

# Soft-Attention Improves Skin Cancer Classification Performance

Soumya Kanti Datta   Mohammad Abuzar Shaikh   Sargur N Srihari   Mingchen Gao  
State University of New York, Buffalo  
Buffalo, New York, USA

{soumyyak, mshaikh2, srihari, mgao8}@buffalo.edu

## Abstract

In clinical applications, neural networks must focus on and highlight the most important parts of an input image. Soft-Attention mechanism enables a neural network to achieve this goal. This paper investigates the effectiveness of Soft-Attention in deep neural architectures. The central aim of Soft-Attention is to boost the value of important features and suppress the noise-inducing features. We compare the performance of VGG, ResNet, Inception ResNet v2 and DenseNet architectures with and without the Soft-Attention mechanism, while classifying skin lesions. The original network when coupled with Soft-Attention outperforms the baseline[16] by 4.7% while achieving a precision of 93.7% on HAM10000 dataset [25]. Additionally, Soft-Attention coupling improves the sensitivity score by 3.8% compared to baseline[31] and achieves 91.6% on ISIC-2017 dataset [2].

## 1. Introduction

Skin cancer is the most common cancer and one of the leading causes of death worldwide. Every day, more than 9500 people<sup>2</sup> in the United States are diagnosed with skin cancer, with 3.6 million people<sup>3</sup> diagnosed with basal cell skin cancer each year. Early diagnosis of the illness has a significant effect on the patients' survival rates. As a result, detecting and classifying skin cancer is important.

It is difficult to distinguish between malignant and benign skin diseases because they look so similar. Although a dermatologist's visual examination is the first step in detecting and diagnosing a suspicious skin lesion, it is usually followed by dermoscopy imaging for further analysis [32]. Dermoscopy images provide a high-resolution magni-

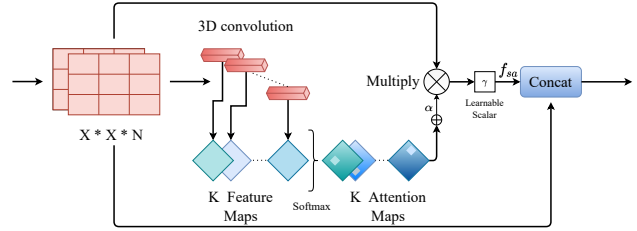


Figure 1. Soft Attention unit

fied image of the infected skin region, but they are not without their drawbacks. Due to the image size being large, it becomes difficult for the feature extractors to extract out the relevant features for classification. Various methods such as Segmentation and detection, Transfer learning, General Adversarial networks, etc. have been used to detect and classify skin cancer. Despite significant progress, skin cancer classification is still a difficult task. This is due to the lack of annotated data and low inter-class variation. Furthermore, the task is complicated by contrast variations, color, shape, and size of the skin lesion, as well as the presence of various artifacts such as hair and veins. Inspired by the work done in [18], this paper studies the effect of soft attention mechanism in deep neural networks. Deep learning architectures identify the image class by learning the salient features and nonlinear interactions. The soft-attention mechanism improves performance by focusing primarily on relevant areas of the input. Moreover, the soft-attention mechanism makes the image classification process transparent to medical personnel, as it maps the parts of the input that the network uses to classify the image, thereby, increasing trust in the classification model.

## 2. Related Work

Following Krichevsky[12], large-scale image classification tasks using deep convolutional neural networks have become common. As reported in the paper[3], the task of skin cancer classification using images has improved rapidly since the implementation of Deep Neural Networks. To make progress, we suggest that soft attention be used to

identify fine-grained variability in the visual features of skin lesions.

Existing art in the field of skin cancer classification used streamlined pipelines based upon current Computer Vision. [4]. Masood et al. in their paper.[13] proposed a general framework from the viewpoint of computer vision, where the methods such as calibration, preprocessing, segmentation, balancing of classes and cross validation are used for automated melanoma screening. In 2018, Valle et al.[26] investigated ten different methodologies to evaluate deep learning models for skin lesion classification. Data augmentation, model architecture, image resolution, input normalization, train dataset, use of segmentation, test data augmentation, additional use of support vector machines, and use of transfer learning are among the ten methodologies they evaluated. They stated that data augmentation had the greatest impact on model efficiency. The same observation is confirmed by Perez’s 2018 paper ”Data Augmentation for Skin Lesion Analysis”[15].

Nonetheless, the problems of low inter-class variance and class imbalance in skin lesion image datasets remain, seriously limiting the capabilities of deep learning models[30]. To fix the lack of annotated data, Zunair et al.[32] proposed the use of adversarial training and Bissoto et al.[1] proposed the use of Generative Adversarial Networks to produce realistic synthetic skin lesion photos.

