

Two-stage Deep Neural Network via Ensemble Learning for Melanoma Classification

Jiaqi Ding¹, Jie Song¹, Jiawei Li¹, Jijun Tang^{3,*} and Fei guo^{2,*}

¹*School of Computer Science and Technology, College of Intelligence and Computing, Tianjin University, Tianjin 300350, China*

²*School of Computer Science and Engineering, Central South University, Changsha 410083, China*

³*Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China*

Correspondence*:

Fei Guo, Jijun Tang

guofei@csu.edu.cn, jj.tang@siat.ac.cn

2 ABSTRACT

3 Melanoma is a skin disease with a high fatality rate. Early diagnosis of melanoma can effectively
4 increase the survival rate of patients. There are three types of dermoscopy images, malignant
5 melanoma, benign nevis and seborrheic keratosis, so using dermoscopy images to classify
6 melanoma is an indispensable task in diagnosis. However, early melanoma classification works
7 can only use the low-level information of images, so the melanoma cannot be classified efficiently;
8 and the recent deep learning methods mainly depend on a single network, although it can extract
9 high-level features, the poor scale and type of the features limited the results of the classification.
10 Therefore, we need an automatic classification method for melanoma, which can make full use
11 of the rich and deep feature information of images for classification. In this study, we propose
12 an ensemble method that can integrate different types of classification networks for melanoma
13 classification. Specifically, we first use U-net to segment the lesion area of images to generate a
14 lesion mask, thus resize images to focus on the lesion; then, we use five excellent classification
15 models to classify dermoscopy images, and adding squeeze-excitation block (SE block) to models
16 to emphasize the informative features; finally we use our proposed new ensemble network to
17 integrate five different classification results. The experimental results prove the validity of our
18 results. We test our method on the ISIC 2017 challenge dataset, and obtain excellent results
19 on multiple metrics, especially, we get 0.909 on ACC. Our classification framework can provide
20 an efficient and accurate way for melanoma classification using dermoscopy images, laying the
21 foundation for early diagnosis and later treatment of melanoma.

22 **Keywords:** Melanoma Classification, Dermoscopy Images, Ensemble Learning, Deep Convolutional Neural Network, Image
23 Segmentation

1 INTRODUCTION

24 Skin cancer is a major public health problem, with more than five million new cases diagnosed annually in
25 the United States(Siegel et al. (2016)Codella et al. (2018)). Melanoma is the fastest-growing and deadliest
26 form of skin cancer in the world, it causes many deaths each year. But it is noticed that melanoma multiplies

27 more slowly in the early stages, so if it is diagnosed early and treated promptly, the survival rates of patients
28 can be greatly improved.

29 Pigmentation lesions occur on the skin surface, and dermoscopic technology was introduced to
30 improve the diagnosis of skin melanoma. Dermoscope is a non-invasive skin imaging technique that
31 can magnify and illuminate skin areas, and then enhance visualization of deep skin by eliminating surface
32 reflections. Compared with standard photography, dermoscopy images can greatly improve the accuracy of
33 diagnosis(Codella et al. (2018)Kittler et al. (2002)). Dermatologists usually use "ABCD" rule to evaluate
34 skin lesions(Stolz (1994)Moura et al. (2019)). This rule analyzes asymmetry, boundary irregularities,
35 color variations and structures of lesions(Xie et al. (2016)). However, the differentiation of skin lesions by
36 dermatologists from dermoscopy images is often time-consuming, subjective and the diagnostic accuracy
37 depends largely on the professional level, so inexperienced dermatologists may not be able to make accurate
38 judgments. Therefore we urgently need an automatic recognition method that is non-subjective and can
39 assist dermatologists to make more accurate diagnosis.

40 However, there are still many challenges in automated recognition of melanoma, we show them in
41 Figure.1. The first column of Figure.1 shows malignant melanoma, the second column shows benign
42 nevis and the third column shows seborrheic keratosis. Firstly, skin lesions have great inter-class similarity
43 and intra-class variation in color, shape and texture, the different classes of skin lesion have high visual
44 similarity. Secondly, the area of skin lesions in dermoscopy images varies greatly, and the boundaries
45 between skin lesions and normal skin are blurred in some images. Thirdly, artifacts such as hair, rulers and
46 texture in dermoscopy images may make it hard to identify melanoma changes. All these factors make
47 automatic recognition more difficult.

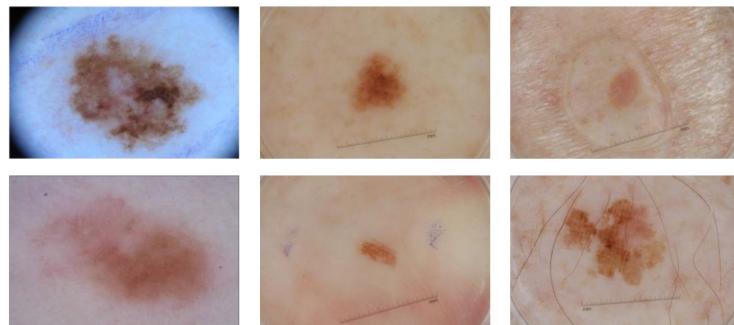


Figure 1. Some samples of dermoscopy images. From left to right: malignant melanoma, benign nevi and seborrheic keratosis.

48 In order to solve these problems, many researches have made attempts. Generally, automatic analysis
49 models includes four steps: image preprocessing, border detection or segmentation, feature extraction
50 and classification. In early works, a large number of studies used shallow models to classify dermoscopy
51 images, mainly using low-level features such as shape, color, texture or their combination(Ganster et al.
52 (2001)Mishra and Celebi (2016)), however, these shallow models for extracting low-level features lack high
53 level representation and powerful generalization capabilities. In recent years, convolutional neural network
54 has made great breakthroughs in image analysis tasks(Long et al. (2015); Krizhevsky et al. (2012); He et al.
55 (2015); Shin et al. (2016); Chen et al. (2017)), especially the deep convolutional neural networks (DCNN),
56 which can extract deep features and have better discrimination ability, have achieved improved performance.
57 So researchers started to apply DCNN to analyze medical images(Myronenko (2018)Roychowdhury et al.

58 (2015)), including image-based melanoma classification. However, deep neural networks still face great
59 challenges in the field of medical image analysis. DCNN requires large datasets to obtain more effective
60 features, while medical image data is often difficult to obtain and the datasets are relatively small. If a
61 small dataset is used directly for deep network training, it will lead to over-fitting of the model. Moreover,
62 a single network may not be able to extract all the informative features, and it is actually difficult to train
63 a model that performs well in all aspects. Therefore, we propose an integrated model based on transfer
64 learning to combine the results of multiple models to get better performance.

65 In this paper, we propose a novel two-stage ensemble method based on deep convolutional neural
66 networks. In the first stage, we perform the image segmentation, we use a segmentation network to generate
67 lesion segmentation masks, then we use these masks to resize the original images so that they are the same
68 size. In the second stage, we implement image classification, we utilize five state-of-the-art networks to
69 extract features, and add Squeeze-and-Excitation Blocks(Hu et al. (2018)) to the network to help emphasize
70 more informative features. Then we construct a new neural network using local connection to integrate the
71 classification results of these models, so that we can obtain the final classification result. We evaluate our
72 method on ISIC 2017 challenge dataset and obtain the best results on some metrics.

