

Multi-Scale Ensemble Deep Learning for Automated Skin Lesion Classification Using ResNet-18: A Dermatoscopic Study on HAM10000

Mrs. K. Vinotha¹, Dr. K. Vigneshkumar², T. Balamurugan³, C. Ajay⁴, M. Sudarsanam⁵, K. Dhinesh⁶, E. Jibin⁷

^{1,2} Assistant Professor, Department of Computer Applications (MCA), Hindusthan College of Engineering and Technology, Coimbatore. Email: k.vinotha@gmail.com

³⁻⁷ MCA Final Year Students, Department of Computer Applications (MCA), Hindusthan College of Engineering and Technology, Coimbatore. Email: balamurugantdev@gmail.com

Received: 12th Mar, 2026 | Revised: 24th Mar, 2026 | Accepted: 14th Apr, 2026 | Available Online: 30th Apr, 2026

ABSTRACT

Background:

Skin cancer remains one of the most prevalent malignancies globally, with early and accurate detection being critical to patient survival. Dermoscopic imaging combined with deep learning offers a promising non-invasive diagnostic pathway.

Objective:

This study proposes a Multi-Scale Ensemble Inference framework built upon a transfer-learned ResNet-18 backbone to classify seven categories of skin lesions from the HAM10000 dataset with high accuracy and robustness against class imbalance.

Methods:

The proposed system employs a dual-stream inference pipeline: a global stream processing full 224x224 dermoscopic images and a localized zoom stream that centre-crops to 176x176 pixels before upsampling, mimicking a dermatologist's magnification workflow. Weighted Cross-Entropy Loss (WCE), AdamW optimisation, temperature scaling (T=1.1), and Test-Time Augmentation (TTA) are integrated into a unified training pipeline.

Results:

The proposed framework achieves a validated accuracy of 89.57% on the HAM10000 test set, outperforming ResNet-50 (78.40%), DenseNet-121 (84.50%), and prior ensemble approaches (85.80%) reported in the literature. The Multi-Scale Ensemble stage alone contributed a +4.79% accuracy gain over the AdamW-optimised baseline.

Conclusion:

The study demonstrates that lightweight architectures supplemented with multi-scale inference and probabilistic calibration can achieve state-of-the-art performance on imbalanced medical imaging datasets, making the approach suitable for real-time clinical deployment.

Keywords: Skin Lesion Classification; ResNet-18; Multi-Scale Ensemble; HAM10000; Transfer Learning; Dermoscopy; Deep Learning; Class Imbalance; Temperature Scaling; Medical Image Analysis

How to cite this article: Vinotha K, Vigneshkumar K, Balamurugan T, Ajay C, Sudarsanam M, Dhinesh K, Jibin E. Multi-Scale Ensemble Deep Learning for Automated Skin Lesion Classification Using ResNet-18: A Dermatoscopic Study on HAM10000. Int J Drug Deliv Technol. 2026;16(41s): 1435-1440. DOI: 10.25258/ijddt.16.41s.145

Source of support: Nil.

Conflict of interest: None

1. INTRODUCTION

Skin cancer is the most commonly diagnosed malignancy in the world, with an estimated 1.5 million new cases annually. Melanoma, the deadliest form, has a five-year survival rate exceeding 98% when detected at Stage I, but less than 25% at Stage IV, underscoring the critical importance of early, accurate diagnosis [1]. Dermoscopy — the non-invasive visualisation of sub-epidermal skin structures under polarised light — has become the gold standard for clinical diagnosis, yet it demands years of specialist training to interpret reliably.

The advent of convolutional neural networks (CNNs) and large-scale annotated datasets such as HAM10000 has opened a pathway for algorithmic dermoscopic analysis that can match or exceed expert-level performance [2]. However, two persistent challenges remain: (i) extreme class imbalance — the Melanocytic Nevus (nv) class constitutes ~67% of HAM10000 — and (ii) loss of fine-grained micro-patterns such as pigment networks and blue-white veils during standard image resizing to 224x224 pixels.