### 3. Experiment Settings And Method

In this paper, five deep neural networks which are ResNet34, ResNet50 [6], Inception ResNet v2[22], DenseNet201[8] and VGG16 [20], are implemented with soft attention mechanism, to classify skin cancer images. ResNet34, ResNet50[6], Inception ResNet v2, DenseNet201[8] and VGG16[20] are all state of the art feature extractors which are trained on ImageNet dataset. The main components and architecture of the proposed approach is described below:

#### 3.1. Dataset

The experiment is performed on two datasets separately. The two datasets are as follows: HAM10000 dataset [25] and ISIC 2017 dataset.

The HAM10000 dataset [25] consists of 10015 dermatoscopic images of a size of  $450 \times 600$ . It consists of 7 diagnostic categories as follows: Melanoma(MEL), Melanocytic Nevi(NV), Basal Cell Carcinoma(BCC), Actinic Keratosis, and Intra-Epithelial Carcinoma(AKIEC), Benign Keratosis(BKL), Dermatofibroma(DF), Vascular lesions(VASC). All the images are resized to  $299 \times 299$  for Inception ResNet v2[22] architecture and  $224 \times 224$  for the other architectures.

The ISIC 2017 dataset consists 2600 images of size  $767 \times 1022$ . In the training dataset there are 2000 images of

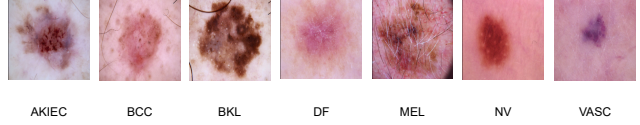


Figure 2. Example of Skin lesions in HAM10000 dataset [25]

3 categories as follows: benign nevi, seborrheic keratosis, and melanoma. The test dataset consist of 600 images. In this experiment we are training our model to classify only benign nevi and seborrheic keratosis. All the images resized to  $224 \times 224$ .

The data in both datasets is then cleaned to remove class imbalances. This is done by the process of over-sampling and under-sampling of data so that there are equal number of images per class. The images are then normalized by dividing each pixel with 255 to keep the pixel values in the range 0 to 1.

#### 3.2. Soft Attention

When it comes to skin lesion images, only a small percentage of pixels are relevant as the rest of the image is filled with various irrelevant artifacts such as veins and hair. So, to focus more on these relevant features of the image, soft attention is implemented. Inspired by the work proposed by Xu et al [28], for image caption generation and the work done by Shaikh et al [18], where they used attention mechanism on images for handwriting verification, in this paper, soft attention is used to classify skin cancer.

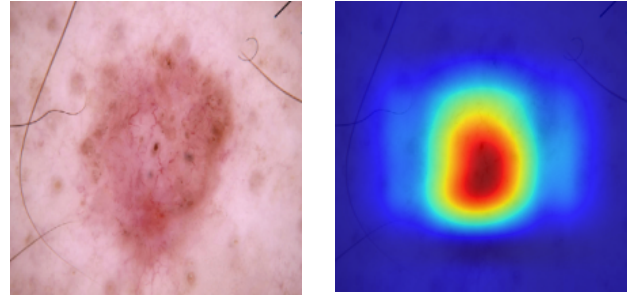


Figure 3. Images with Soft Attention

In Figure [3], we can see that areas with higher attention are red in color . This is because soft attention discredits irrelevant areas of the image by multiplying the corresponding feature maps with low weights. Thus the low attention areas have weights closer to 0. With more focused information, the model performs better.

In the soft attention module as discussed in paper [18] and [23], the feature tensor (t) which flows down the deep neural network is used as input.

$$f_{sa} = \gamma t \left( \left( \sum_{k=1}^K softmax(W_k * t) \right) \right) \quad (1)$$

This feature tensor  $t \in \mathbb{R}^{h \times w \times d}$  is input to a 3D convolution layer[24] with weights  $W_k \in \mathbb{R}^{h \times w \times d \times K}$ , where  $K$  is the number of 3D weights. The output of this convolution is normalized using softmax function to generate  $K = 16$  attention maps. As shown in Figure 1, these attention maps are aggregated to produce a unified attention map that acts as a weighting function  $\alpha$ . This  $\alpha$  is then multiplied with  $t$  to attentively scale the salient feature values, which is further scaled by  $\gamma$  a learnable scalar. Finally, the attentively scaled features ( $f_{sa}$ ) are concatenated with the original feature  $t$  in form of a residual branch. During training we initialize  $\gamma$  from 0.01 so that the network can slowly learn to regulate the amount of attention required by the network.

