

Fusing Fine-tuned Deep Features for Skin Lesion Classification

Amirreza Mahbod^{a,b,*}, Gerald Schaefer^c, Isabella Ellinger^a, Rupert Ecker^b,
Alain Pitiot^d, Chunliang Wang^e

^a*Institute of Pathophysiology and Allergy Research, Medical University of Vienna, Vienna, Austria*

^b*Research and Development Department of TissueGnostics GmbH, Vienna, Austria*

^c*Department of Computer Science, Loughborough University, Loughborough, United Kingdom*

^d*Laboratory of Image and Data Analysis, Ilixa Limited, Nottingham, United Kingdom*

^e*School of Technology and Health, KTH Royal Institute of Technology, Stockholm, Sweden*

Abstract

Malignant melanoma is one of the most aggressive forms of skin cancer. Early detection is important as it significantly improves survival rates. Consequently, accurate discrimination of malignant skin lesions from benign lesions such as seborrheic keratoses or benign nevi is crucial, while accurate computerised classification of skin lesion images is of great interest to support diagnosis. In this paper, we propose a fully automatic computerised method to classify skin lesions from dermoscopic images. Our approach ensembles deep features from several well-established convolutional neural networks (CNNs) at different abstraction levels in combination with a support vector machine classifier to distinguish malignant melanomas from benign lesions. Importantly, the CNNs are pre-trained on a common natural image database and then fine-tuned on a limited set of dermoscopic skin lesion images. Finally, prediction probability classification vectors obtained from different models with different training settings are fused to provide improved classification performance. Evaluated on the 600 test images of the ISIC 2017 skin lesion classification challenge, the proposed algorithm yields an area under receiver operating characteristic curve (AUC) of

*Corresponding author

Email address: amirreza.mahbod@tissuegnostics.com (Amirreza Mahbod)

87.3% for melanoma classification and an AUC of 95.5% for seborrheic keratosis classification, outperforming the top-ranked methods of the challenge while being simpler compared to them. The obtained results convincingly demonstrate our proposed approach to represent a reliable and robust method for feature extraction, model fusion and classification of dermoscopic skin lesion images.

Keywords: Dermatology, skin cancer, melanoma, dermoscopy, medical image analysis, deep learning.

1. Introduction

Malignant melanoma (MM) is a very aggressive form of skin cancer. Although occurrences of non-melanoma skin cancer are far more common (MM represents less than 5% of all skin cancers), 70% of skin cancer deaths are due to MM. 132,000 melanoma skin cancers occur globally each year [1], and both
5 incidence and mortality rates have increased throughout most of the developed world over the past 30 years [2]. Prevention as well as early detection are crucial to reverse this trend [3]. If identified early enough, skin cancer can be cured through a simple excision, while diagnosis at later stages is associated with a
10 greater risk of death - the estimated 5-year survival rate is over 95% for early stage diagnosis, but below 20% for late stage detection [4, 5].

Seborrheic keratosis (SK) is one of the most common benign skin lesions. SKs can exhibit wide variations in its clinical features, and some types of SK resemble melanomas or other skin tumors. Moreover, melanomas may appear adjacent to or within SKs. Therefore, it can be difficult to distinguish melanomas
15 from SKs. Likewise, benign nevi (BN), which are pigmented skin growths with no current signs of pathology, can appear similar to melanomas, while patients with numerous nevi have a significantly higher risk of developing skin cancer [6].

Pathological analysis of a biopsy specimen enables differentiation between
20 different types of skin lesions with certainty, but this type of analysis is both time and labour intensive and not always possible. Dermoscopy, in contrast, is a non-invasive, microscopy-based diagnostic method, which allows for en-

hanced visualisation of the internal structures of lesions [7]. When performed by well-trained and experienced dermatologists, dermoscopy supports a diagnostic accuracy of about 80% [8, 9] and leads to a reduced number of unnecessary excisions [6]. However, visual inspection of dermoscopic images by dermatologists requires training and experience, since the diagnostic accuracy achieved by non-experts using dermoscopy is no better than with the unaided eye [10].

Despite the definition of commonly employed diagnostic schemes such as the ABCD rule [11] or the 7-point checklist [12], due to the difficulty and subjectivity of human interpretation as well as the variety of lesions and confounding factors encountered in practice, computerised analysis of dermoscopic images has become an important research area to support diagnosis [13, 14]. Conventional computer-aided methods for dermoscopic lesion classification typically involve three main stages: segmenting the lesion, extracting hand-crafted image feature from the lesion area and its border, and classification [15, 16]. In addition, often extensive pre-processing is involved to improve image contrast [16, 17], perform white balancing based on colour constancy algorithms [18], apply colour normalisation [19] or calibration [20], colour space transformation [16], illumination correction [16], or remove image artefacts such as hairs [13, 21] or bubbles [13].

Accurate segmentation of the lesion area is considered important, since the shape of the lesion gives crucial clues for diagnosis, while the subsequent processing steps rely on a precise division between lesion and skin areas. A variety of segmentation algorithms have been developed for border detection [22, 16] including thresholding-based methods [23], region merging approaches [24], clustering techniques [25], active contours [26] and machine learning techniques such as artificial neural networks [17]. Based on the segmented lesion area, domain specific features are then extracted. These features can relate to lesion type (primarily morphological features), lesion configuration (secondary morphological features), colour, shape, texture and lesion border [15, 27]. In order to select the most relevant features and to reduce the dimensionality of the feature space, a number of feature selection methods can be utilised, which in turn can lead to improved classification performance and lower training and testing time [28].

In supervised approaches, where the ground truth of a subset of data is
55 available, the selected features together with the corresponding labels are used
to train a classifier (such as support vector machines (SVMs), random forest
classifiers or multi-layer perceptrons (MLP) [28, 4, 29]), while the extracted
features can also be utilised in unsupervised learning approaches such as k -
means clustering or principle component analysis [30]. For both supervised and
60 unsupervised approaches, the trained model can then be employed for classifying
new skin lesion images. An overview of classifiers that have been used for skin
lesion classification can be found in [28] and shows that a SVM is a common
choice due to its relatively good generalisation properties [28], the possibility
to incorporate of kernel functions to simplify and enhance the classification
65 of non-linear feature distributions in high-dimensional spaces, and competitive
classification performance compared to the more complex classifiers [28, 31].

The main drawback of conventional approaches is a lack of generalisation
capability due to high variations in dermoscopic images, different artefacts and
insufficient training data. Variations in dermoscopic images are due to differ-
70 ent zooming configurations, lighting conditions, instruments or operators, while
common artifacts in dermoscopic images include not just skin hair and bubbles
but also, among others, dark corners/borders, light reflections or shadows, skin
lines, ruler or calibration chart artefacts or ink markings, which can lead to
failures of segmentation algorithms, changes in extracted image features and
75 consequently a negative effect on classification accuracy [32, 16].

Deep neural networks (DNNs), in particular convolutional neural networks
(CNNs), are superior to other methods for tasks such as object detection and
natural image classification [33, 34]. To achieve high accuracy, well established
CNN architectures such as AlexNet [33], VGGNet [34] and ResNet [35] are
80 typically trained on large image databases such as ImageNet [36] which comprise
millions of heterogenous images. However in medicine, access to validated data is
heavily restricted and expensive to obtain, which makes training such networks
from scratch problematic [37]. One way to address this problem is to use transfer
learning, which employs a pre-trained network (i.e., one trained on other tasks

85 such as generic image classification) and adapt it to the problem at hand. This pre-training allows the network to identify useful features even when training samples are limited [27].