2 RELATED WORKS

73 2.1 Traditional Methods

74 Traditional methods are usually based on manually extracted features to classify dermoscopy images,
75 including features of color and texture. The "ABCD" rule is the standard used by dermatologists, and
76 there are many automatic classification methods are based on this rule. Barata et al. (2013) introduced
77 two different dermoscopy image detection systems, one used a global approach to classify skin lesions
78 and the other used local features and a bag-of-features(BoF) classifier. Ganster et al. (2001) used manual
79 features containing shape, boundary and radiometric features to describe lesions, and then used KNN(K-
80 Nearest Neighbor) to classify melanoma. Celebi et al. (2007) extracted descriptors related to shape, color
81 and texture from dermoscopy images and used nonlinear support vector machines to classify melanoma
82 lesions. Capdehourat et al. (2011) firstly preprocessed the image with hair removal, then used segmentation
83 algorithm to segment each image, finally trained the AdaBoost classifier with descriptors containing shape
84 and color information.

85 2.2 Deep CNN Models

86 In recent years, convolutional neural network(CNN) has been widely used in image
87 segmentation(Myronenko (2018); Roychowdhury et al. (2015); Dai et al. (2016)) and classification(Krizhevsky
88 et al. (2012)Simonyan and Zisserman (2014); Szegedy et al. (2015, 2016); He et al. (2016); Szegedy et al.
89 (2017); Huang et al. (2017); Chollet (2017)), object detection(He et al. (2015)Redmon et al. (2016); Liu et al.
90 (2016)) and other scopes of computer vision(Xiao et al. (2021); Chen et al. (2021b)). CNN models have
91 multiple layers to extract features. The network extractor mainly has two parts: convolutional layers and
92 pooling layers, and the network classifier is fully connected layer. Convolutional layers use convolutional
93 kernels to carry out convolution operation with input images to extract features. Kernels obtains features
94 of the whole image by sliding on it as a window. And the convolution operation of each kernel is only
95 connected to local area called receptive field of the input. Receptive field and weight sharing are important
96 parts of convolution neural network, they can effectively change the amount of training parameters. Pooling
97 operation is a kind of down sampling, its purpose is to reduce the training time, increase the receptive field
98 and prevent over-fitting, including widely used max pooling and average pooling. In addition, the fully

99 connected layer maps the learned feature representation to the label space for classification. If you need to
100 classify the samples into n classes, there are n neurons in the last fully connected layer.

101 Many CNN models have great performance on computer vision tasks(Cao et al. (2021); Chen et al.
102 (2021a); Feng et al. (2021)). Studies have shown that increase the number of layers in network can
103 significantly improve the performance(Simonyan and Zisserman (2014); Szegedy et al. (2015)). In recent
104 years, deep CNN has been proposed and performed well in the field of dermoscopy recognition. Codella
105 et al. (2015) used integrate CNN, sparse coding and SVM for melanoma classification. Yu et al. (2016)
106 proposed an automatic recognition method based on DCNN and residual learning, which firstly segmented
107 skin lesions and identified melanoma with two classifiers. Yu et al. (2018) proposed a network based on
108 DCNN and used feature coding strategy to generate representative features. Xie et al. (2016) processed
109 the incomplete inclusion of lesions in dermoscopy images and proposed a new boundary feature that can
110 describe boundary characteristics of complete and incomplete lesions. Lai and Deng (2018) combined
111 the extracted low-level features (color, texture) with the extracted high-level features of the convolutional
112 neural network for classification. González-Díaz (2019) proposed a CAD system called DermaKNet to help
113 dermatologists in their diagnosis.DermaKNet was divided into four parts, first segmenting the lesions in the
114 dermoscopic images using the Lesion Segmentation Network (LSN), then using the segmented masks to
115 perform data augmentation on the original data, and next the Dermoscopic Structure Segmentation Network
116 (DSSN) was used to segment the global and local features of the image, finally the image classification
117 is performed using the ResNet50-based network. Xie et al. (2020) proposed MB-DCNN to perform
118 segmentation and classification of dermoscopic images. They first used a coarse segmentation network
119 (coarse-SN) to generate a coarse lesion mask, which was used to assist the mask-guided classification
120 network (mask-CN) to locate and classify lesions, and the localized lesion regions were fed into the
121 enhanced segmentation network (enhanced-SN) to obtain a fine-grained lesion segmentation map. They
122 also proposed a new rank loss to alleviate the sample class imbalance problem. Gessert et al. (2020)
123 proposed a patch-based attention architecture in order to classify high-resolution dermoscopic images,
124 which was able to provide global contextual information to improve the accuracy of classification. In
125 addition, they proposed a new weighting loss to address the class imbalance in the data. Zunair and Hamza
126 (2020) first performed conditional image synthesis by learning inter-class mapping and synthesizing samples
127 of under-represented classes from over-represented classes using unpaired image-to-image translations,
128 thereby exploiting inter-class variation in the data distribution. Then the set of these synthetic and original
129 data was used to train a deep convolutional neural network for skin lesion classification. Bdair et al. (2021)
130 proposed FedPerl, a semi-supervised federated learning approach, which used peer learning and ensemble
131 averaging to build communities and encourage their members to learn from each other so that they can
132 generate more accurate pseudo labels. They also proposed the peer anonymization (PA) technique as a
133 core component of FedPerl. Datta et al. (2021) explored the goal of Soft-Attention to emphasize the value
134 of important features and to suppress features that cause noise. Then they compared the performance
135 of VGG, ResNet, Inception ResNet v2 and DenseNet architectures for classifying skin lesions with and
136 without the Soft-Attention mechanism. The results showed that the Soft-Attention mechanism improved
137 the performance of the baseline networks.

3 MATERIALS AND METHODS

138 In this section, we introduce our proposed two-stage ensemble network model. Firstly, in the first stage,
139 we train a segmentation network to segment skin lesions to get the lesion mask, and resize the mask area
140 to generate lesion image with the same size. Then, in the second stage, we use five networks with good

141 classification results on ImageNet to classify dermoscopy images, respectively. And we propose a new
 142 neural network to integrate the five results. The entire framework is shown in Figure.2.

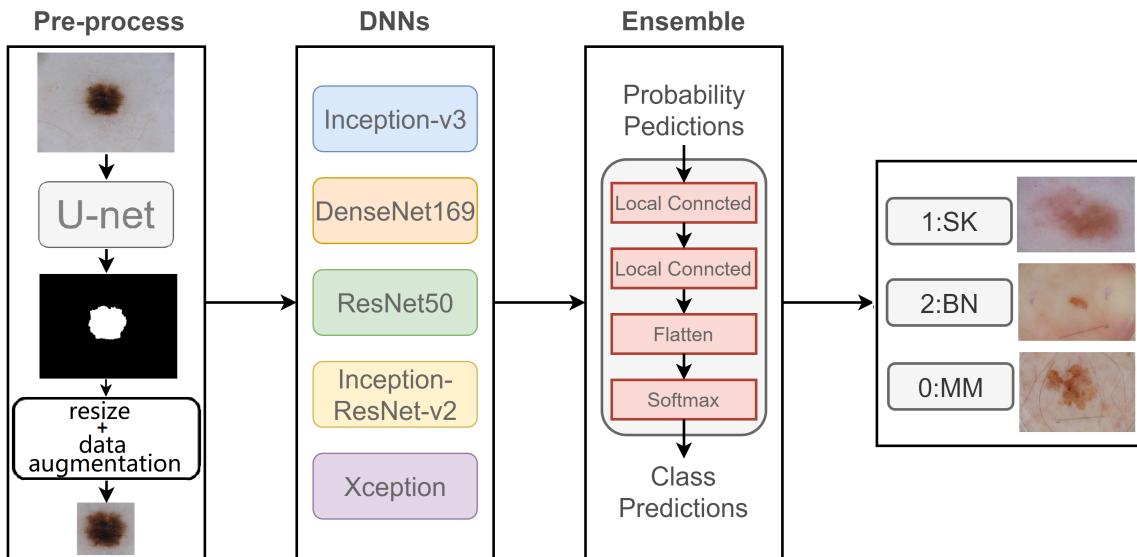


Figure 2. Flowchart of our proposed model.

