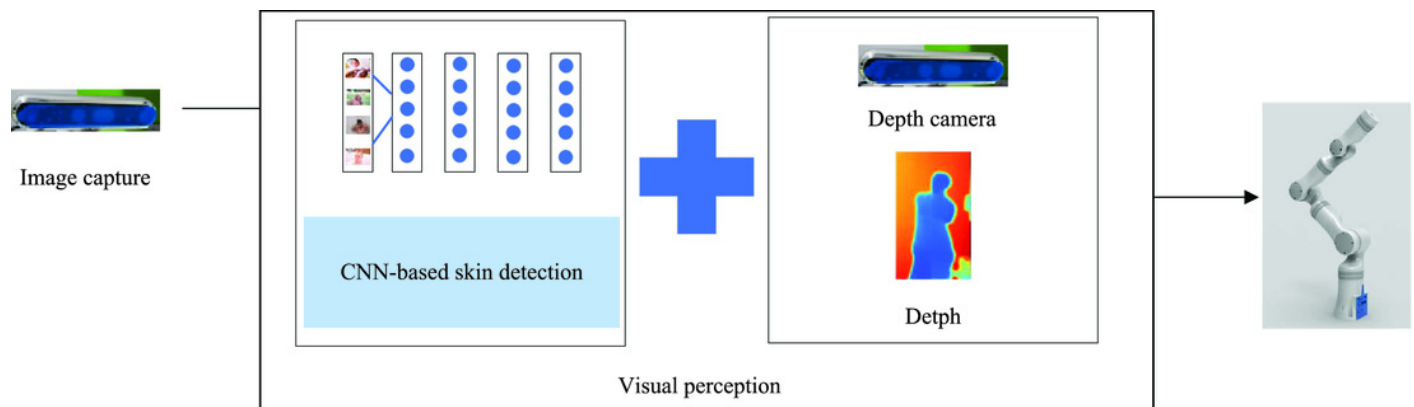The perception process in the bathing tasks: achieving the three-dimensional positioning of the target

**Table 1**(on next page)

Summary of IOU, GIOU, DIOU, CIOU

1
2
3

| | Features | Shortcomings |
|---|---|---|
| IOU | Representing the ratio of intersection and union of the GT box and the prediction box | When the prediction box and the GT box do not intersect, the loss function is not differentiable, leading losses cannot propagate |
| GIOU | scale invariant | Slow convergence speed and low positioning accuracy |
| DIOU | Overlapping area and center point distance are taken into account | Widely used in post-processing |
| CIOU | The consistency of aspect ratio is considered on the basis of DIOU | Widely used in post-processing |

4

**Table 2**(on next page)

Models information

1

2

| Model | Attention | mAP | Weight file(MB) |
|---|---|---|---|
| YOLOv4-tiny | - | 52.5% | 22.4 |
| SE_YOLOv4-tiny | SE | 51.6% | 22.6 |
| CBAM_YOLOv4-tiny | CBAM | **57.2%** | 22.8 |
| ECA_YOLOv4-tiny | ECA | 53.6% | 22.4 |

3

**Table 3**(on next page)

*A*, *B*, and *W* of all models

1

2

| Models | Attention | A | B | W |
|---|---|---|---|---|
| SE_YOLOv4-tiny | SE | 0.2 | -0.9 | -0.74 |
| CBAM_YOLOv4-tiny | CBAM | 0.4 | 4.7 | **3.15** |
| ECA_YOLOv4-tiny | ECA | 0 | 1.1 | 1.1 |

3

**Table 4**(on next page)

AP values for the six categories

1
2

|  | Face_skin | Hand_skin | Upperlimb_skin | Lowerlimb_skin | Trunk_skin | Foot_skin |
|---|---|---|---|---|---|---|
| - | 0.97 | 0.68 | 0.57 | 0.54 | 0.21 | 0.18 |
| SE | 0.96 | 0.67 | 0.55 | 0.55 | 0.17 | 0.21 |
| CBAM | **0.97** | **0.72** | 0.58 | **0.56** | **0.3** | **0.3** |
| ECA | 0.97 | 0.7 | **0.6** | 0.51 | 0.21 | 0.23 |

3