In medical image analysis, transfer learning has been used for a variety of applications including radiology, cardiology, ultrasound imaging, gastroenterol-
90 ogy, retinopathy, microscopic imaging as well as dermoscopy [38, 27]. So far, mainly two different approaches of transfer learning were used for medical image analysis and in particular for skin lesion classification [38]. On the one hand, pre-trained CNNs were used as feature generators. In this setting, images are fed to pre-trained models and deep features extracted from a certain
95 fully connected (FC) layer or convolution layer. The generated extracted features are then used to train a classical classifier such as an SVM [39, 27]. In some extended studies, these features were encoded to more invariant and discriminative representations [40] or combined with other hand-crafted feature descriptors [41, 29] to enhance classification performance. On the other hand,
100 trained models can be adapted to the problem at hand by fine-tuning. To fine-tune deep models, FC layers of the pre-trained networks were typically replaced by one or more new logistic layers and then the networks re-trained to adapt the weights of the newly added layers for classifying skin lesions [42]. The pre-trained models used in both approaches for skin lesion classification
105 varied in different studies and include AlexNet [39, 32, 41], VGG16 [32, 27, 43], VGG19 [32], GoogleNet [44, 45, 43], ResNet-50 [43, 46, 47], ResNet-101 [48], ResNet-152 [49, 50] Inception-v3 [42, 51], Inception-v4 [48, 49], variations of DenseNets [31, 49], SeNets [31, 50] and PolyNets [31]. Moreover, ensembles of fine-tuned deep networks [48, 46] and fusing outputs of classical and deep
110 models [41, 40] were utilised to boost classification performance.

In this paper, in contrast to former studies, we utilise both schemes of transfer learning in one single approach. We exploit several well-known CNNs pre-trained on ImageNet and fine-tune them on a limited dataset of dermoscopic lesion images. We ensemble deep features, that is the outputs of the last few
115 fully-connected layers, in an SVM classifier that then gives the classification of

the lesion type. Unlike previous works using deep features for skin lesion classification [27, 29, 39], which were limited to exploit specific network architectures or using specific layers for extracting features, in our approach, we hypothesise that extracting features from different layers of different abstraction levels and from different deep models can improve the classification results. More importantly, we fine-tune pre-trained networks for feature extraction to achieve better classification performance for skin lesion categorisation. Moreover, compared to conventional methods and some fusing approaches, we avoided using extensive pre-processing steps, lesion segmentation masks or engineered hand-crafted feature descriptors to potentially increase the generalisation ability and at the same time its adaptability to be extended for other classification tasks. Finally, we perform a thorough investigation of the performance of each component of our proposed method to justify our approach and to provide a useful guideline for further developments of CNN-based algorithms for skin lesion classification.

2. Materials and Methods

Our proposed skin lesion classification method consists of the following major steps: image pre-processing, deep neural network fine-tuning and feature extraction to train a SVM classifier, and ensembling the model outputs. In the following, we describe the utilised datasets, and cover in detail each of the stages of our approach.

2.1. Dataset

We used the training, validation and test images of the ISIC 2016 challenge [52] as well as the training and validation images of the ISIC 2017 challenge [53]. These probably represent the most challenging skin lesions datasets that are publicly available to date for ternary skin lesion classification. From these two datasets, 2,187 training images were extracted for training which included 441 MMs, 296 SKs and 1450 BN images. We tested our trained model on the 600 images that comprise the test set of the ISIC 2017 challenge and which

were not used in the training phase. All training and test images are 24-bit
145 RGB images of various sizes (ranging from $1,022 \times 767$ to $6,748 \times 4,499$ pixels), perspectives, and lighting conditions, while a significant number of images contained various artefacts.

2.2. Pre-processing

In our proposed pipeline, we aimed to keep the pre-processing steps to a
150 minimum to support better generalisation ability when tested on other datasets. Three pre-processing steps were applied in our approach where only one was task specific (related to skin lesion classification) while the other two were standard pre-processing steps to prepare the images before feeding them to deep networks.

2.2.1. Colour standardisation

155 As the images were acquired under different lighting conditions and with different devices, we performed colour normalisation using the gray world colour constancy algorithm, which has been reported to support improved skin lesion classification [18, 46] .

2.2.2. Normalisation

160 In order to utilise pre-trained deep networks, a common normalisation technique is to subtract the mean RGB value of the ImageNet dataset from all training and test images [33]. Other approaches were also tested, including subtracting mean RGB values computed over each individual image and subtracting mean RGB values computed over the whole training dataset from all
165 training and test images as suggested in [39] and [40].

2.2.3. Resizing

Since all pre-trained networks used in our implementation expect the input images to be of the same size defined during training, we resized all images to the appropriate size (227×227 and 224×224 pixels) using bicubic interpolation.
170 For non-square images, the aspect ratio was changed during this resizing step.

2.3. Pre-trained deep learning models and fine-tuning

In order to extract optimised features from the images, we used well-established CNN architectures, namely AlexNet [33], VGGNet [34] and two variations of ResNet [35] which have shown excellent performance in previous classification tasks such as the Image Large Scale Visual Recognition Challenge (VGGNet was the runner-up of the challenge in 2014, while AlexNet and ResNet were the winners of the challenge in 2012 and 2015, respectively [54]). While AlexNet has a well-established architecture with 5 convolutional layers and 3 FC layers, the original implementations of VGGNet and ResNet come with several variations. In our work, we used VGG16, which has 16 weight layers, 13 convolutional and 3 FC layers as well as ResNet-18 and ResNet-101 which exhibit different depths. In general, ResNet’s architecture consists of special building blocks called residual blocks and one FC layer on top which performs the classification.

In order to extract features from these DNNs, one approach is to simply run the images through the pre-trained networks and take the output of the FC layers as was done in some previous works [27, 39]. However, we hypothesise that fine-tuning of pre-trained networks using skin lesion images should lead to higher quality features from the images.

Fine-tuning of the selected networks was performed as follows. First, the last FC layer and the output layer of all pre-trained networks were removed and replaced by two new FC layers with 64 and 3 nodes to solve the ternary (MM/SK/BN) classification problem, as shown in Fig. 1 with ResNet as an example. The weights of the added fully connected layers were randomly generated from a Gaussian distribution with zero mean and standard deviation of 0.01. In order to prevent overfitting and to speed up the training, we froze the weights of the initial layers of the deep models. For AlexNet and VGG16, we froze the initial layers up to the 4-th and 10-th convolutional layers, while we froze the layers up to the 4-th and 30-th residual blocks for ResNet-18 and ResNet-101, respectively.

We tested different optimisers with regularisation terms for the loss function in order to perform fine-tuning. In particular, we utilised stochastic gradient

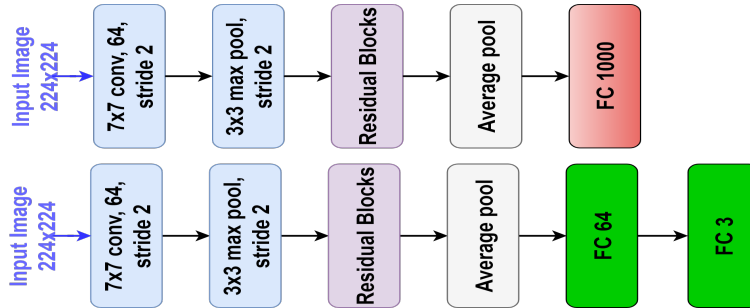


Figure 1: Generic structure of the original ResNet (top) and the modified architecture adapted for fine-tuning in our proposed approach (bottom). The final FC layer of the original architecture (red block) is replaced by two FC layers (green blocks).

descent with momentum (SGDM) [55, 56], root mean square propagation (RMSProp) [57] and adaptive moment estimation (Adam) [58] optimisers in our experiments.