143 3.1 Data Pre-processing

144 The deep network model needs a large amount of training data to better fit the real data distribution,
 145 and the lack of training data may lead to over-fitting and other problems, which will seriously affect the
 146 classification ability of the model. However, most medical image datasets do not have much data, which is
 147 one of the biggest challenges of medical image analysis. Data augmentation is one of common solutions to
 148 increase the amount of training data, and it can improve the model generalization ability. Therefore, we use
 149 different data augmentation methods on the original dataset, including rotation transform with 180 degrees,
 150 flip the images horizontally and vertically, and move the image height and width direction by 10%, so that
 151 each original image generates five new samples.

152 3.2 Skin Lesion Segmentation

153 Lesion segmentation plays an important role in the automatic analysis of skin lesion. It can separate the
 154 lesion from the normal skin, therefore the classifier can better identify the lesion features.

155 Unlike the classification network, which takes the images of fixed size as input and then outputs the class
 156 of each image, it gradually reduces the resolution of original images through convolution and max-pooling,
 157 and the feature maps it finally obtain are much smaller than the original image, then it classifies the
 158 feature maps through several fully connected layers. However, the output of segmentation network is
 159 the equal-sized prediction maps with input images. In the segmentation network, each pixel is a sample
 160 that needs to be classified into positive or negative. Therefore, the segmentation network needs decoder
 161 to compensate for the loss of feature resolution that caused by max-pooling. In our experiment, we use
 162 deconvolution operation in the decoder to obtain a prediction mask with the same size as the input image.

163 U-net(Ronneberger et al. (2015)) is an end-to-end deep convolutional neural network, which does not
 164 contain fully connected layer, but is composed of convolution layers and up-sampling layers. U-net has an
 165 encoder and a decoder. Encoder reduces the dimension of images and extracts feature, it is composed of
 166 four blocks, each of which consists two 3×3 convolution layers followed by a ReLU activation function,

167 and one max-pooling layer with stride of 2. Decoder also has four blocks, each contains a deconvolution
 168 layer, which double the size of feature maps, and two 3×3 convolution layer. So as for up-sampling
 169 operation in the decoder, U-net combines the output of up-sampling layer with feature map of symmetric
 170 encoder using skip-connection, so that the final output of network can consider both the shallow spatial
 171 information and deep semantic information. In this way, the outputs of the same size of the corresponding
 172 blocks in the encoder and decoder can be concatenated for segmentation and then the final prediction map
 173 is generated through a 1×1 convolution layer.

174 We train a U-net network to segment the original images and generate segmentation masks to show the
 175 lesion. These segmentation masks are used to crop the original images to help the classification network
 176 better focus on lesion features.

177 3.3 Skin Lesion Classification

178 The skin lesions have great inter-class similar visual effects, if we train our classification network use the
 179 original images, the results will be less effective. So we divide our classification model into three stages.
 180 First we segment skin lesions from original images using segmentation network and then resize them into
 181 a fixed size. Next we use five classification network with SE block to classify dermoscopy images. And
 182 finally we construct a convolution neural network to ensemble five results.

183 3.3.1 Resize

184 The size of lesions varies greatly, and in most dermoscopy images, the lesion area only occupies a small
 185 part of the image, and most parts are non-lesion areas that may affect classification. In this case, if the
 186 original images are directly classified, the size of skin lesion will seriously affect the performance of
 187 network. Therefore, we firstly segment skin lesions from the dermoscopy images, then adjust the segmented
 188 lesion to a fixed size. Compared with the network trained on original dermoscopy images, the network
 189 trained on segmented and resized images can better extract features and has better performance.

190 3.3.2 SE Block

191 The features extracted by a convolutional neural network can directly affect the results of subsequent
 192 tasks, either segmentation or classification. Therefore, improving the quality of the feature representation
 193 of the network is crucial to improve the final classification results. The role of the Squeeze-and-Excitation
 194 block(Hu et al. (2018)) is to further improve the classification accuracy by emphasizing the more important
 195 and informative features in the feature map. The SE block can be seen as a channel-wise attention
 196 mechanism, which emphasizes the importance of some features in the task by giving them greater weights.
 197 The specific strategy is shown below.

198 SE block is primarily concerned with the dependencies between feature channels. SE block do squeeze
 199 and excitation operation on feature maps $U(H \times W \times C)$. The squeeze operation include a global average
 200 pooling, it can map feature maps to feature vectors. For c -th feature map can be expressed as:

$$z_c = F_{sq}(u_c) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j) \quad (1)$$

201 where H and W represent the height and width of feature map separately. Then the excitation operation
 202 include two fully connected layer, a ReLU activation and sigmoid activation, so that it is able to fit complex
 203 correlations between channels by adding nonlinear processing through dimensional changes. The formula
 204 can be expressed as:

$$s = F_{ex}(z, W) = \sigma(W_2 \delta(W_1 z)) \quad (2)$$

205 where δ represents ReLU function and σ means Sigmoid, W_1 and W_2 are the weights of the first and
 206 second fully connected layer separately. In this way, the values in this feature vector are mapped to 0 – 1.
 207 Then the vector s can be multiplied as a channel descriptor with the original feature map to obtain the
 208 weighted feature map:

$$\tilde{x}_c = F_{scale}(u_c, s_c) = s_c u_c \quad (3)$$

209 Therefore, SE block is used to standardize feature maps according to their importance and highlight more
 210 informative feature maps, thus it can improve the network performance effectively. The schematic of
 211 adding SE Block to the five networks is shown in Figure 3. We add the SE Block in the same position in
 212 each network, that is, after feature extraction(orange box in Figure 3) and before final classification of each
 213 network.

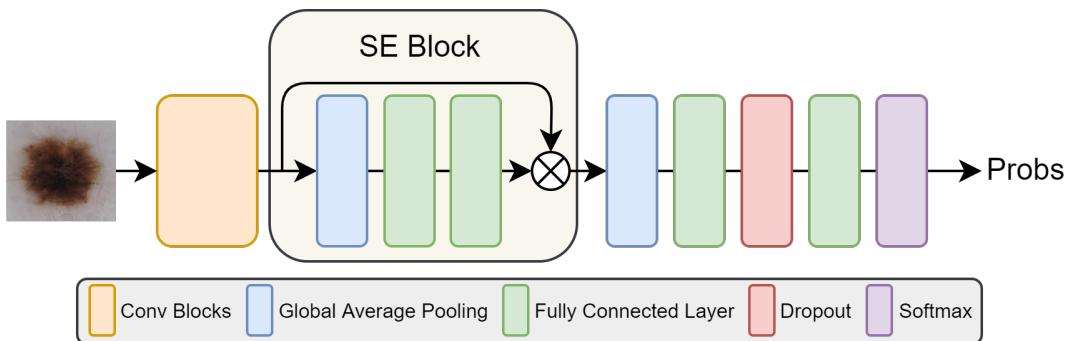


Figure 3. The illustration of five network structures after adding SE Blocks.

214 3.3.3 Network Model

215 For ensemble problems, in addition to the ensemble method, the basic model of integration is also
 216 important. We use five state-of-the-art networks as basic network for our integration, which are Inception-
 217 v3, Densenet169, ResNet50, Inception-ResNet-v2 and Xception. These networks all have good performance
 218 on image classification tasks.