### 3.3. Model Setup

In this section, the detailed architecture of all the models is discussed. For all experiments, to train the networks, Adam optimizer[11] of 0.01 learning rate and 0.1 epsilon is used. A batch normalization[10] layer is added after each layer in all the networks to introduce some regularization. For the HAM10000 dataset [25], since there are 7 classes of skin cancer, an output layer with 7 hidden units is implemented, followed by a softmax activation unit. All the experiments were executed on the Keras framework.

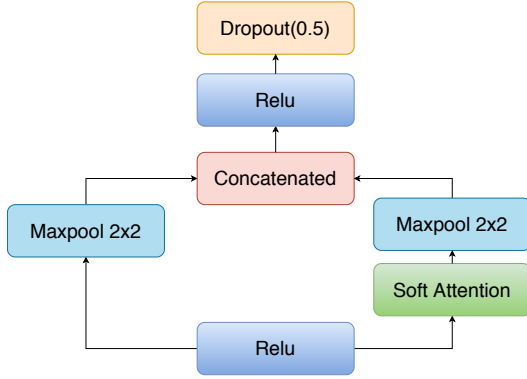


Figure 4. The schema for Soft Attention Block

#### 3.3.1 Inception ResNet v2

In Inception ResNet v2[22], the soft attention layer is added to the Inception Resnet C block of the model where the feature size of the image is  $8 \times 8$  as shown in Figure [5a]. In this case, the soft attention layer is followed by a maxpool layer with a pool size of  $2 \times 2$ , which is then concatenated with the filter concatenate layer of the inception block. The concatenate layer is then followed by a relu activation unit. To regularize the output of the attention layer, the activation

unit is followed by a 0.5 dropout layer[21] as in Figure [4]. The network is trained for 150 epochs with early stopping patience of 30. The overall network is shown in Figure [5a].

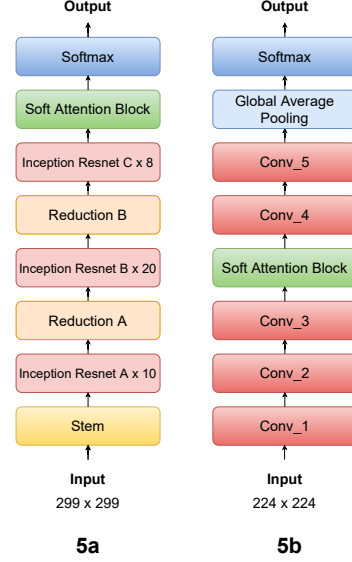


Figure 5. **5a.** End to end architecture Of Inception ResNet v2[22] with Soft Attention Block. **5b.** End to end architecture Of ResNet34[6] with Soft Attention Block . conv\_x indicates convolution blocks, where x is the block number.

#### 3.3.2 DenseNet201

In DenseNet201[8], the soft attention layer is added to the 4th dense block where the size of feature map of the image is  $7 \times 7$  as shown in Figure[6]. Like in the previous model, the soft attention layer is integrated with the same procedure as it was integrated with the Inception ResNet V2[22] architecture.[4]. The network is trained for 150 epochs with early stopping patience of 35.

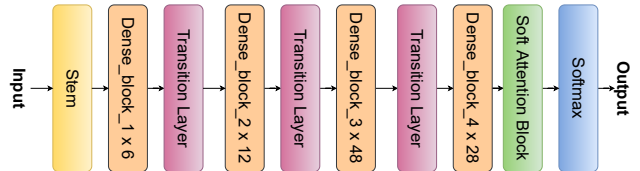


Figure 6. End to end schema of DenseNet201[8] with Soft Attention Block.