205 The SGDM optimiser updates the weights and biases of the network in each iteration in order to minimise the error (i.e., minimise the loss function output) by taking small steps in the negative direction of the gradient. We used the momentum term in order to prevent oscillations along the steepest descent path. The general SGDM term employed in our approach is

$$\theta_{i+1} = \theta_i - \alpha \nabla E_R(\theta_i) + \gamma(\theta_i - \theta_{i-1}), \quad (1)$$

210 where θ is the parameter vector of the network, i represents the iteration number, α is the learning rate, E_R indicates the loss function, and γ is the momentum term which determines the contribution of the previous gradient step to the current iteration. We used the cross-entropy loss function in the optimisation process as

$$E(\theta) = - \sum_l \sum_{m=1}^k t_{lm} \ln(y_m(x_l, \theta)) \quad (2)$$

215 and

$$E_R(\theta) = E(\theta) + \lambda \Omega(w), \quad (3)$$

where k is the number of classes, t_{lm} indicates that the l -th sample belongs to the m -th class and $y_m(x_l, \theta)$ is the network output of the l -th sample. The

added term in Equ. (3) is the regularisation term, where w is the weight factor, λ is the regularisation factor coefficient and $\Omega(w)$ is the regularisation function
 220 defined as

$$\Omega(w) = \frac{1}{2}w^T w. \quad (4)$$

RMSProp minimises the loss function based on

$$\theta_{i+1} = \theta_i - \frac{\alpha \nabla E(\theta_i)}{\sqrt{v_i} + \epsilon}, \quad (5)$$

where v_i is

$$v_i = \beta_2 v_{i-1} + (1 - \beta_2) [\nabla E(\theta_i)]^2, \quad (6)$$

and β_2 is the decay rate which needs to be set as an hyperparameter (ϵ is a very small number and prevents division by zero).

225 While SGDM uses a single learning rate for updating the parameters, RMSProp tries to adapt the learning rate for different parameters based on the loss function being optimised. In the RMSProp optimisation approach, the learning rate of the parameters with large gradients will be reduced and the learning rate of the parameters with relatively small gradients will be increased.

230 The Adam optimiser, similar to RMSProp, adapts the learning rate for optimisation but with a momentum term as

$$\theta_{i+1} = \theta_i - \frac{\alpha m_i}{\sqrt{v_i} + \epsilon}, \quad (7)$$

where m is

$$m_i = \beta_1 m_{i-1} + (1 - \beta_1) \nabla E(\theta_i), \quad (8)$$

and v is as in Equ. (6). β_1 is the gradient decay factor, another hyperparameter. The added momentum term in Adam controls the parameter updates. If the
 235 gradients over many iterations are similar, the updates will be larger and if the gradient varies a lot (e.g. through noise) then the updates will be small.

In our experiments, we set the initial learning rate to 0.001 for the SGDM optimiser and to 0.0001 for RMSProp and Adam, but we kept the learning rate of the new FC layers 10 times bigger compared to all other learnable layers.
 240 Weight decay was set to 0.0001 and the momentum term for SGDM was set to

0.9. β_1 and β_2 in Equ. (6) and Equ. (8) were set to 0.9 and 0.999, respectively. For AlexNet, the batch size was set to 128, for VGG16 to 32, and for the ResNets to 16 in order to fit into GPU memory. The learning rate was dropped by a factor of 10 after 6 epochs and we retrained all models for 12 epochs.

245 In order to prevent overfitting of the networks to our limited training dataset, we artificially increased the training size by data augmentation. For this, we used rotation (90, 180 and 270 degrees) and horizontal flipping as main data augmentation techniques. Moreover, the images randomly underwent small changes in each iteration in the training process. These changes included random ro-
250 tations (-5 to 5 degrees), random scaling (0.9 to 1.1) and random shearing (-2 to 2 degrees). From the derived modified training data, we randomly split the dataset to 90% for training and 10% for validation.

2.4. *Ensembling deep features and fusion of networks*

The deep features are the outputs of the FC layers from the pre-trained
255 or fine-tuned DNNs. We tested two strategies to extract deep features from DNNs. The first was to use the output of only the first FC layer following the convolutional layers. The second was to concatenate the outputs of all FC layers. For the fine-tuned networks, we also included the outputs of the two added/replaced layers in the modified networks, i.e. the FC64 and FC3 outputs
260 in Fig. 1.

The extracted deep features along with the corresponding labels identifying the lesion types were used to train a ternary SVM classifier. We tested both linear and radial basis function (RBF) kernels and observed slightly better performance with the RBF kernel, similar to others [16, 31]. We therefore utilised
265 one-versus-all multi-class SVM classifier with RBF kernels in our final models. The SVM scores were mapped to probabilities using logistic regression [59], and the classification results were the probabilistic prediction vectors derived from the trained SVMs for the three different classes, which can also be used to identify the predicted lesion type. Data augmentation, similar to that employed
270 during the DNN fine-tuning step, was also performed. During the inference

stage on testing data, 8 copies of a single test image (0, 90, 180 and 270 degree rotation, with and without horizontal flipping) were fed to the pipeline. The final classification for each individual test image was based on the average probabilities of the 8 results for each model.

275 Finally, we employed an extensive yet straight-forward ensembling approach to boost our classification performance and to improve the robustness of our approach. For each architecture, we took the average over different prediction vectors which were acquired from the same model architecture, but with different training parameters. The varied parameters in the ensembling approach
280 were the normalisation technique (ImageNet mean subtraction or training mean subtraction) and the optimisers (SGDM, RMSProp and Adam). Moreover, we trained each model 3 times and took the average over the results. Hence, the final results of a single architecture (e.g. ResNet-18) were acquired from 18 different models.

285 *2.5. Evaluation*

Evaluation of the proposed method was performed by calculating the area under the receiver operating characteristics curve (AUC) which is the main evaluation metric in the ISIC 2017 challenge [53].

Since the ISIC 2017 challenge evaluation was based on two binary classification tasks (MM vs. all and SK vs. all), we converted our three elemental
290 prediction vectors to two elemental binary vectors by a one-versus-all approach. For these binary tasks, we also evaluated the results based on the accuracy at the threshold of 50%. Moreover optimal sensitivity and specificity of our best performing approach were calculated using Youden index method [60].

295 **3. Results**

The obtained results are derived from the 600 test images of the ISIC 2017 challenge. These are comprised of 117 MMs, 90 SKs, and 393 BN images not used in the training phase. All test images underwent the same pre-processing steps that were applied to the training images.

Table 1: Effects of gray world color constancy (using fine-tuned ResNet-18).

	AUC MM (%)	AUC SK (%)	average AUC (%)
no standardisation	80.23 ± 1.77	89.64 ± 0.99	84.93 ± 0.56
color constancy	83.48 ± 0.74	91.39 ± 1.33	87.44 ± 0.57

Table 2: Effects of various normalisation techniques (using fine-tuned ResNet-18).

	AUC MM (%)	AUC SK (%)	average AUC (%)
no normalisation	74.38 ± 0.19	86.00 ± 1.59	80.19 ± 0.89
ImageNet mean	83.48 ± 0.74	91.39 ± 1.33	87.44 ± 0.57
image mean	75.89 ± 0.51	83.53 ± 1.44	79.70 ± 0.97
training mean	84.36 ± 0.45	91.88 ± 0.85	88.12 ± 0.61

300 For most of the hyperparameter searches and to show the effect of the individual components of the proposed methods on the classification results, we utilise the ResNet-18 model since its single model performance is very competitive (see Table 7) and as due to its shallower depth compared to ResNet-101 its training is faster. In all experiments, we use the RMSprob optimiser, ImageNet
 305 mean subtraction, gray world normalisation and feature extraction from all FC layers, unless stated otherwise in the text.