219 3.3.3.1 Inception-v3

220 Inception module(Szegedy et al. (2015)) used 1×1 , 3×3 , 5×5 convolution layer at the same time,
 221 then concatenated three kinds of outputs and transmitted it to next module. In this way, it can consider
 222 information of different scales at the same time by increasing the width of the network. In addition,
 223 Inception module also can split channel-wise and spatial-wise correlation and small size of convolution
 224 kernel can greatly reduce the parameters. On the basis of Inception module, Inception-v3(Szegedy et al.
 225 (2016)) replaced the 5×5 convolution layer in the original Inception network with two 3×3 convolution
 226 layers in order to further reduce the amount of parameters while maintaining the receptive field and
 227 increasing the ability of representation. Furthermore, another innovation of Inception-v3 was to decompose
 228 a large $n \times n$ convolution kernel (for example, a 7×7 convolution kernel) into two one-dimensional
 229 convolution kernels with the size of $n \times 1$ and $1 \times n$ respectively. This can increase the model's nonlinear
 230 representation capability while reducing the risk of over-fitting.

231 3.3.3.2 ResNet-50

232 ResNet(He et al. (2016)) appeared to alleviate the problem of vanishing/exploding gradients. ResNet was
 233 composed of a set of residual blocks, each of which is composed of several layers, including convolutional
 234 layer, ReLU layer and batch normalization layer. And for each residual block, its input was directly added
 235 to its output via identity, a short connection that allowed us to perform residual learning, this is the key to

236 solve gradients problem when training deep networks. A residual block can be formulated as:

$$H_l = H_{l-1} + F(H_{l-1}) \quad (4)$$

237 where H_l and H_{l-1} are the output and input of the l -th residual block respectively. $F(x)$ represent the
238 residual mapping function of stacked layers. It is obvious that the dimensions of H_{l-1} and $F(H_{l-1})$
239 should be equal. But convolution operation usually change the dimensions, so a linear projection W_s is
240 used to match the dimensions. So the Eq.4 can be converted to:

$$H_l = W_s H_{l-1} + F(H_{l-1}) \quad (5)$$

241 Therefore, ResNet-50 was obtained by stacking the residual blocks to make the final network layer count
242 to 50.

243 3.3.3.3 Densenet169

244 Densenet(Huang et al. (2017)) was inspired by Resnet. It also used connections to alleviate the problem
245 of vanishing gradients, but it did not use residual blocks to achieve this goal. Densenet was composed of
246 dense blocks. In each dense block, as shown in Figure 4, the input of the n -th layer was the result of the
247 concatenating of all the previous $n-1$ layers. In this way, when performing related operations on the n -th
248 layer, the utilization of the features of all the previous layers can be maximized. This feature reuse method
249 can make the features work better while reducing the amount of parameters.

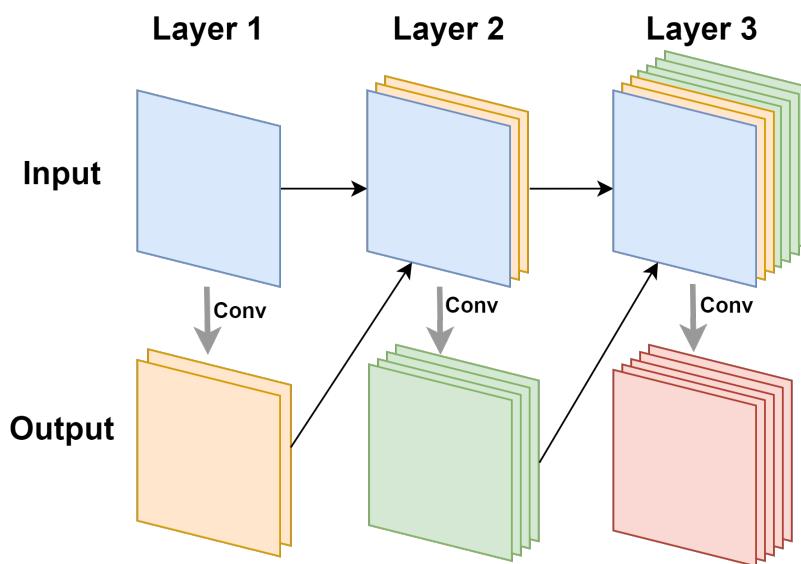


Figure 4. The illustration of feature reuse of dense block.

250 3.3.3.4 Inception-ResNet-v2

251 Inception-ResNet-v2(Szegedy et al. (2017)) combined Inception module with residual learning. It was
252 based on Inception-v4, which was deeper and better than Inception-v3, but had more parameters. Inception-
253 ResNet-v2 added residual identities to different types of Inception modules of Inception-v4, so that the
254 network converged faster, and the training time of the network was shortened.

255 3.3.3.5 Xception

256 Xception(Chollet (2017)) was an improvement to Inception-v3. It mainly replaced ordinary convolution
257 in Inception-v3 with depthwise separable convolution. The multiple convolution kernels of depthwise

258 separable convolution only processed part of feature maps produced by the previous layer. For example, for
 259 the result of 1×1 convolution output from the Inception module, depthwise separable convolution referred
 260 to using three 3×3 convolution kernels to operate on one-third of the channel of this result, and finally
 261 three results from three 3×3 convolution kernels were concatenated together. In this way, the amount
 262 of parameters can be greatly reduced. And the author believed that Xception can decouple the channel
 263 correlation and spatial correlation of the features, thereby producing better computational results.

264 We use these five pre-trained networks on ImageNet as feature extractors, then add SE blocks after every
 265 extractors to emphasize more informative features. Then, a full connected layer of 128-dimension is used
 266 to generate the final feature vector, and finally we use softmax classifier to obtain class predictions.

267 3.3.4 Ensemble Learning

268 There are usually two ways to ensemble multiple networks: averaging and voting. Averaging refers to the
 269 average results of multiple networks, with each network accounting for the same proportion, so that having
 270 the same influence on the final result. But for each class, some networks produce better results, and some
 271 have relative worse effect, taking the average directly would reduce the advantage of good networks.

272 For voting ensemble, we can implement it through neural networks. In detail, the neural network we build
 273 for ensemble learning is equivalent to a new classifier, whose input is the classification probabilities from
 274 five networks, and whose output is the final classification result. The reason we chose to build the classifier
 275 with local connected layer instead of full connected layer is that full connected layer will be connected
 276 to all the outputs of the previous layer, while local connected layer will only be connected to parts of the
 277 previous layer. In this case, the part of output of the ensemble network will only be determined by a specific
 278 input, and the prediction of one class will not be influenced by the other two classes because the local
 279 connection layer extracts features for each class separately, so the network will produce more accurate
 280 classification results. This new network is used to integrate the results of the five networks, consisting of
 281 two local connected layers and a softmax layer, as shown in Figure 2. The result has an improvement over
 282 the averaging ensemble method.

4 RESULTS

283 4.1 Dataset

Table 1. Details of ISIC 2017 challenge dataset.

Subsets	MM	SK	BN	Total
Training	374	254	1372	2000
Validation	30	42	78	150
Testing	117	90	393	600

284 The dataset we use to evaluate our method was provided by ISIC 2017 challenge organized by The
 285 International Society for Digital Imaging of the Skin(Codella et al. (2018)). It includes 2750 dermoscopy
 286 images and is divided into three subsets: 2000 for training, 150 for validation and 600 for testing. The images
 287 in the dataset are classified as three classes: benign nevi(BN), seborrheic keratosis(SK) or melanoma(MM).
 288 The details of ISIC 2017 challenge dataset is shown in TABLE 1, MM refers to melanoma, SK refers to
 289 seborrheic keratosis and BN refers to benign nevi. Also, we can see from the Figure. 5 that the distribution
 290 of training, validation and test sets is very uneven, the images of BN are far more than the images of the
 291 other two classes in three subsets. In addition, the ISIC 2017 dataset also provides dermoscopy images
 292 with their binary masks as their segmentation ground truth.