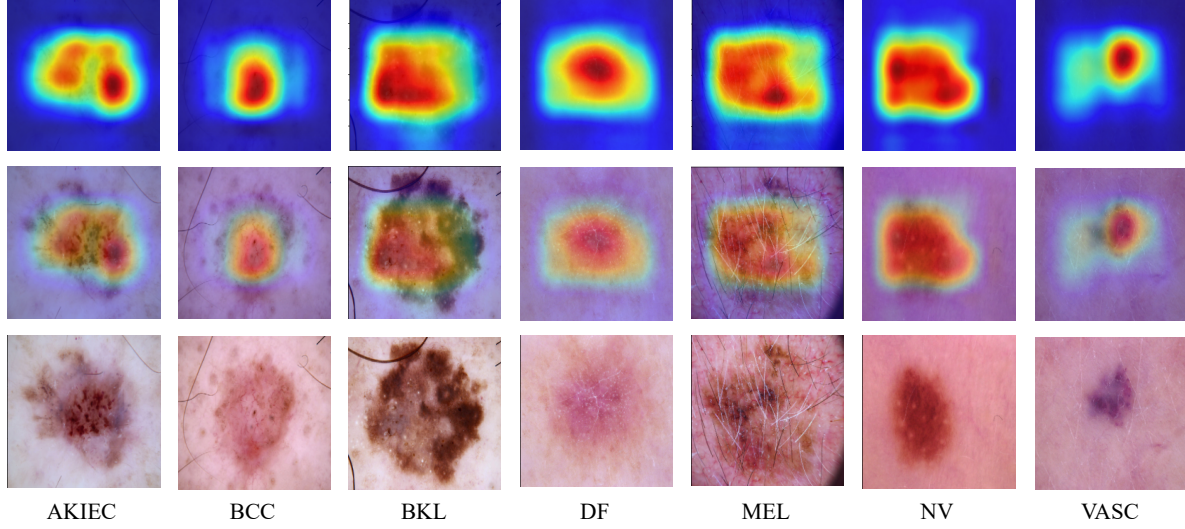


Figure 7. Soft Attention maps of Skin lesion in Inception ResNet V2 on HAM10000 dataset [25]

### 3.3.3 ResNet34 and ResNet50

In ResNet34[6], a soft attention layer is added after the 3rd convolution block where the size of feature map is 28 x 28 as shown in [5b] whereas, in the ResNet50[6], the soft attention layer is added after the 5th convolution block where the size of feature map is 7 x 7. In both cases, the soft attention layer is followed by a maxpool layer with a pool size of 2x2, which is then concatenated with the standard maxpool layer of the architecture, as shown in Figure [4]. The concatenate layer is then followed by a relu activation unit. To regularize the output of the attention layer, the activation unit is followed by a 0.5 dropout[21] layer. This is the same approach as to how the soft attention module was integrated with the Inception ResNet V2[22] architecture. The overall architecture for ResNet 34 model is shown in Figure [5b].

### 3.3.4 VGG16

In VGG16[20], the soft attention layer is added after the conv\_layer\_4 of the VGG16 architecture where the size of feature map is 28 x 28. Like in the previous model, the soft attention layer is integrated with the same procedure as it was integrated with the ResNet[6] and Inception ResNet V2[22] architecture.[4]. The network is trained for 300 epochs with early stopping patience of 65. The overall architecture for the model is shown in Figure [8].

In Figure [8], a Conv\_layer block consists of two to three convolution layers with filters of sizes ranging from 64 to 512, followed by a maxpool layer. Conv\_layer\_1, and Conv\_layer\_2 consists of two convolution layers each with 64, and 128 filters respectively, and Conv\_layer\_3, Conv\_layer\_4 and Conv\_layer\_5 consists of three convolution layers each with 256, 512 and 512 filters respectively.

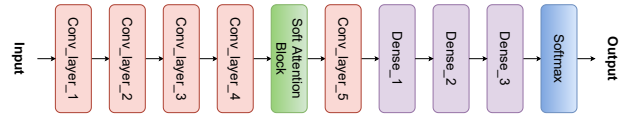


Figure 8. End to end schema of VGG16[20] with Soft Attention Block. conv x indicates convolution layer with x filters.

### 3.4. Loss Function

In this experiment, there are seven different classes of skin cancer. Hence , categorical cross entropy loss ( $L_{CCE}$ ) is used to optimize the neural network.

$$L_{CCE} = - \sum_{i=1}^C t_i \log(f(s)_i) \quad (2)$$

where

$$f(s)_i = \frac{e^{s_i}}{\sum_{j=1}^C e^{s_j}} \quad (3)$$

Here, as there are seven classes,  $C \in [0..6]$ , where  $t_i$  is the ground truth and  $s_i$  is the CNN score for each class  $i$  in  $C$ .  $f(s)_i$  is the softmax activation function applied to the scores.

### 3.5. Evaluation Metrics

In this paper, the model is evaluated using  $Precision = \frac{TP}{TP+FP}$ ,  $Accuracy = \frac{TP+TN}{T}$ ,  $Sensitivity = \frac{TP}{TP+FN}$ ,  $Specificity = \frac{TN}{TN+FP}$  and AUC scores[9]. Here TN, TP, FP, FN, T mean, True Negatives, True Positives, False Positives, False Negatives, Total Number respectively.