We started our experiments by examining the effect of colour standardisation and normalising the images prior to feature extraction as described in Section 2.2. The obtained results are given in Table 1 and Table 2 where
 310 the average and standard deviation were calculated by running each setting 3 times. Since we observed better performance using ImageNet normalisation and training mean subtraction normalisation, we did not use the other settings in subsequent experiments. Similarly, as colour constancy was found to be beneficial, subsequent experiments always incorporated the colour standardisation
 315 step.

In the next experiment, we investigated the effect of optimiser on the classification performance. Table 3 shows the results of this comparison, i.e. the

Table 3: Effects of various optimisers (using fine-tune ResNet-18).

	AUC MM (%)	AUC SK (%)	average AUC (%)
SGDM	83.30 ± 0.64	91.64 ± 0.99	87.47 ± 0.81
RMSProp	83.48 ± 0.74	91.39 ± 1.33	87.44 ± 0.57
Adam	84.38 ± 0.41	91.81 ± 0.64	88.10 ± 0.50

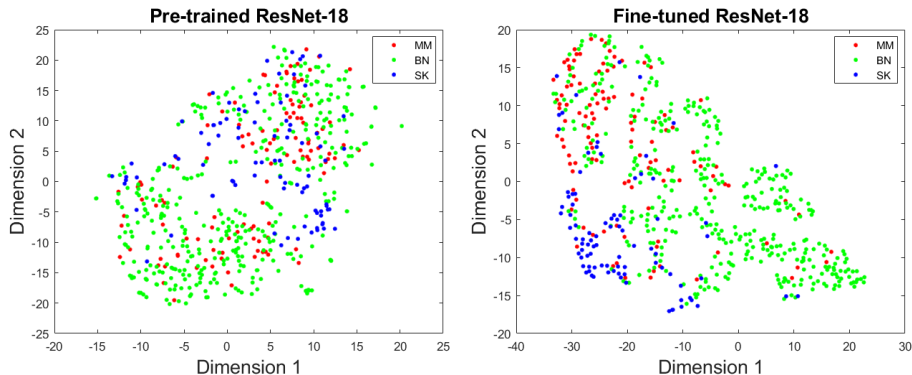


Figure 2: t-SNE visualisation of the extracted features for the pre-trained (left) and fine-tuned (right) ResNet-18 models.

results of using the SGDM, RMSProp and Adam optimiser.

In order to investigate the generalisability of the employed transfer learning approach (i.e., extracting features from the pre-trained and the fine-tuned DNNs), we performed dimensionality reduction to two dimensions using t-distributed stochastic neighbor embedding (t-SNE) [61]. This method allows to visualise the natural clusters of the high-dimensional features which we use. We used the extracted features from the first FC layer of the pre-trained network and first FC layer of the modified fine-tuned network, and utilised the Barnes-Hut Variation of t-SNE [62] to speed up the algorithm while setting the dimensionality of the principal component analysis to 50. The obtained results for pre-trained and fine-tuned ResNet-18 architectures, based on the test dataset, are shown in Fig. 2.

Moreover, we performed experiments to fine-tune ResNet-18 with the same

Table 4: Effect of weight initialisation on performance of ResNet-18 model.

	AUC MM (%)	AUC SK (%)	average AUC (%)
ImageNet	83.48	91.39	87.44
random	64.07	84.25	74.16

Table 5: Classification results from the fine-tuned networks from different abstraction levels for ResNet-18

	AUC MM (%)	AUC SK (%)	average AUC (%)
single FC	82.17	90.97	86.57
all FCs	83.48	91.39	87.44

model architecture but with random weight initialisation in order to compare the obtained performance with ImageNet weight initialisation. The same initialisation method as described in Section 2.3 was used for random weight initialisation. The results of this experiments are shown in Table 4.

335 In the next experiment, we evaluated the effect of feature extraction from different abstraction levels of the fine-tuned ResNet-18 model. Table 5 shows the obtained results and allows to compare the performance of using features from a single FC and from all FCs.

As the results confirm, there is a level of variation in all results when running 340 the experiments multiple times. Moreover, the models with different parameters (e.g., different optimisers) lead to slightly different but yet competitive classification results. Therefore, as explained in Section 2.4, to achieve more robust and improved classification performance, we took the average over 18 models of a single architecture. The results of this fusion scheme are given in Table 6 for 345 Res-18 networks. We performed the same fusion approach for the other deep models (i.e., AlexNet, VGG16 and ResNet-101).

Table 7 compares the performance of different deep feature extraction strategies and fusion schemes, showing the results obtained based on deep features from pre-trained single networks (plain AlexNet, plain VGG16, plain ResNet-

Table 6: Fusion scheme over 18 ResNet-18 models.

optimiser	normalisation	AUC MM(%)	AUC SK(%)	average AUC (%)
			(average over 3 runs)	
Adam	ImageNet mean	85.24	93.20	89.22
RMSProp	ImageNet mean	84.70	93.18	88.94
SGDM	ImageNet mean	84.18	92.85	88.52
Adam	training mean	84.28	93.23	88.76
RMSProp	training mean	85.02	93.09	89.05
SGDM	training mean	85.54	92.93	89.23
average over above models		85.65	94.04	89.85

18, and plain ResNet-101), from fine-tuned single networks (fine-tuned AlexNet, fine-tuned VGG16, fine-tuned ResNet-18m and fine-tuned ResNet-101) as well as the results obtained based on the fusion scheme of the networks. Receiver operating characteristics (ROC) curves of the fusion models (fusion of plain pre-trained networks and fusion of the fine-tuned networks) are shown in Fig. 3 and Fig. 4 for the MM and SK classification problems, respectively.

We also investigated the contribution of each single model to the final classification results. To do so, we removed one of the model at a time in the fusion scheme, calculated the resulting AUC, and report the results in Table 8.

Table 9 summarises the performance of the best performing approach of our proposed method (i.e. fusion of all fine-tuned network, the last row in Table 7) and compares it to the top three teams that participated in the ISIC 2017 challenge (ranked based on average AUC), as well as an earlier approach of our work that was submitted to the final classification phase of the ISIC 2017 challenge and that was obtained by feature extraction and combination of VGG16 and AlexNet pre-trained models.

The top-ranked approach by Matsunaga *et al.* [46] used colour constancy [18] as a main pre-processing step and a variation of fine-tuned ResNet-50 networks to obtain the final classification. Gonzalez-Diaz [47], the runner-up, performed

Table 7: Classification results from plain pre-trained networks, fine-tuned networks, and fusion of networks.