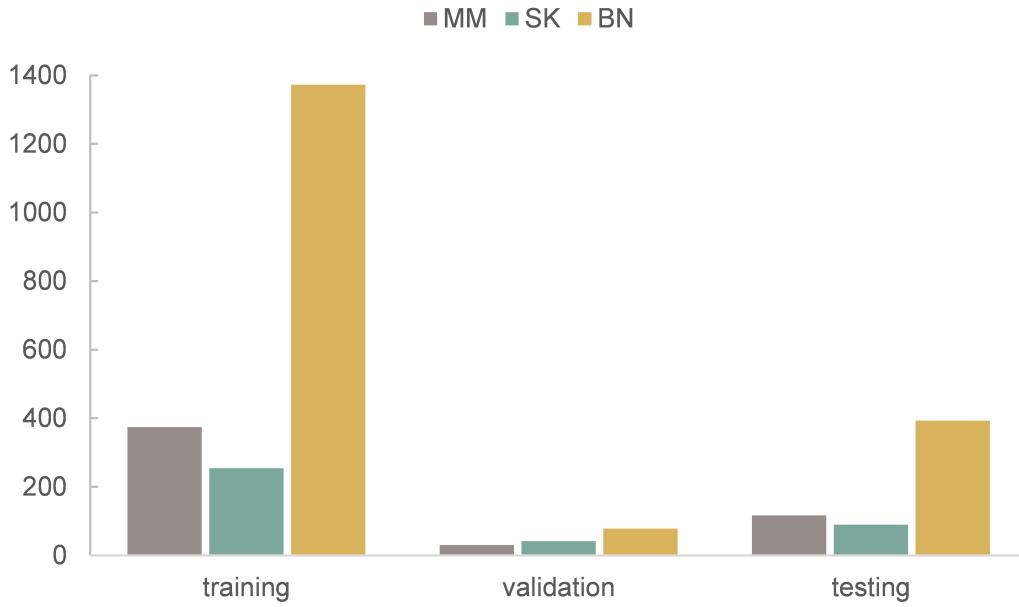


Figure 5. The distribution of training, validation and test sets of ISIC 2017 challenge dataset.

293 The ISIC 2017 challenge consists of two binary classification subtasks: melanoma or others and seborrheic
294 keratosis or others.

295 4.2 Implementation

296 Our method is implemented with Keras on a computer with GeForce RTX 2080Ti GPU. The images with
297 the size of 224×224 are taken as input of model, so all dermoscopy images are resized to 224×224 after
298 segmentation. We use Adam algorithm as optimizer, the learning rate is set as 0.0001 initially. Our epoch
299 number is set to 100 initially. To prevent over-fitting, we use early stopping method with patience of 10
300 epochs.

301 4.3 Metrics

302 We use accuracy(ACC), recall, precision, F1-score and AUC (Area Under ROC Curve) as classification
303 metrics. They are defined as:

$$304 ACC = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

$$305 recall = \frac{TP}{TP + FN} \quad (7)$$

$$306 precision = \frac{TP}{TP + FP} \quad (8)$$

$$f1score = \frac{2 \times precision \times recall}{precision + recall} \quad (9)$$

307 where TP, TN, FP and FN denote the number of true positive, true negative, false positive and false
308 negative. The number of three classes in our dataset are imbalanced, so in this case, ACC cannot well
309 reflect the performance of our classifier, therefore we use AUC, the same indicator as ISIC classification
310 challenge Codella et al. (2018), as the main metric.

311 4.4 Performance on Multi-class Classification

312 Our method is divided into three parts. After segmenting and cropping the original dermoscopy images,
 313 five pre-trained models are used to do classification, and then ensemble the results of these models to
 314 generate final result. To verify our method, in this section, we modify the dataset, and convert the two
 315 binary classification tasks into a multi-classification task. Then we compare the performance with and
 316 without segmentation and resize, and the performance before and after ensemble. TABLE 2 shows the
 317 experimental results with and without segmentation under one pre-trained network called Inception-v3. It
 318 can be seen that the network has better performance running on the segmented images than on the original
 319 images. As shown in Figure 6, especially on ACC and AUC, the results of network with segmentation gets
 320 0.791 and 0.883 respectively, which are much higher than that of network without segmentation. This is
 321 because the size of skin lesions varies greatly, and there are some interference factors such as artificial
 322 rulers in the original dermoscopy images. Segmentation can remove these interference factors to some
 323 extent, so that the network can better identify features.

Table 2. Classification results with or without segmentation.

Methods	ACC	precision	recall	f1 score	AUC
Without segmentation	0.698	0.598	0.622	0.592	0.781
With segmentation	0.791	0.634	0.688	0.659	0.883

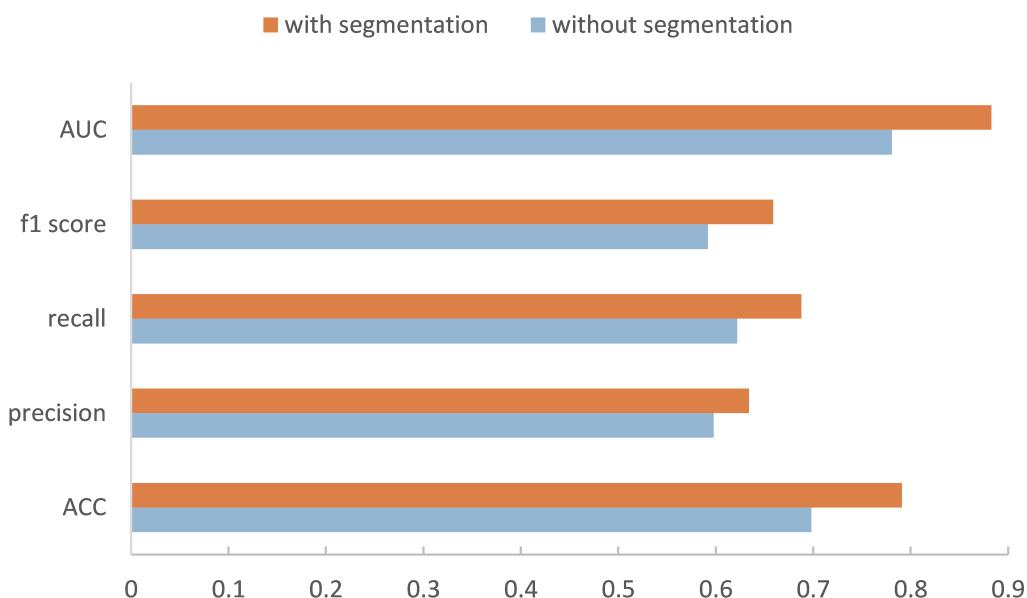


Figure 6. Performance of our method with or without segmentation.

324 In the ensemble stage, we construct a neural network model with two local connected layers with softmax
 325 classifier to fuse the results of five basic networks. Our new ensemble method can further improve the
 326 performance, and is better than the commonly used ensemble method. TABLE 3 lists the results of the
 327 five pre-trained models we use and the results of averaging ensemble and our ensemble method. It can be
 328 seen that the fusion model have better performance than any single network and average method on most
 329 metrics. For the recall and f1 scores, our ensemble method is 0.033 and 0.007 lower than Xception, but it is

330 higher than other methods in other metrics. Especially, it has a 2% improvement on AUC over the result of
 331 best network, i.e.Xception. Also our ensemble method is better than traditional average ensemble method
 332 on all metrics except for recall.