Dis.	Precision										AUC										#
	[22]	[22] + SA	[8]	[8] + SA	[20]	[20] + SA	[6]50	[6]50 + SA	[6]34	[6]34 + SA	[22]	[22] + SA	[8]	[8] + SA	[20]	[20] + SA	[6]50	[6]50 + SA	[6]34	[6]34 + SA	
AKIEC	0.830	<b>1.000</b>	<b>1.000</b>	0.920	0.620	0.700	0.740	0.670	0.670	0.500	<b>0.993</b>	0.981	0.975	0.967	0.949	0.964	0.980	0.981	0.969	0.970	23
BCC	0.850	0.880	0.830	0.800	0.540	0.620	<b>0.910</b>	0.880	0.660	0.880	0.997	<b>0.998</b>	0.993	0.994	0.977	0.984	0.997	0.996	0.991	0.993	26
BKL	<b>0.850</b>	0.720	0.690	0.730	0.570	0.630	0.670	0.670	0.510	0.520	0.970	<b>0.982</b>	0.960	0.964	0.930	0.900	0.948	0.964	0.904	0.916	66
DF	0.670	<b>1.000</b>	0.500	<b>1.000</b>	0.250	0.500	0.800	<b>1.000</b>	0.400	0.330	0.973	<b>0.982</b>	0.851	0.921	0.847	0.809	0.973	0.971	0.925	0.949	6
MEL	0.700	0.670	0.540	0.530	0.500	0.430	0.520	<b>0.730</b>	0.420	0.540	0.965	0.974	0.963	<b>0.976</b>	0.925	0.956	0.961	0.973	0.910	0.953	34
NV	0.930	<b>0.970</b>	0.950	0.950	0.930	0.950	0.950	0.950	0.930	0.930	<b>0.984</b>	<b>0.984</b>	0.975	0.976	0.954	0.951	0.974	0.979	0.944	0.958	663
VASC	<b>1.000</b>	<b>1.000</b>	0.900	0.830	<b>1.000</b>	<b>1.000</b>	0.900	<b>1.000</b>	0.910	0.820	<b>1.000</b>	<b>1.000</b>	0.993	0.999	0.972	0.999	0.995	0.999	0.999	0.996	10
Avg	0.832	<b>0.892</b>	0.771	0.824	0.631	0.690	0.783	0.841	0.642	0.646	0.983	<b>0.984</b>	0.959	0.971	0.936	0.937	0.975	0.980	0.949	0.962	828
W. Avg	0.905	<b>0.937</b>	0.904	0.909	0.862	0.882	0.898	0.910	0.857	0.865	0.982	<b>0.984</b>	0.974	0.975	0.951	0.948	0.972	0.978	0.942	0.957	828

Table 1. Ablation results for choosing the best model on HAM10000 dataset [25]. [22] refers to IRv2 architecture, [8] refers to DenseNet 201 architecture, [20] refers to VGG 16 architecture, and [6] refers to ResNet architecture.

## 4. Discussion

### 4.1. Ablation Analysis

Table 1 lists, the performance of all the models in terms of precision, and AUC score on HAM10000 dataset [25]. In this table (+SA) stands for models with soft attention. IRv2 stands for Inception ResNet v2[22], [6]34 stands for ResNet34[6] and [6]50 stands for ResNet50[6]. From the table, it can be observed that IRv2 when coupled with SA (IRv2+SA) shows significant improvements in results, with a precision and AUC score of 93.7% and 98.4% respectively, which are also the highest scores amongst all models. Furthermore, we can see that Soft Attention (SA) boosts the performance of IRv2 by 3.2% in terms of precision as compared to the original IRv2 model. This phenomenon is true for VGG16, ResNet34, ResNet50 and DenseNet201 as well. For instance, Soft Attention (SA) boosts the precision of DenseNet201[8], ResNet34[6], ResNet50[6], and VGG16[20] by 0.5%, 0.8%, 1.2% and 2% respectively. We see a similar behaviour for the AUC scores when SA block is integrated in to the networks, such as, the performance of ResNet50[6], and ResNet34[6] has grown by 0.6% and 1.5% respectively and the performance of DenseNet201[8], and VGG16[20] is on par with the original models.

Although IRv2+SA performs the best in terms of weighted average(W.Avg), when we look at it's class wise performance, we can see that Soft Attention enhances the efficiency of the original IRv2 while categorizing AKIEC, BCC, DF and NV by 17%, 3%, 33% and 4% respectively in terms of precision. Moreover, when comparing AUC scores, the IRv2+SA performs better for BKL and MEL by 1.2% and 0.9% respectively, while, for BCC, NV and VASC, IRv2+SA performs as good as original model.

We thus select IRv2 coupled with SA (IRv2+SA) for our experiments, also the SA block consistently boosts the performance of it's original counterpart, hence, we can justify the integration of Soft Attention to the networks.