	AUC MM (%)	AUC SK (%)	average AUC (%)
Plain AlexNet	72.04	91.43	81.73
Plain VGG16	69.85	89.71	79.78
Plain ResNet-18	72.51	89.72	81.11
Plain ResNet-101	74.31	91.90	83.10
Fine-tuned AlexNet	80.31	88.49	84.40
Fine-tuned VGG16	84.16	93.51	88.83
Fine-tuned ResNet-18	85.65	94.04	89.85
Fine-tuned ResNet-101	85.54	92.24	88.89
Fusion of all pre-trained networks	73.19	93.02	83.10
Fusion of all fine-tuned networks	87.26	95.52	91.39

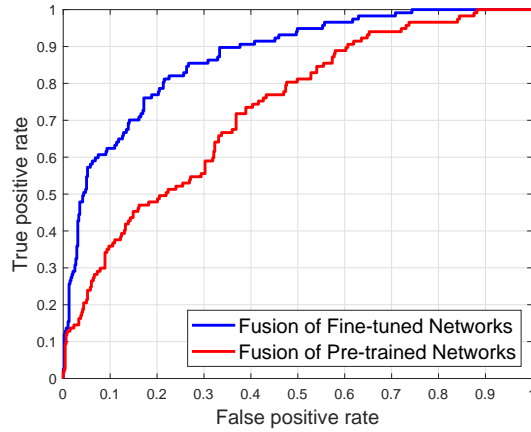


Figure 3: ROC curve of MM vs. all classification.

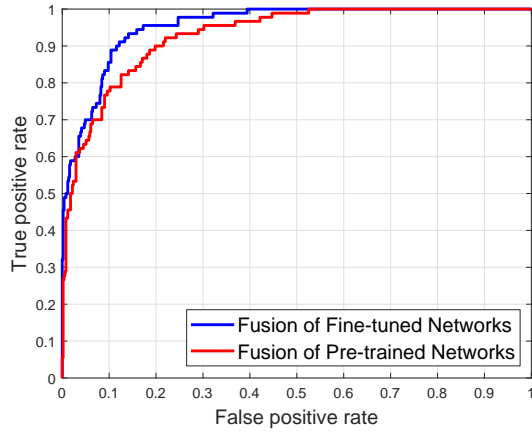


Figure 4: ROC curve of SK vs. all classification.

Table 8: Effects of removing a model in the fusion scheme.

fused networks (dropped model)	AUC MM (%)	AUC SK (%)	average AUC (%)
ResNet-18+ResNet-101+VGG16 (AlexNet)	87.01	95.36	91.18
ResNet-18+ResNet-101+AlexNet (VGG16)	87.14	94.84	90.99
ResNet-101+AlexNet+VGG16 (ResNet-18)	86.59	94.13	90.36
ResNet-18+AlexNet + VGG16 (ResNet-101)	86.76	95.43	91.09
all (none)	87.26	95.52	91.39

Table 9: Comparison of selected algorithms on the ISIC 2017 challenge.

authors	approach	AUC MM (%)	AUC SK (%)	average AUC (%)	average accuracy (%)
Matsunaga <i>et al.</i> [46]	ResNet-50 Ensemble	86.8	95.3	91.1	81.6
Gonzalez-Diaz [47]	ResNet-50 + Segmentation	85.6	96.5	91.0	84.9
Menegola <i>et al.</i> [48]	ResNet-101 + Inception-v4	87.4	94.3	90.8	88.3
Mahbod <i>et al.</i> [32]	pre-trained AlexNet + VGG	71.5	90.8	81.1	81.1
Proposed approach	see Table 7	87.3	95.5	91.4	87.7

lesion segmentation using a fully convolutional network [63] and trained a structure segmentation network to produce a set of eight global and local structures which were assumed to be beneficial for dermatologists in their routine diagnosis procedure. In a final step, the produced set of structures along with augmented data were fed to a modified ResNet-50 network for classification. Menegola *et al.* [48], whose approach was ranked third, utilised extensive data sources for fine-tuning an ensemble of seven models, six based on Inception-v4 [64] and one based on ResNet101 [35]. As the comparison shows, our proposed approach outperforms all other algorithms submitted, while it would rank 2-nd both for the MM vs. all and for the SK vs. all classification tasks among 23 participating teams in the final test phase of the ISIC 2017 challenge [53].

Figs. 5 and 6 show examples of skin lesion images correctly and incorrectly classified by our best performing approach. Moreover, in Fig. 7, the effect of fine-tuning and model fusion in terms of accuracy for MM classification is illustrated. Here, the fusion approaches from Table 7 are selected for comparison.

The algorithm was implemented in MatLab (versions 2017b and 2018a) using the MatConvNet framework [65] and the MatLab Neural Network Toolbox. All experiments were performed on a single desktop computer. For the pre-processing steps an Intel Corei5-6600k 3.50 GHz CPU was utilised. The model training was performed on a single NVIDIA GTX 1070 with 8 GB of installed memory. The training of the models took around 25 minutes, 90 minutes, 70 minutes, and 230 minutes for the AlexNet, VGG16, ResNet-18 and ResNet-101,

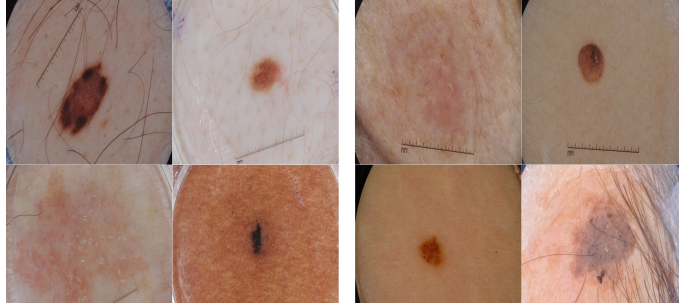


Figure 5: Examples of correctly classified images for MM vs. all (left) and SK vs. all (right) tasks.

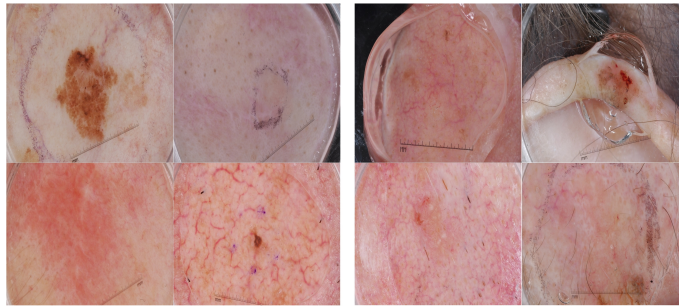


Figure 6: Examples of incorrectly classified images for MM vs. all (left) and SK vs. all (right) tasks.

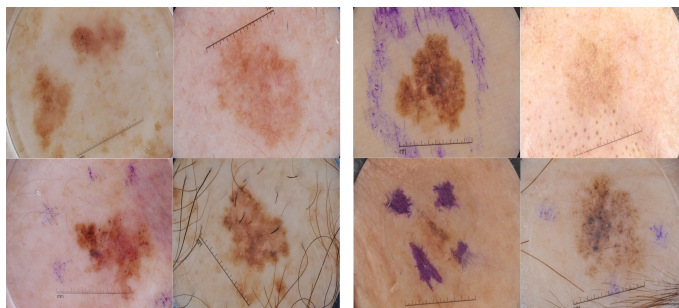


Figure 7: Comparison of different fusion approaches - fusion of plain pre-trained networks, and fusion of fine-tuned networks - for MM classification: MM examples that are correctly classified by both fusion approaches (left) . Challenging MM examples that are only correctly classified by fusion of fine-tuned networks but not by fusion of pre-trained networks (right).

respectively.

4. Discussion

The main contribution of our approach is proposing a hybrid DNN method for skin lesion classification by extracting deep features through multiple DNNs from lesion images and ensembling features in an SVM classifier that yields very accurate results without requiring extensive pre-processing or segmentation of the lesion area. By transferring deep features that were trained on a large image database of 1.4 million natural images and fine-tuning them on a relatively small skin lesion dataset, we show that it is feasible to train a reliable DNN-based classifier on a small number of domain specific sample images.