333 We also compare the amount of parameters and training time of different networks (including our
 334 ensemble network). From the TABLE 4, we can see that the classification networks have more parameters,
 335 especially Inception-Resnet-v2, which has up to 54.87M. But compared to these classification networks, our
 336 ensemble network has very few parameters, only 423. For training time, since the classification networks
 337 have been pre-trained on ImageNet, we just need to fine-tune the networks during training, and our training
 338 set is small, so we can see that the training time of each network is relatively short (when training 100
 339 epochs). At the same time, we can also notice that the training time of the network is not entirely determined
 340 by their parameters, but is also related to the parallelism of the model and the memory access cost. In
 341 addition, these five classification networks are independent of each other, so they can be trained at the
 342 same time, which can also greatly reduce training time. Finally, our ensemble network requires very little
 343 training time, only 20s.

Table 3. Results of different networks and two ensemble methods on multi-classification task.

Methods	ACC	precision	recall	f1 score	AUC
Inception-v3	0.792	0.634	0.688	0.659	0.883
Densenet169	0.800	0.739	0.727	0.722	0.881
Resnet50	0.762	0.676	0.678	0.672	0.864
Inception-Resnet-v2	0.800	0.736	0.726	0.725	0.873
Xception	0.810	0.75	0.748	0.748	0.896
Average	0.793	0.724	0.724	0.719	0.880
Ensemble	0.851	0.769	0.715	0.741	0.913

Table 4. The amount of parameters and the training time of each network.

Networks	Inception-v3	Densenet169	Resnet50	Inception-Resnet-v2	Xception	Ensemble
Params	22.56M	13.22M	24.32M	54.87M	21.59M	423
Time(s)	1900	3200	1900	3000	2700	20

344 4.5 Performance on Binary Classification

345 ISIC 2017 challenge has two binary classification tasks: melanoma or others and seborrheic keratosis or
 346 others, so we also carry out the experiment regarding to challenge tasks. We show the results of melanoma
 347 classification and seborrheic keratosis classification in the form of radar diagrams, as shown in Figure. 7.
 348 Polar coordinates represent different metrics and each line represents a network. It can be seen that
 349 our method performs pretty well on both tasks. For the classification of melanoma, it is clear that our
 350 performance is the highest in all metrics, especially in precision, where we outperform the second highest,
 351 DenseNet, by more than 10%; secondly, for the f1 score, which can take into account both positive and
 352 negative samples, our method also outperforms the rest of the networks by about 5%; finally, for our main
 353 metric, AUC, we also surpass the other networks by a large margin. As for the classification of seborrheic
 354 keratosis, although the advantage of our method is not as obvious as when classifying melanoma, it still
 355 performs well. Firstly, our method still outperforms the other networks in terms of AUC, which is our
 356 main metric; secondly, for precision and ACC, our method leads by a small margin; and for recall and f1,

357 we are slightly below the performance of Inception-Resnet-v2 and Xception. In general, our method is
 358 very efficient for classifying melanoma, although it is not significantly superior for classifying seborrheic
 359 keratosis, so it can improve the accuracy of classification in this task in general.

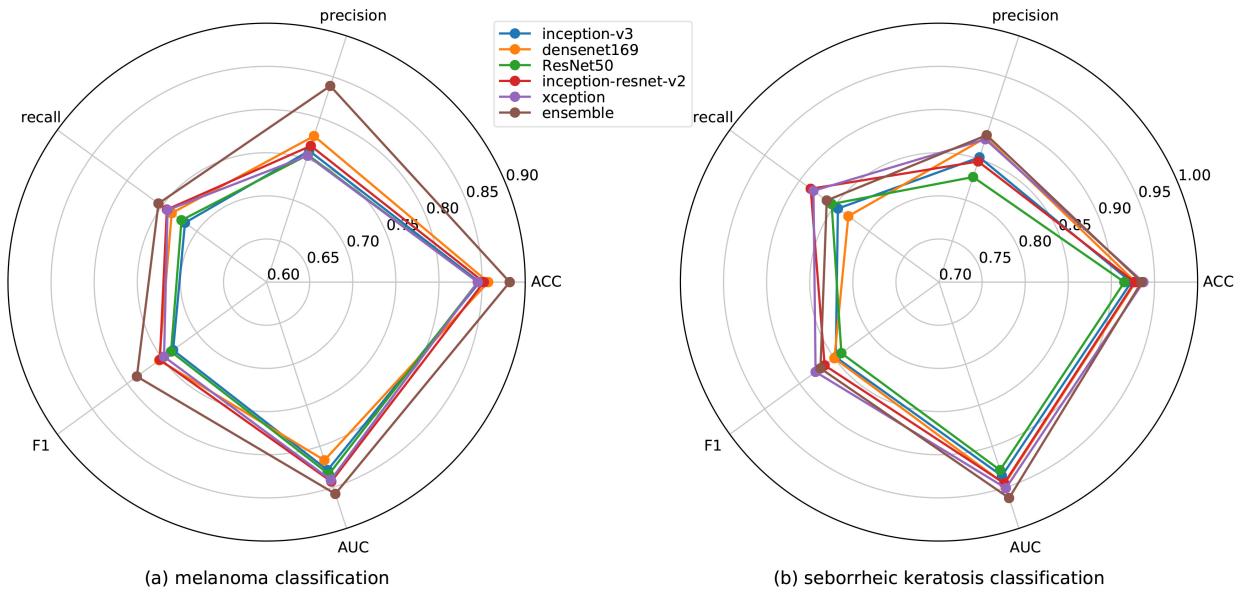


Figure 7. Results of melanoma and seborrheic keratosis classification for different networks.

Table 5. Average results of two skin lesion classifications of different networks.

Methods	ACC	precision	recall	f1 score	AUC
Inception-v3	0.885	0.806	0.781	0.791	0.883
Densenet169	0.893	0.827	0.783	0.802	0.882
Resnet50	0.88	0.792	0.788	0.789	0.882
Inception-Resnet-v2	0.89	0.807	0.814	0.809	0.894
Xception	0.891	0.814	0.811	0.812	0.896
SVC ¹	0.911	0.798	0.66	0.719	0.813
Random forest	0.912	0.802	0.664	0.721	0.816
Extra-Trees	0.911	0.805	0.65	0.716	0.809
KNN	0.908	0.782	0.657	0.709	0.81
GBDT ²	0.91	0.808	0.644	0.71	0.807
Ensemble	0.909	0.859	0.808	0.828	0.911

¹ Support Vector Classification

² Gradient Boost Decision Tree

360 We average the performance of all networks and ensemble methods on two binary tasks and show them
 361 in TABLE 5. When compared to a single network, it can be seen that our ensemble method can effectively
 362 improve the performance, especially the AUC is 1% better than the best single network, i.e.Xception.
 363 At the same time, for precision and f1 score, our ensemble network is also the highest one. In addition,
 364 when compared to other ensemble methods, we use several machine learning classifier to do ensemble

365 as comparison. We can see that except that ACC is 0.003 lower than Random forest, we are significantly
 366 better than machine learning methods on other metrics. We also illustrate this comparison in Figure 8, so
 367 we can more intuitively see the advantages of our ensemble method in various metrics.