### 4.2. Quantitative Analysis

When we tested the model with different train-test splits on the HAM10000 dataset [25], we discovered that the model with 85 % training data outperforms the model with 80 % and 70 % training data by 2.2 % and 2.6 % respectively, as shown in Table 2. Hence we select 85/15% training/testing split for performing our experiments.

Disease	split = 15			split = 20			split = 30		
	Support	Precision	AUC	Support	Precision	Auc	Support	Precision	Auc
AKIEC	23	1.000	0.981	30	0.750	0.958	45	0.690	0.972
BCC	26	0.880	0.998	35	0.880	0.992	53	0.830	0.995
BKL	66	0.720	0.982	88	0.790	0.972	132	0.720	0.960
DF	6	1.000	0.973	8	1.000	0.998	12	0.710	0.989
MEL	34	0.670	0.974	46	0.490	0.953	69	0.600	0.946
NV	663	0.970	0.984	883	0.960	0.981	1325	0.960	0.978
VASC	10	1.000	1.000	13	0.920	0.999	19	0.860	0.983
Avg	828	0.892	0.984	1103	0.827	0.9793	1655	0.766	0.975
W. Avg	828	<b>0.937</b>	0.984	1103	0.915	0.9797	1655	0.911	0.976

Table 2. Comparison with Models with different train-test split on HAM10000 dataset [25]

Furthermore, the proposed approach is compared with state-of-the-art models for skin cancer classification on the HAM10000 dataset [25] in Table 3. Our Soft Attention-based approach outperforms the baseline[16] by 4.7% in terms of precision. In terms of AUC scores, our Soft Attention-based approach clearly outperforms them all by 0.5% to 4.3%.

Model	Avg AUC	Precision	Accuracy
Loss balancing and ensemble[5]	0.941	-	0.926
Single Model Deep Learning[29]	0.974	-	0.864
Data classification augmentation[19]	0.975	-	0.853
Two path CNN model[14]	-	-	0.886
Various Deep CNN (Baseline) [16]	0.979	0.890	-
<b>IRv2+SA(Proposed Approach)</b>	<b>0.984</b>	<b>0.937</b>	<b>0.934</b>

Table 3. Comparison with state-of-the-art-Model in terms of Average AUC score on HAM10000 dataset [25]

In Table 4, the performance of the proposed approach Inception Resnet V2[22] (IRv2<sub>5x5</sub>+SA and IRv2<sub>12x12</sub>+SA) with soft attention is measured on ISIC-2017 dataset [2] on basis of AUC scores, Accuracy , Sensitivity and Specificity