The results in Table 1 and Table 2 show the effects of pre-processing schemes on the classification results. From Table 1, it can be seen that a colour constancy algorithm can improve the performance and we hence used the colour corrected images in the remainder of the experiments. Table 2 shows the effects of different normalisation approaches to prepare the images before feeding them to the selected deep models. Among these normalisation techniques, ImageNet mean subtraction and training mean subtraction delivered better results compared to no normalisation and per image mean normalisation. Instead of choosing one of them (which delivers slightly better performance), we used both of them in our ensembling approach leads to an improvement in the classification results.

In order to validate the generalisability of the extracted features from the pre-trained and fine-tuned networks visually, we mapped the high-dimensional feature maps to two dimensions as shown in Fig. 2. As can be seen, the extracted features from the plain pre-trained ResNet-18 model is distinguishable between the different classes to some extent even without any training. Hence, it can be inferred that ImageNet features are indeed well-generalised to our ternary classification task for skin lesions. It can further be observed that, although not completely separable, by fine-tuning the network using only a limited dataset, the three skin lesion classes become more distinguishable. While Fig. 2 illus-

420 trates the applicability of the extracted features initiated by ImageNet weights visually, Table 4 confirms this quantitatively. As the results demonstrate, ImageNet weight initialisation clearly yields better performance compared to random weight initialisation which is in agreement with former studies [38].

The results in Table 5 suggest that ensembling deep features from all FC
425 layers in an SVM classifier delivers better performance compared to extracting features only from the first FC layer. Different FC layers are often thought to represent different levels of abstraction. Hence, our data suggest that combining features of different abstraction levels leads to improved classification accuracy.

In order to improve the robustness of the model as well as the classification
430 performance, we fused the probabilities of 18 different models from a single architecture as shown in Table 6. As the results clearly show, averaging over the models' outputs yields better performance compared to individual models. Moreover, it reduces the chances of degradation in results which can be caused by random weight initialisation or other factors.

435 From Table 7, we can observe that the performance of the SVM classifier when trained on features from fine-tuned networks is better compared to pre-trained networks for skin lesion classification, which is in agreement with our hypothesis. Comparing the results from different architectures shows that, although all single models delivers quite impressive classification performance, the
440 results of the ResNet-18 model are slightly better compared to the other models. This is probably because of the shallower depth of ResNet-18 compared to the depth of ResNet-101. While generally ResNet-101 should deliver better performance [35], our training data size is relatively small, and the deeper model thus likely overfits to our limited dataset while the shallower network shows a
445 better generalisation ability under these circumstances. Compared to AlexNet and VGG16, although ResNet-18 is still deeper, it consists of residual blocks which in general deliver better performance compared to regular convolutional blocks.

Fusion of deep fine-tuned DNNs is demonstrated to deliver even better re-
450 sults. Since the depths of the networks are different for AlexNet, VGG16,

ResNet18 and ResNet-101, we can anticipate that deep features from different networks may provide information complementary to each other. Moreover, from Table 8 we can see that dropping each model from the fusion scheme results in a slight degradation in classification performance. However, one can
455 use fewer networks in order to reduce the computational complexity with only a relatively small performance drop.

From Fig. 5 it can be seen that even challenging lesion images are correctly classified, while instances where an incorrect classification is obtained often include samples where the lesion is difficult to make out as illustrated in Fig. 6.
460 It can also be observed from Fig. 7 that more challenging examples with vague lesion borders, low contrast and more severe artefacts can be correctly classified when fusion of fine-tuned networks is employed.

As shown in Table 9, in comparison to other methods evaluated on the same dataset, our best performing approach delivers better performance compared
465 to the ISIC2017 competition winner and clearly outperforms the results of our earlier submission to the contest. However, the DNN models in our methods have lower complexity compared with those of the top 3 teams in Table 9, and were not trained on extensive external data sources. Direct comparison of the methodologies of the algorithms is challenging since different teams implemented
470 different pre-processing steps and used various training schemes. Moreover, our algorithm can be easily used for other classification task with minimal changes in the models.

The last column of Table 9 shows the average accuracy of our best performing approach and other state-of-the-art algorithms. It should be noted that the
475 accuracy numbers in this column are derived with mapping the score vectors to binary numbers using a probability of 50% which may not be the optimal thresholding. Optimal thresholding can be derived from the ROC curve of our best performing approach. Using Youden index method, our best performing approach yields a sensitivity of 81.20% and a specificity of 78.47% for MM vs.
480 all classification. Likewise, a sensitivity of 93.33 % and a specificity of 85.88% can be driven for SK vs. all classification from the ROC curve. However, by

considering the clinical importance of not missing any MM lesions, it is possible to choose a threshold from ROC curve that improves the sensitivity of the MM vs. all classification at the expense of reduced specificity. From the ROC
485 curve of MM vs. all, a sensitivity of 85%, 90% and 95% can be reached with corresponding specificity of 73.29%, 62.32% and 44.72%, respectively.

As stated in Section 2.1, the training images in the three classes are not well-balanced as there are relatively few MM and SK images in comparison to BN lesions. While it is common practice to balance a dataset in such cases using
490 e.g. boot strapping, class-balanced cost functions or through resampling, we have not gained any improvement in performance by balancing the dataset (we performed resampling of the minority classes to deal with class imbalance), while the training time drastically increased. This appears to confirm experiments on weighting strategies reported in [48]. We there do not explicitly address class
495 imbalance in our approach, nor do [46, 18, 48] i.e. the three top teams of the ISIC contest.

There are some limitations of our current approach that can be explored in future work. First, even though we show that fusing deep features from different DNNs can improve the classification accuracy, the number of networks
500 investigated is limited. Extending this study by incorporating other DNN architectures, such as GoogleNet [67] or DensNet [68], may result in further improvements. Moreover, ensembling hand-crafted feature descriptors as used in conventional methods alongside proposed fused deep features could lead to better classification performance [69, 41], but also increases the complexity of the
505 method. Second, the employed training data are limited. The amount of training data is important for appropriately training or fine-tuning DNNs. Hence, having access to additional reliable skin lesion data sources can lead to better results. Third, using pre-trained networks for skin lesion classification requires the images to be resized to a certain dimension that is pre-defined for other
510 image classification tasks. Some valuable information may be lost during the downsampling step. Although in some works, images were resized to higher resolutions (e.g. 339×339 pixels in [39], 448×448 pixels in [42] and up to

512 × 512 in [40]), they are still significantly smaller compared to their original sizes. However, these input sizes are still significantly bigger compared to our
515 approach and they may hence capture more useful information. Finally, using more extensive pre-processing steps or data augmentation techniques might further improve the classification performance. [16, 19, 20].

5. Conclusions

In this paper, a fully automatic computerised method with minimal pre-
520 and post-processing operations is proposed for accurate skin lesion classification. The proposed algorithm ensembles deep features from multiple pre-trained and fine-tuned DNNs at multiple abstraction levels and fuses the prediction probability vectors of different models. The obtained results show that such fusion of features provides better discrimination ability and is complementary to the
525 individual networks. The general performance of the proposed method is competitive with other state-of-the-art algorithms, while the generalisation ability of the proposed approach for other medical imaging classification tasks is subject for future work.

Acknowledgments

530 This work was supported by the European Union Horizon 2020 Research and Innovation Program ("CaSR Biomedicine", 675228). The authors appreciate the help of TissueGnostics Support team <http://tissuegnostics.com/en/> for their valuable comments and feedback. Moreover, we would like to thank Prof. Örjan Smedby since part of this study was conducted in his research group.