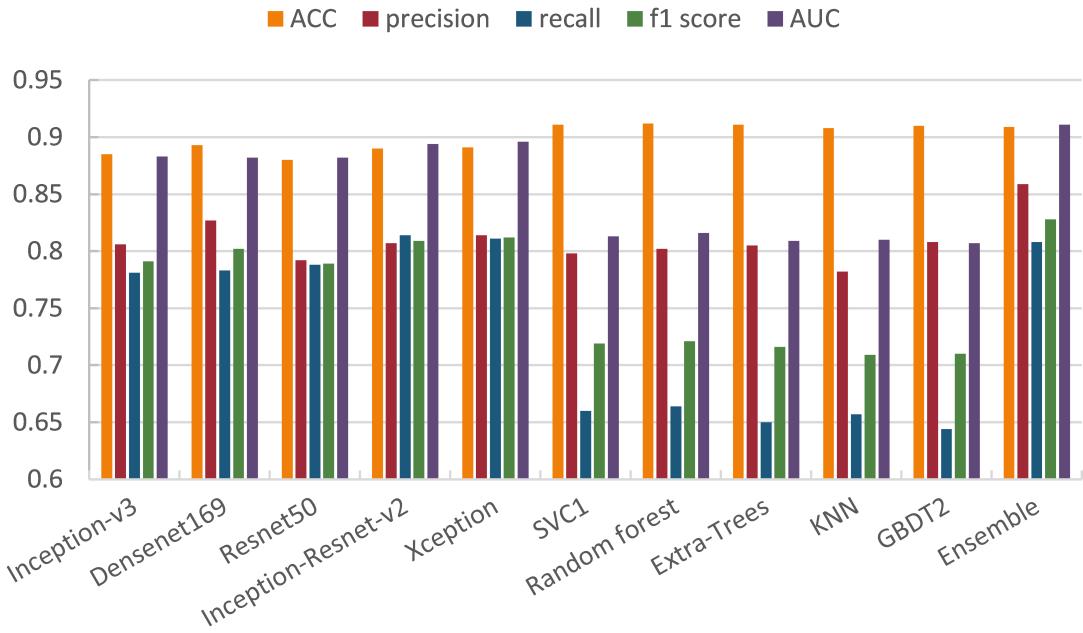


Figure 8. Comparison of different methods on skin lesion classification.

368 4.6 Comparison of various predictors

369 In TABLE 6, we compare our method to the top five performance in the ISIC 2017 challenge skin lesion
 370 classification task(Matsunaga et al. (2017); Díaz (2017); Menegola et al. (2017); Bi et al. (2017); Yang et al.
 371 (2017)) and some excellent methods in recent years. Most of the networks participating in the chanllenge
 372 used external images which we do not do this. From In TABLE 6, It can be seen that our method achieves
 373 0.909 and 0.859 on ACC and precision, which are highest on these metrics. Besides, we get 0.911 on AUC,
 374 which is 0.048 lower than the Datta et al. (2021). For f1 score, our method obtains 0.828, which is 0.023
 375 lower than the best score. But for recall, our model's performance is a bit unsatisfactory, which shows that
 376 our model still has some shortcomings in classifying positive samples.

5 CONCLUSION

377 In this paper, we have the following innovations. 1)we propose new two-stage ensemble method that
 378 integrates five excellent classification models to classify skin melanoma; 2)we also propose a new method
 379 of segmenting the lesion area of the dermoscopy image to generate a mask of the lesion area, so that the
 380 image can be resized to focus on the lesion; 3)we propose a new ensemble network that can use local
 381 connected layers to effectively integrate the classification results from the five classification networks. We
 382 test our method on the ISIC 2017 challenge dataset and get pretty good results. In future work, we will
 383 explore more effective classification methods based on the characteristics of dermoscopy images and the
 384 association of different classes of dermoscopy images, especially in process of pre-processing, because the
 385 experimental results show that our segmented images can largely improve the accuracy of classification.

Table 6. Comparison among our method, some existing methods sand the top five ISIC2017 classification challenge.

Method	ACC	precision	recall	f1 score	AUC
Top 1	0.816	0.748	0.856	0.851	0.911
Top 2	0.849	0.747	0.140	0.236	0.910
Top 3	0.883	0.752	0.451	0.564	0.908
Top 4	0.888	0.732	0.508	0.600	0.896
Top 5	0.873	0.665	0.568	0.613	0.886
Zhang et al. (2019)	0.868	-	0.878	-	0.958
González-Díaz (2019)	-	-	-	-	0.917
Xie et al. (2020)	0.904	-	0.786	-	0.938
Datta et al. (2021)	0.833	-	0.916	-	0.959
Ours	0.909	0.859	0.808	0.828	0.911

FUNDING

386 This work is supported by a grant from the National Natural Science Foundation of China (NSFC 62172296,
 387 61972280) and National Key R&D Program of China (2020YFA0908400). This work was also supported
 388 by the Shenzhen Science and Technology Program (No. KQTD20200820113106007).

CONFLICT OF INTEREST STATEMENT

389 The authors declare that the research was conducted in the absence of any commercial or financial
 390 relationships that could be construed as a potential conflict of interest.

DATA AVAILABILITY STATEMENT

391 The original contributions presented in the study are included in the article; further inquiries can be directed
 392 to the corresponding author.

REFERENCES

- 393 Barata, C., Ruela, M., Francisco, M., Mendonça, T., and Marques, J. S. (2013). Two systems for the
 394 detection of melanomas in dermoscopy images using texture and color features. *IEEE Systems Journal*
 395 8, 965–979
- 396 Bdair, T. M., Navab, N., and Albarqouni, S. (2021). Peer learning for skin lesion classification. *CoRR*
 397 abs/2103.03703
- 398 Bi, L., Kim, J., Ahn, E., and Feng, D. (2017). Automatic skin lesion analysis using large-scale dermoscopy
 399 images and deep residual networks. *arXiv preprint arXiv:1703.04197*
- 400 Cao, Z., Sun, C., Wang, W., Zheng, X., Wu, J., and Gao, H. (2021). Multi-modality fusion learning
 401 for the automatic diagnosis of optic neuropathy. *Pattern Recognition Letters* 142, 58–64. doi:<https://doi.org/10.1016/j.patrec.2020.12.009>
- 403 Capdehourat, G., Corez, A., Bazzano, A., Alonso, R., and Musé, P. (2011). Toward a combined tool
 404 to assist dermatologists in melanoma detection from dermoscopic images of pigmented skin lesions.
 405 *Pattern Recognition Letters* 32, 2187–2196
- 406 Celebi, M. E., Kingravi, H. A., Uddin, B., Iyatomi, H., Aslandogan, Y. A., Stoecker, W. V., et al. (2007). A
 407 methodological approach to the classification of dermoscopy images. *Computerized Medical imaging
 408 and graphics* 31, 362–373