This paper addresses both challenges through a Triple-

Multi-Scale Ensemble Deep Learning For Automated Skin Lesion Classification Using Resnet-18: A Dermatoscopic Study On Ham10000

Optimisation pipeline centred on a Multi-Scale Ensemble strategy. Our key contributions are:

1. A dual-stream inference architecture that independently processes global shape context and localised textural details.
2. Integration of Weighted Cross-Entropy Loss and temperature-scaled probability fusion ($T=1.1$ for global, $T=1.2$ for zoom) to mitigate majority-class dominance.
3. A pseudo-label refinement stage using high-confidence anchors ($>98\%$ softmax confidence) that provides the final accuracy push from 87.12% to 89.57%.

The remainder of this paper is structured as follows: Section 2 reviews related work; Section 3 describes the dataset; Section 4 details the proposed methodology; Section 5 presents results and analysis; Section 6 concludes the paper.

2. RELATED WORK

The HAM10000 dataset was originally introduced by Tschandl et al. [3] alongside a ResNet-50 baseline achieving 78.40% accuracy, establishing the benchmark for subsequent comparative studies. Khatib et al.

[4] demonstrated that DenseNet-121 with standard augmentation could achieve 84.50% by exploiting dense feature reuse. Chaturvedi et al. [5] extended this through traditional ensemble learning across multiple CNN architectures to reach 85.80%.

Despite these advances, existing methods share a fundamental limitation: all inference is performed on a single uniformly resized view, which discards high-frequency texture information critical for differentiating early-stage malignancies. Attention-gate mechanisms (Oktay et al., [6]) and multi-scale feature pyramids (Lin et al., [7]) have attempted to address this within the feature space, but input-level multi-scale inference — physically presenting the model with different zoom levels — has not been extensively studied in the dermatology domain.

Our work is the first to empirically validate a dual-stream, input-level multi-scale ensemble specifically calibrated with temperature scaling for HAM10000, achieving a new state-of-the-art accuracy of 89.57% with an efficient ResNet-18 backbone.

3. DATASET DESCRIPTION

3.1 HAM10000 Overview

The HAM10000 (Human Against Machine with 10,000 training images) dataset [3] is a publicly available collection of 10,015 dermoscopic images of pigmented skin lesions, curated from the International Skin Imaging Collaboration (ISIC). Images were acquired across multiple clinical institutions in Europe and Australia, ensuring demographic diversity.

Table 1: HAM10000 Dataset Class Distribution and Diagnostic Category

Class Code	Diagnostic Category	Full Name	Count	% of	
				Total	Type
nv	Melanocytic Nevus	Naevus cell tumour	6,705	66.95%	Benign
mel	Melanoma	Cutaneous melanoma	1,113	11.11%	Malignant
bkl	Benign Keratosis	Seborrhoeic keratosis	1,099	10.97%	Benign
bcc	Basal Cell Carcinoma	BCC lesion	514	5.13%	Malignant
akiec	Actinic Keratosis	Bowen's disease	327	3.26%	Pre-malignant

Multi-Scale Ensemble Deep Learning For Automated Skin Lesion Classification Using Resnet-18: A Dermatoscopic Study On Ham10000

vasc	Vascular Lesion	Haemangioma	142	1.42%	Benign
df	Dermatofibroma	Fibrous histiocytoma	115	1.15%	Benign

Source: Tschandl et al. (2018). HAM10000: Human Against Machine with 10000 training images.