535 References

References

- [1] WHO, Ultraviolet radiation and the INTERSUN Programme (Data Accessed May 11, 2018).
URL <http://www.who.int/uv/faq/skincancer/en/index1.html>

- 540 [2] Z. Apalla, D. Nashan, R. B. Weller, X. Castellsagué, Skin Cancer: Epidemiology, Disease Burden, Pathophysiology, Diagnosis, and Therapeutic Approaches, *Dermatology and Therapy* 7 (1) (2017) 5–19.
- [3] R. Shellenberger, M. Nabhan, S. Kakaraparthi, Melanoma screening: A plan for improving early detection, *Annals of Medicine* 48 (3) (2016) 142–
545 148.
- [4] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, S. Thrun, Dermatologist-level classification of skin cancer with deep neural networks, *Nature* 542 (7639) (2017) 115–118.
- [5] C. M. Balch, J. E. Gershenwald, S.-j. Soong, J. F. Thompson, M. B. Atkins,
550 D. R. Byrd, A. C. Buzaid, A. J. Cochran, D. G. Coit, S. Ding, Final version of 2009 AJCC melanoma staging and classification, *Journal of Clinical Oncology* 27 (36) (2009) 6199–6206.
- [6] L. Thomas, S. Puig, Dermoscopy, Digital Dermoscopy and Other Diagnostic Tools in the Early Detection of Melanoma and Follow-up of High-risk
555 Skin Cancer Patients., *Acta Dermato-Venereologica* 97 (2017) 14–21.
- [7] G. Argenziano, H. P. Soyer, V. De Giorgi, D. Piccolo, P. Carli, M. Delfino, *Dermoscopy: a tutorial*, EDRA, Medical Publishing & New Media, 2002.
- [8] P. Carli, E. Quercioli, S. Sestini, M. Stante, L. Ricci, G. Brunasso, V. De
560 Giorgi, Pattern analysis, not simplified algorithms, is the most reliable method for teaching dermoscopy for melanoma diagnosis to residents in dermatology, *British Journal of Dermatology* 148 (5) (2003) 981–984.
- [9] M. E. Vestergaard, P. Macaskill, P. E. Holt, S. W. Menzies, Dermoscopy compared with naked eye examination for the diagnosis of primary melanoma: a metaanalysis of studies performed in a clinical setting,
565 *British Journal of Dermatology* 159 (3) (2008) 669–676.
- [10] H. Kittler, H. Pehamberger, K. Wolff, M. Binder, Diagnostic accuracy of dermoscopy, *Lancet Oncology* 3 (3) (2002) 159–165.

- [11] W. Stolz, A. Riemann, A. B. Cagnetta, L. Pillet, W. Abmayr, D. Holzel, P. Bilek, F. Nachbar, M. Landthaler, O. Braun-Falco, {ABCD} rule of dermatoscopy: A new practical method for early recognition of malignant melanoma, *European Journal of Dermatology* 4 (7) (1994) 521–527.
- [12] G. Argenziano, G. Fabbrocini, P. Carli, V. De Giorgi, E. Sammarco, M. Delfino, Epiluminescence Microscopy for the Diagnosis of Doubtful Melanocytic Skin Lesions. Comparison of the {ABCD} Rule of Dermatoscopy and a New 7-Point Checklist Based on Pattern Analysis., *Archives of Dermatology* 134 (12) (1998) 1536–1570.
- [13] M. G. Fleming, C. Steger, J. Zhang, J. Gao, A. B. Cagnetta, C. R. Dyer, Techniques for a structural analysis of dermatoscopic imagery, *Computerized Medical Imaging and Graphics* 22 (5) (1998) 375–389.
- [14] C. Carrera, M. A. Marchetti, S. W. Dusza, G. Argenziano, R. P. Braun, A. C. Halpern, N. Jaimes, H. J. Kittler, J. Malvehy, S. W. Menzies, Validity and Reliability of Dermoscopic Criteria Used to Differentiate Nevi From Melanoma: A Web-Based International Dermoscopy Society Study, *JAMA Dermatology* 152 (7) (2016) 798–806.
- [15] M. E. Celebi, H. Kingravi, B. Uddin, H. Iyatomi, A. Aslandogan, W. V. Stoecker, R. H. Moss, A Methodological Approach to the Classification of Dermoscopy Images, *Computerized Medical Imaging and Graphics* 31 (6) (2007) 362–373.
- [16] R. B. Oliveira, E. Mercedes Filho, Z. Ma, J. P. Papa, A. S. Pereira, J. M. R. S. Tavares, Computational methods for the image segmentation of pigmented skin lesions: A review, *Computer Methods and Programs in Biomedicine* 131 (2016) 127–141.
- [17] S. M. Jaisakthi, A. Chandrabose, P. Mirunalini, Automatic Skin Lesion Segmentation using Semi-supervised Learning Technique, arXiv preprint arXiv:1703.04301.

- [18] C. Barata, M. E. Celebi, J. S. Marques, Improving dermoscopy image classification using color constancy, *IEEE Journal of Biomedical and Health Informatics* 19 (3) (2015) 1146–1152.
- [19] G. Schaefer, M. I. Rajab, M. E. Celebi, H. Iyatomi, Colour and contrast enhancement for improved skin lesion segmentation, *Computerized Medical Imaging and Graphics* 35 (2) (2011) 99–104.
- [20] H. Iyatomi, M. E. Celebi, G. Schaefer, M. Tanaka, Automated color calibration method for dermoscopy images, *Computerized Medical Imaging and Graphics* 35 (2) (2011) 89–98.
- [21] Q. Abbas, M. E. Celebi, I. F. Garcia, Hair removal methods: A comparative study for dermoscopy images, *Biomedical Signal Processing and Control* 6 (4) (2011) 395–404.
- [22] M. Celebi, H. Iyatomi, G. Schaefer, W. Stoecker, Lesion border detection in dermoscopy images, *Computerized Medical Imaging and Graphics* 33 (2) (2009) 148–153.
- [23] M. E. Celebi, Q. Wen, S. Hwang, H. Iyatomi, G. Schaefer, Lesion border detection in dermoscopy images using ensembles of thresholding methods, *Skin Research and Technology* 19 (1) (2013) e252–e258.
- [24] M. E. Celebi, H. A. Kingravi, H. Iyatomi, Y. A. Aslandogan, W. V. Stoecker, R. H. Moss, J. M. Malters, J. M. Grichnik, A. A. Marghoob, H. S. Rabinovitz, S. W. Menzies, Border detection in dermoscopy images using statistical region merging, *Skin Research and Technology* 14 (3) (2008) 347–353.
- [25] H. Zhou, G. Schaefer, A. Sadka, M. E. Celebi, Anisotropic Mean Shift Based Fuzzy C-Means Segmentation of Dermoscopy Images, *IEEE Journal of Selected Topics in Signal Processing* 3 (1) (2009) 26–34.