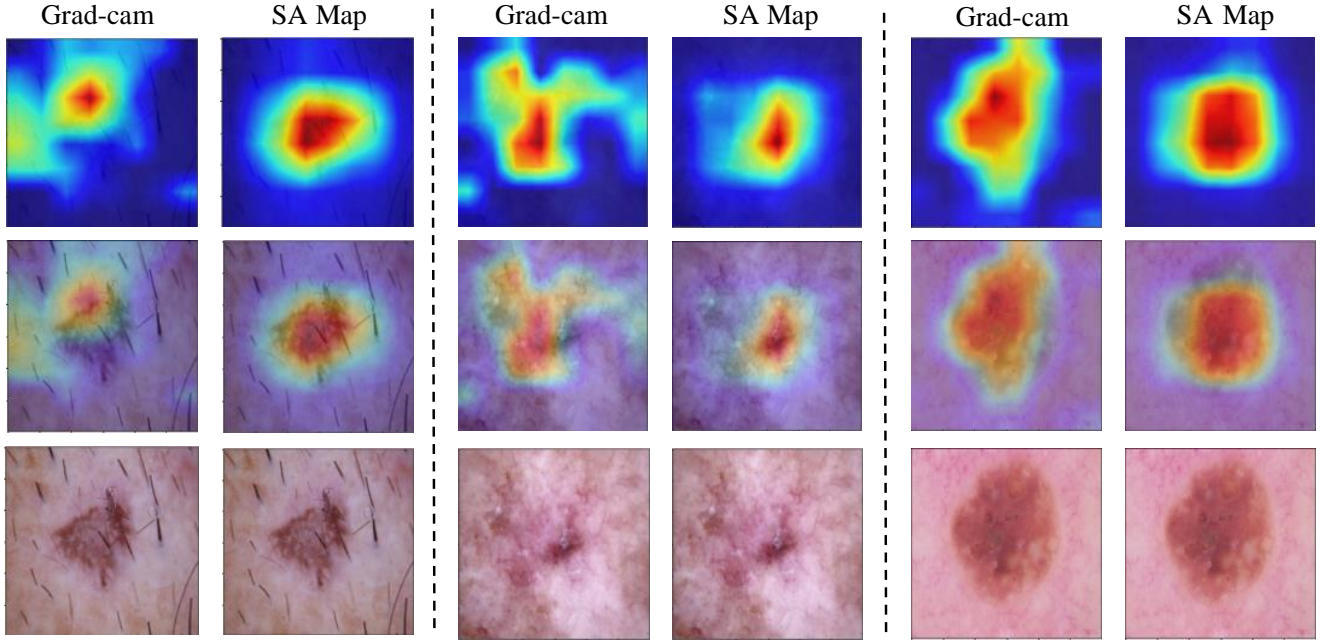


Figure 9. Comparison of GradCAM [17] heatmaps with our Soft Attention (SA) maps on HAM10000 dataset [25]

with the state-of-the-art models.

Networks	AUC	Accuracy	Sensitivity	Specificity
ResNet50 [6]	0.948	0.842	0.867	0.837
RAN50 [27]	0.942	0.862	0.878	0.859
SEnet50 [7]	0.952	0.863	0.856	0.865
ARL-CNN50[31]	0.958	0.868	0.878	<b>0.867</b>
IRv2 <sub>12x12</sub> +SA	0.935	0.898	<b>0.945</b>	0.711
<b>IRv2<sub>5x5</sub>+SA</b>	<b>0.959</b>	<b>0.904</b>	0.916	0.833

Table 4. Comparison with state-of-the-art-Model in terms of AUC, Accuracy, sensitivity and specificity score on ISIC-2017 dataset [2]

From Table 4, it can be observed that in IRv2<sub>5x5</sub>+SA, and in IRv2<sub>12x12</sub>+SA, the attention layer was added when the feature map size is 5x5 and 12x12 respectively. Out of the two models with soft attention, the model IRv2<sub>5x5</sub>+SA outperforms IRv2<sub>12x12</sub>+SA in terms of AUC scores, Accuracy, and Specificity by a percentage of 2.4%, 0.6%, and 12.2% respectively whereas IRv2<sub>12x12</sub>+SA outperforms IRv2<sub>5x5</sub>+SA in terms of Sensitivity by 2.9%. In this case, the attention layer was added when the feature size is 5x5. When IRv2<sub>5x5</sub>+SA is compared with the ARL-CNN50[31] (baseline model), it performs on par with it in terms of AUC score but our model outperforms it when it comes to accuracy and Sensitivity by 3.6% and 3.8% respectively. But ARL-CNN50[31] takes the upper hand when it comes to Specificity by 3.4%. Since sensitivity measures the proportion of correctly identified positives and specificity measures the proportion of correctly identified negatives, we

are prioritizing Sensitivity because classifying a person with cancer as not having cancer is riskier than vice versa.

### 4.3. Qualitative Analysis

Fig.7 displays the Soft Attention heat maps from the IRv2+SA model. In the Fig.7, the images on the bottom row are the input images of the seven skin cancer categories. The images in the middle row show the Soft Attention maps superimposed on input images to show where the model is focusing and the images of the top row are attention maps themselves.

In Fig.9, we show pairs of comparison between the Soft Attention maps with Grad-CAM [17] heatmaps. In the first pair, the SA map focuses on the main part of the lesion area whereas the Grad-cam heatmap is slightly shifted towards top left and is also spread out on the uninfected area of skin. We have similar observations for the second and third pairs as well. From this observation it is evident that the Soft Attention maps are focused more on the relevant locations of the image compared to Grad-CAM[17] heatmaps.

## 5. Conclusion

In this paper, we present the implementation and utility of Soft Attention mechanism being applied while image encoding to tackle the problem of high-resolution skin cancer image classification. The model outperformed the current state-of-the-art approaches on the HAM10000 dataset [25] and the ISIC-2017 dataset [2]. This demonstrates the Soft Attention based deep learning architecture’s potential and

effectiveness in image analysis. The Soft Attention mechanism also eliminates the need of using external mechanisms like GradCAM [17], and internally provides the location of where the model focuses while categorizing a disease, while also boosting the performance of the main network. Soft Attention has the added advantage of naturally dealing with image noise internally. In the future, this model can be implemented in dermoscopy systems to assist dermatologists. This mechanism can be easily implemented to classify data from other medical databases as well.