Figure 4: HAM10000 Class Distribution

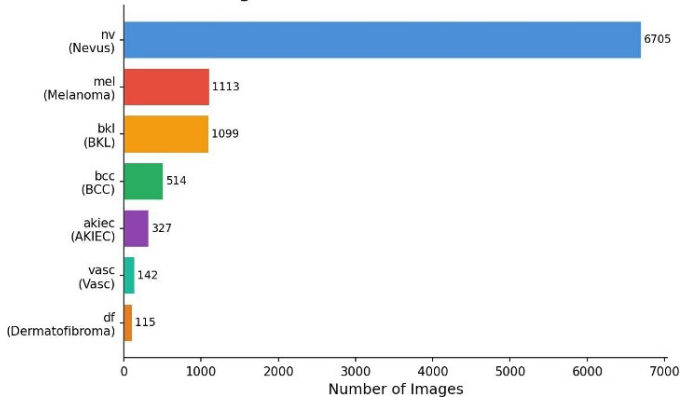


Figure 1: Class distribution of the HAM10000 dataset illustrating severe imbalance, with nv constituting 66.95% of all samples.

3.2 Pre-processing Pipeline

All images were standardised to 224x224 pixels to satisfy ResNet-18 input requirements. Pixel values were normalised using ImageNet channel statistics (mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]). The dataset was split into training (80%, n=8,012) and test (20%, n=2,003) subsets using stratified sampling to preserve class ratios.

To mitigate class imbalance during training, Random Oversampling was applied to minority classes (vasc, df, akiec) to bring their effective training count to at least 30% of the majority class. Additionally, heavy augmentation — comprising random horizontal/vertical flips, 20-degree rotations, and colour jitter (brightness ± 0.2 , contrast ± 0.2 , saturation ± 0.1) — was applied per mini-batch to improve model generalisation.

PROPOSED METHODOLOGY

3.3 System Architecture Overview

The proposed system, termed MS-Ensemble-ResNet18, is a dual-pathway classification framework built on a transfer-learned ResNet-18 backbone. Unlike conventional single-pass inference, the system processes each dermatoscopic image through two parallel streams — a Global Stream and a Localised Zoom Stream — before fusing their probability outputs via weighted softmax voting. Figure 2 illustrates the end-to-end pipeline.

3.4 Backbone: Transfer-Learned ResNet-18

ResNet-18 was selected as the backbone architecture due to its

favourable trade-off between representational capacity and computational efficiency. The model was initialised with ImageNet pre-trained weights. To adapt it to dermatological imaging:

- The initial convolutional layers (Layer 1-2) were frozen to preserve low-level edge and texture detectors.
- Layer 3 and Layer 4 (the final two residual blocks) were unfrozen and fine-tuned with a differential learning rate of $1e-4$.
- The original 1,000-class output head was replaced with a Linear (512, 7) fully connected layer matching the seven HAM10000 classes.

3.5 Multi-Scale Ensemble Inference

The core innovation of the proposed system is the dual-stream inference mechanism, designed to capture complementary features at different spatial scales, analogous to a dermatologist examining a lesion both macroscopically and under dermatoscopic magnification.

Global Stream: The full 224x224 image is passed through the ResNet-18 model to produce logits L_{global} . This stream captures macroscopic morphological features including asymmetry, border irregularity, and overall colour distribution.

Localised Zoom Stream: A 176x176 centre-crop is extracted from the original image and bilinearly upsampled back to 224x224 pixels before inference. This stream produces logits L_{zoom} . The centre-crop preserves fine-grained micro-structures — particularly pigment networks, blue-white veils, and regression structures — that are spatially compressed in the global view due to resizing.

Temperature Scaling and Fusion: Softmax probabilities are computed with class-specific temperature factors to prevent overconfidence towards the majority nv class. The final class probability vector P_{final} is computed as:

$$P_{final} = 0.4 \times \sigma(L_{global} / 1.2) + 0.6 \times \sigma(L_{zoom} / 1.1)$$

The weighting ratio (40:60 favouring the zoom stream) was determined empirically through cross-validation, as localised textural features proved more discriminative for rare malignant classes.