- [26] H. Zhou, G. Schaefer, M. E. Celebi, F. Lin, T. Liu, Gradient vector flow with mean shift for skin lesion segmentation, *Computerized Medical Imaging and Graphics* 35 (2) (2011) 121–127.
- 625 [27] A. R. Lopez, X. Giro-i Nieto, J. Burdick, O. Marques, Skin lesion classification from dermoscopic images using deep learning techniques, in: 13th IASTED International Conference on Biomedical Engineering (BioMed), IEEE, 2017, pp. 49–54.
- [28] R. B. Oliveira, J. P. Papa, A. S. Pereira, J. M. R. S. Tavares, Computational methods for pigmented skin lesion classification in images: review and future trends, *Neural Computing and Applications* 29 (3) (2018) 613–636.
- 630 [29] N. Codella, J. Cai, M. Abedini, R. Garnavi, A. Halpern, J. R. Smith, Deep Learning, Sparse Coding, and SVM for Melanoma Recognition in Dermoscopy Images, in: *International Workshop on Machine Learning in Medical Imaging*, Springer, 2015, pp. 118–126.
- [30] Z. Ma, J. M. R. S. Tavares, A Review of the Quantification and Classification of Pigmented Skin Lesions: From Dedicated to Hand-Held Devices, *Journal of Medical Systems* 39 (11) (2015) 1–12.
- 640 [31] N. Gessert, T. Sentker, F. Madesta, R. Schmitz, H. Kniep, I. Baltruschat, R. Werner, A. Schlaefer, Skin lesion diagnosis using ensembles, unscaled multi-crop evaluation and loss weighting, arXiv preprint arXiv:1808.01694.
- [32] A. Mahbod, R. Ecker, I. Ellinger, Skin Lesion Classification Using Hybrid Deep Neural Networks, arXiv preprint arXiv:1702.08434.
- 645 [33] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in Neural Information Processing Systems* 25, Curran Associates, Inc., 2012, pp. 1097–1105.
- [34] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, arXiv preprint arXiv:1409.1556.

- 650 [35] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2016, pp. 770–778.
- [36] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, ImageNet: A large-scale hierarchical image database, in: IEEE Conference on Computer
655 Vision and Pattern Recognition, IEEE, 2009, pp. 248–255.
- [37] C. N. Vasconcelos, B. N. Vasconcelos, Increasing Deep Learning Melanoma Classification by Classical And Expert Knowledge Based Image Transforms, arXiv preprint arXiv:1702.07025.
- [38] N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall,
660 M. B. Gotway, J. Liang, Convolutional Neural Networks for Medical Image Analysis: Full Training or Fine Tuning?, IEEE Transactions on Medical Imaging 35 (5) (2016) 1299–1312.
- [39] J. Kawahara, A. BenTaieb, G. Hamarneh, Deep features to classify skin lesions, in: 13th International Symposium on Biomedical Imaging, IEEE,
665 2016, pp. 1397–1400.
- [40] Z. Yu, X. Jiang, T. Wang, B. Lei, Aggregating Deep Convolutional Features for Melanoma Recognition in Dermoscopy Images, in: International Workshop on Machine Learning in Medical Imaging, Springer, 2017, pp. 238–246.
- 670 [41] N. Codella, Q.-B. Nguyen, S. Pankanti, D. Gutman, B. Helba, A. Halpern, J. R. Smith, Deep learning ensembles for melanoma recognition in dermoscopy images, arXiv preprint arXiv:1610.04662.
- [42] T. DeVries, D. Ramachandram, Skin Lesion Classification Using Deep Multi-scale Convolutional Neural Networks, arXiv preprint
675 arXiv:1703.01402.
- [43] B. Harangi, Skin lesion detection based on an ensemble of deep convolutional neural networks, arXiv preprint arXiv:1705.03360 (2015) 1–4.

- [44] X. Yang, Z. Zeng, S. Y. Yeo, C. Tan, H. L. Tey, Y. Su, A Novel Multi-task Deep Learning Model for Skin Lesion Segmentation and Classification, arXiv preprint arXiv:1703.01025.
- [45] C. Nader, B. Nader, Experiments using deep learning for dermoscopy image analysis, *Pattern Recognition Letters* 0 (2018) 1–9.
- [46] K. Matsunaga, A. Hamada, A. Minagawa, H. Koga, Image Classification of Melanoma, Nevus and Seborrheic Keratosis by Deep Neural Network Ensemble, arXiv preprint arXiv:1703.03108.
- [47] I. G. Díaz, Incorporating the Knowledge of Dermatologists to Convolutional Neural Networks for the Diagnosis of Skin Lesions, arXiv preprint arXiv:1703.01976.
- [48] A. Menegola, J. Tavares, M. Fornaciali, L. T. Li, S. Avila, E. Valle, RECOD Titans at ISIC Challenge 2017, arXiv preprint arXiv:1703.04819.
- [49] K. M. Li, E. C. Li, Skin lesion analysis towards melanoma detection via end-to-end deep learning of convolutional neural networks, arXiv preprint arXiv:1807.08332.
- [50] Y. Li, L. Shen, Skin lesion analysis towards melanoma detection using deep learning network, *Sensors* 18 (2) (2018) 556.
- [51] P. Mirunalini, A. Chandrabose, V. Gokul, S. M. Jaisakthi, Deep Learning for Skin Lesion Classification, arXiv preprint arXiv:1703.04364.
- [52] D. Gutman, N. C. F. Codella, E. Celebi, B. Helba, M. Marchetti, N. Mishra, A. Halpern, Skin Lesion Analysis toward Melanoma Detection: A Challenge at the International Symposium on Biomedical Imaging (ISBI) 2016, hosted by the International Skin Imaging Collaboration (ISIC), arXiv preprint arXiv:1605.01397.
- [53] N. C. F. Codella, D. Gutman, M. E. Celebi, B. Helba, M. A. Marchetti, S. W. Dusza, A. Kalloo, K. Liopyris, N. Mishra, H. Kittler, Skin Lesion

- 705 Analysis Toward Melanoma Detection: A Challenge at the 2017 International Symposium on Biomedical Imaging (ISBI), Hosted by the International Skin Imaging Collaboration (ISIC), arXiv preprint arXiv:1710.05006.
- [54] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, ImageNet Large Scale Visual Recognition Challenge, International Journal of Computer Vision 115 (3) (2015) 710 211–252.
- [55] C. M. Bishop, Pattern Recognition and Machine Learning, springer, 2006.
- [56] K. P. Murphy, Machine Learning: A Probabilistic Perspective, MIT press, 2012.
- 715 [57] T. Tieleman, G. Hinton, Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude, COURSERA: Neural networks for machine learning 4 (2) (2012) 26–31.
- [58] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980.
- 720 [59] J. Platt, Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods, Advances in Large Margin Classifiers 10 (3) (1999) 61–74.
- [60] Y. WJ., Index for rating diagnostic tests, CANCER 3 (1) (1950) 32–35.
- [61] G. Hinton, Visualizing Data using t-SNE, Journal of Machine Learning 725 Research 9 (Nov) (2008) 2579–2605.
- [62] L. Van Der Maaten, Barnes-hut-sne, arXiv preprint arXiv:1301.3342.
- [63] J. Long, E. Shelhamer, T. Darrell, Fully Convolutional Networks for Semantic Segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3431–3440.

- 730 [64] C. Szegedy, S. Ioffe, V. Vanhoucke, A. A. Alemi, Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning, in: Association for the Advancement of Artificial Intelligence, 2017, pp. 4278–4284.
- [65] A. Vedaldi, K. Lenc, MatConvNet: Convolutional Neural Networks for MATLAB, in: Proceedings of the 23rd ACM International Conference on
735 Multimedia, ACM, 2015, pp. 689–692.
- [66] A. Mahbod, I. Ellinger, R. Ecker, Ö. Smedby, C. Wang, Breast cancer histological image classification using fine-tuned deep network fusion, in: International Conference Image Analysis and Recognition, Springer, 2018, pp. 754–762.
- 740 [67] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going Deeper with Convolutions, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2015, pp. 1–9.
- [68] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely con-
745 nected convolutional networks., in: CVPR, Vol. 1, 2017, p. 3.
- [69] R. B. Oliveira, A. S. Pereira, J. M. R. S. Tavares, Skin lesion computational diagnosis of dermoscopic images: Ensemble models based on input feature manipulation, *Computer Methods and Programs in Biomedicine* 149 (2017) 43–53.