- 409 Chen, J., Ying, H., Liu, X., Gu, J., Feng, R., Chen, T., et al. (2021a). A transfer learning based super-
410 resolution microscopy for biopsy slice images: The joint methods perspective. *IEEE/ACM Transactions
411 on Computational Biology and Bioinformatics* 18, 103–113. doi:10.1109/TCBB.2020.2991173
- 412 Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., and Yuille, A. L. (2017). Deeplab: Semantic
413 image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE
414 transactions on pattern analysis and machine intelligence* 40, 834–848
- 415 Chen, T., Liu, X., Feng, R., Wang, W., Yuan, C., Lu, W., et al. (2021b). Discriminative cervical lesion
416 detection in colposcopic images with global class activation and local bin excitation. *IEEE Journal of
417 Biomedical and Health Informatics*, 1–1doi:10.1109/JBHI.2021.3100367
- 418 Chollet, F. (2017). Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the
419 IEEE conference on computer vision and pattern recognition*. 1251–1258
- 420 Codella, N., Cai, J., Abedini, M., Garnavi, R., Halpern, A., and Smith, J. R. (2015). Deep learning,
421 sparse coding, and svm for melanoma recognition in dermoscopy images. In *International workshop on
422 machine learning in medical imaging* (Springer), 118–126
- 423 Codella, N. C., Gutman, D., Celebi, M. E., Helba, B., Marchetti, M. A., Dusza, S. W., et al. (2018).
424 Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on
425 biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In *2018 IEEE
426 15th International Symposium on Biomedical Imaging (ISBI 2018)* (IEEE), 168–172
- 427 Dai, J., He, K., and Sun, J. (2016). Instance-aware semantic segmentation via multi-task network cascades.
428 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3150–3158
- 429 Datta, S. K., Shaikh, M. A., Srihari, S. N., and Gao, M. (2021). Soft-attention improves skin cancer
430 classification performance
- 431 Díaz, I. G. (2017). Incorporating the knowledge of dermatologists to convolutional neural networks for the
432 diagnosis of skin lesions. *arXiv preprint arXiv:1703.01976*
- 433 Feng, R., Liu, X., Chen, J., Chen, D. Z., Gao, H., and Wu, J. (2021). A deep learning approach
434 for colonoscopy pathology wsi analysis: Accurate segmentation and classification. *IEEE Journal of
435 Biomedical and Health Informatics* 25, 3700–3708. doi:10.1109/JBHI.2020.3040269
- 436 Ganster, H., Pinz, P., Rohrer, R., Wildling, E., Binder, M., and Kittler, H. (2001). Automated melanoma
437 recognition. *IEEE transactions on medical imaging* 20, 233–239
- 438 Gessert, N., Sentker, T., Madesta, F., Schmitz, R., Kniep, H., Baltruschat, I., et al. (2020). Skin
439 lesion classification using cnns with patch-based attention and diagnosis-guided loss weighting. *IEEE
440 Transactions on Biomedical Engineering* 67, 495–503. doi:10.1109/TBME.2019.2915839
- 441 González-Díaz, I. (2019). Dermaknet: Incorporating the knowledge of dermatologists to convolutional
442 neural networks for skin lesion diagnosis. *IEEE Journal of Biomedical and Health Informatics* 23,
443 547–559. doi:10.1109/JBHI.2018.2806962
- 444 He, K., Zhang, X., Ren, S., and Sun, J. (2015). Spatial pyramid pooling in deep convolutional networks for
445 visual recognition. *IEEE transactions on pattern analysis and machine intelligence* 37, 1904–1916
- 446 He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In
447 *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778
- 448 Hu, J., Shen, L., and Sun, G. (2018). Squeeze-and-excitation networks. In *Proceedings of the IEEE
449 conference on computer vision and pattern recognition*. 7132–7141
- 450 Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional
451 networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4700–
452 4708

- 453 Kittler, H., Pehamberger, H., Wolff, K., and Binder, M. (2002). Diagnostic accuracy of dermoscopy. *The
454 lancet oncology* 3, 159–165
- 455 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional
456 neural networks. In *Advances in neural information processing systems*. 1097–1105
- 457 Lai, Z. and Deng, H. (2018). Medical image classification based on deep features extracted by deep model
458 and statistic feature fusion with multilayer perceptron? *Computational Intelligence and Neuroscience*
459 2018
- 460 Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., et al. (2016). Ssd: Single shot multibox
461 detector. In *European conference on computer vision* (Springer), 21–37
- 462 Long, J., Shelhamer, E., and Darrell, T. (2015). Fully convolutional networks for semantic segmentation.
463 In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3431–3440
- 464 Matsunaga, K., Hamada, A., Minagawa, A., and Koga, H. (2017). Image classification of melanoma, nevus
465 and seborrheic keratosis by deep neural network ensemble. *arXiv preprint arXiv:1703.03108*
- 466 Menegola, A., Tavares, J., Fornaciali, M., Li, L. T., Avila, S., and Valle, E. (2017). Recod titans at isic
467 challenge 2017. *arXiv preprint arXiv:1703.04819*
- 468 Mishra, N. K. and Celebi, M. E. (2016). An overview of melanoma detection in dermoscopy images using
469 image processing and machine learning. *arXiv preprint arXiv:1601.07843*
- 470 Moura, N., Veras, R., Aires, K., Machado, V., Silva, R., Araújo, F., et al. (2019). Abcd rule and pre-trained
471 cnns for melanoma diagnosis. *Multimedia Tools and Applications* 78, 6869–6888
- 472 Myronenko, A. (2018). 3d mri brain tumor segmentation using autoencoder regularization. In *International
473 MICCAI Brainlesion Workshop* (Springer), 311–320
- 474 Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). You only look once: Unified, real-time
475 object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
476 779–788
- 477 Ronneberger, O., Fischer, P., and Brox, T. (2015). U-net: Convolutional networks for biomedical
478 image segmentation. In *International Conference on Medical image computing and computer-assisted
479 intervention* (Springer), 234–241
- 480 Roychowdhury, S., Koozekanani, D. D., and Parhi, K. K. (2015). Iterative vessel segmentation of fundus
481 images. *IEEE Transactions on Biomedical Engineering* 62, 1738–1749
- 482 Shin, H.-C., Roth, H. R., Gao, M., Lu, L., Xu, Z., Nogues, I., et al. (2016). Deep convolutional neural
483 networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning.
484 *IEEE transactions on medical imaging* 35, 1285–1298
- 485 Siegel, R. L., Miller, K. D., and Jemal, A. (2016). Cancer statistics, 2016. *CA: a cancer journal for
486 clinicians* 66, 7–30
- 487 Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image
488 recognition. *arXiv preprint arXiv:1409.1556*
- 489 Stolz, W. (1994). Abcd rule of dermatoscopy: a new practical method for early recognition of malignant
490 melanoma. *Eur. J. Dermatol.* 4, 521–527
- 491 Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. A. (2017). Inception-v4, inception-resnet and the
492 impact of residual connections on learning. In *Thirty-first AAAI conference on artificial intelligence*
- 493 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., et al. (2015). Going deeper with
494 convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1–9
- 495 Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2016). Rethinking the inception
496 architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern
497 recognition*. 2818–2826

- 498 Xiao, J., Xu, H., Gao, H., Bian, M., and Li, Y. (2021). A weakly supervised semantic segmentation
499 network by aggregating seed cues: The multi-object proposal generation perspective. *ACM Transactions*
500 *on Multimedia Computing Communications and Applications* 17, 1–19
- 501 Xie, F., Fan, H., Li, Y., Jiang, Z., Meng, R., and Bovik, A. (2016). Melanoma classification on dermoscopy
502 images using a neural network ensemble model. *IEEE transactions on medical imaging* 36, 849–858
- 503 Xie, Y., Zhang, J., Xia, Y., and Shen, C. (2020). A mutual bootstrapping model for automated skin lesion
504 segmentation and classification. *IEEE Transactions on Medical Imaging* 39, 2482–2493. doi:10.1109/
505 TMI.2020.2972964
- 506 Yang, X., Zeng, Z., Yeo, S. Y., Tan, C., Tey, H. L., and Su, Y. (2017). A novel multi-task deep learning
507 model for skin lesion segmentation and classification. *arXiv preprint arXiv:1703.01025*
- 508 Yu, L., Chen, H., Dou, Q., Qin, J., and Heng, P.-A. (2016). Automated melanoma recognition in dermoscopy
509 images via very deep residual networks. *IEEE transactions on medical imaging* 36, 994–1004
- 510 Yu, Z., Jiang, X., Zhou, F., Qin, J., Ni, D., Chen, S., et al. (2018). Melanoma recognition in dermoscopy
511 images via aggregated deep convolutional features. *IEEE Transactions on Biomedical Engineering* 66,
512 1006–1016
- 513 Zhang, J., Xie, Y., Xia, Y., and Shen, C. (2019). Attention residual learning for skin lesion classification.
514 *IEEE Transactions on Medical Imaging* 38, 2092–2103. doi:10.1109/TMI.2019.2893944
- 515 Zunair, H. and Hamza, A. B. (2020). Melanoma detection using adversarial training and deep transfer
516 learning 65, 135005. doi:10.1088/1361-6560/ab86d3