## References

- [1] Alceu Bissoto, Fábio Perez, Eduardo Valle, and Sandra Avila. Skin lesion synthesis with generative adversarial networks. In *OR 2.0 Context-Aware Operating Theaters, Computer Assisted Robotic Endoscopy, Clinical Image-Based Procedures, and Skin Image Analysis*, pages 294–302. Springer, 2018. 2
- [2] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin K. Mishra, Harald Kittler, and Allan Halpern. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (ISIC). *CoRR*, abs/1710.05006, 2017. 1, 5, 6
- [3] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017. 1
- [4] Michel Fornaciali, Micael Carvalho, Flávia Vasques Bittencourt, Sandra Avila, and Eduardo Valle. Towards automated melanoma screening: Proper computer vision & reliable results. *arXiv preprint arXiv:1604.04024*, 2016. 2
- [5] Nils Gessert, Maximilian Nielsen, Mohsin Shaikh, René Werner, and Alexander Schlaef. Skin lesion classification using ensembles of multi-resolution efficientnets with meta data. *MethodsX*, page 100864, 2020. 5
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 3, 4, 5, 6
- [7] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 6
- [8] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 2, 3, 5
- [9] Jin Huang and Charles X Ling. Using auc and accuracy in evaluating learning algorithms. *IEEE Transactions on knowledge and Data Engineering*, 17(3):299–310, 2005. 4
- [10] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 3
- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 3
- [12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 1
- [13] Ammara Masood and Adel Ali Al-Jumaily. Computer aided diagnostic support system for skin cancer: a review of techniques and algorithms. *International journal of biomedical imaging*, 2013, 2013. 2
- [14] Hemanth Nadipineni. Method to classify skin lesions using dermoscopic images. *arXiv preprint arXiv:2008.09418*, 2020. 5
- [15] Fábio Perez, Cristina Vasconcelos, Sandra Avila, and Eduardo Valle. Data augmentation for skin lesion analysis. In *OR 2.0 Context-Aware Operating Theaters, Computer Assisted Robotic Endoscopy, Clinical Image-Based Procedures, and Skin Image Analysis*, pages 303–311. Springer, 2018. 2
- [16] Amirreza Rezvantalab, Habib Safigholi, and Somayeh Karimijeshni. Dermatologist level dermoscopy skin cancer classification using different deep learning convolutional neural networks algorithms. *arXiv preprint arXiv:1810.10348*, 2018. 1, 5
- [17] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 6, 7
- [18] Mohammad Abuzar Shaikh, Tiehang Duan, Mihir Chauhan, and Sargur N. Srihari. Attention based writer independent verification. *2020 17th International Conference on Frontiers in Handwriting Recognition (ICFHR)*, Sep 2020. 1, 2
- [19] Shuwei Shen, Mengjuan Xu, Fan Zhang, Pengfei Shao, Honghong Liu, Liang Xu, Chi Zhang, Peng Liu, Zhihong Zhang, Peng Yao, et al. Low-cost and high-performance data augmentation for deep-learning-based skin lesion classification. *arXiv preprint arXiv:2101.02353*, 2021. 5
- [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2, 4, 5
- [21] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014. 3, 4
- [22] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261*, 2016. 2, 3, 4, 5
- [23] Naofumi Tomita, Behnaz Abdollahi, Jason Wei, Bing Ren, Arief Suriawinata, and Saeed Hassanpour. Attention-based deep neural networks for detection of cancerous and precancerous esophagus tissue on histopathological slides. *JAMA network open*, 2(11):e1914645–e1914645, 2019. 2
- [24] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with

- 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015. 3
- [25] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018. 1, 2, 3, 4, 5, 6
- [26] Eduardo Valle, Michel Fornaciali, Afonso Menegola, Julia Tavares, Flávia Vasques Bittencourt, Lin Tzy Li, and Sandra Avila. Data, depth, and design: Learning reliable models for skin lesion analysis. *Neurocomputing*, 383:303–313, 2020. 2
- [27] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification, 2017. 6
- [28] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015. 2
- [29] Peng Yao, Shuwei Shen, Mengjuan Xu, Peng Liu, Fan Zhang, Jinyu Xing, Pengfei Shao, Benjamin Kaffenberger, and Ronald X Xu. Single model deep learning on imbalanced small datasets for skin lesion classification. *arXiv preprint arXiv:2102.01284*, 2021. 5
- [30] Lequan Yu, Hao Chen, Qi Dou, Jing Qin, and Pheng-Ann Heng. Automated melanoma recognition in dermoscopy images via very deep residual networks. *IEEE transactions on medical imaging*, 36(4):994–1004, 2016. 2
- [31] Jianpeng Zhang, Yutong Xie, Yong Xia, and Chunhua Shen. Attention residual learning for skin lesion classification. *IEEE transactions on medical imaging*, 38(9):2092–2103, 2019. 1, 6
- [32] Hasib Zunair and A Ben Hamza. Melanoma detection using adversarial training and deep transfer learning. *Physics in Medicine & Biology*, 2020. 1, 2