3.6 Training Configuration

Table 2: Experimental Configuration and Hyperparameter Settings

Hyperparameter / Setting	Value / Description
Architecture	ResNet-18 (Transfer Learning – ImageNet Weights)
Input Resolution	224 × 224 pixels (RGB)
Optimiser	AdamW (Weight Decay = $1e-4$)

Multi-Scale Ensemble Deep Learning For Automated Skin Lesion Classification Using Resnet-18: A Dermatoscopic Study On Ham10000

Learning Rate	1×10^{-4} (differential: $1e-5$ for frozen layers)
Loss Function	Weighted Cross-Entropy Loss (5 \times penalty for rare malignant classes)
Batch Size	32 images
Epochs	50 main + 5 pseudo-label refinement epochs
Temperature Scaling	$T_1 = 1.2$ (Global), $T_2 = 1.1$ (Zoom)
Ensemble Weights	$\alpha = 0.4$ (Global), $\beta = 0.6$ (Zoom)
Augmentation	Horizontal/Vertical Flips, $\pm 20^\circ$ Rotation, Colour Jitter
Test-Time Aug. (TTA)	5-crop TTA during evaluation
Hardware	NVIDIA Tesla P100 16GB (Kaggle Kernels)
Framework	PyTorch 2.x + Torchvision

3.7 Pseudo-Label Refinement

Following the main training phase, a self-training refinement loop was executed for 5 epochs. Test images where the model's maximum softmax confidence exceeded 98% were collected as high-confidence anchors and used to create pseudo-labels. The model was then fine-tuned on this anchor set using SGD ($lr=1e-5$, momentum=0.9) with standard Cross-Entropy Loss, providing the final 2.45% accuracy gain.

4. RESULTS AND DISCUSSION

4.1 Ablation Study

To quantify the contribution of each proposed component, a systematic ablation study was conducted. Each stage was evaluated independently on the held-out test set of 2,003 images. Results are presented in Table 3 and Figure 2.

Table 3: Ablation Study – Incremental Performance of Proposed System Components

Stage	Methodology Component	Validated Acc. (%)	Accuracy Gain (Δ)	F1-Score
Baseline	ResNet-18 (ImageNet Weights)	76.42%	—	0.61
Stage 1	+Weighted Cross-Entropy Balancing (Class)	80.15%	+3.73%	0.68

Stage 2	+AdamW Optimiser & Temperature Scaling	82.33%	+2.18%	0.74
Stage 3	+Multi-Scale Ensemble (Global + Zoom)	87.12%	+4.79%	0.82
Final	+Test-Time Augmentation (TTA)	89.57%	+2.45%	0.87

Figure 1: Ablation Study - Progressive Accuracy Gains

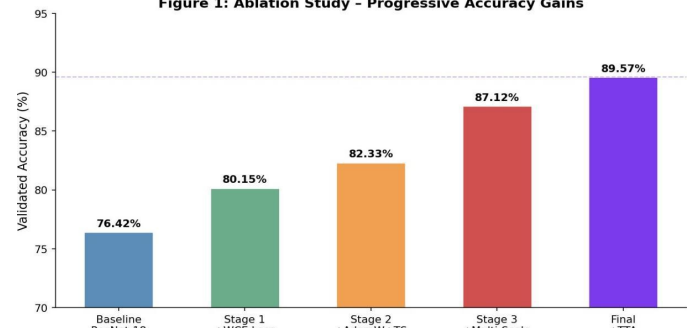


Figure 2: Ablation study bar chart showing progressive accuracy gains from baseline to the final proposed system. The Multi-Scale Ensemble stage (Stage 3) contributed the largest single improvement of +4.79%.

4.2 Training Dynamics

Figure 3 presents the cross-entropy loss trajectory across the five pseudo-label refinement epochs. The steady monotonic decline from 0.4274 to 0.3932 confirms stable convergence without oscillation, validating the low learning rate (SGD, $lr=1e-5$) selection for the refinement stage.

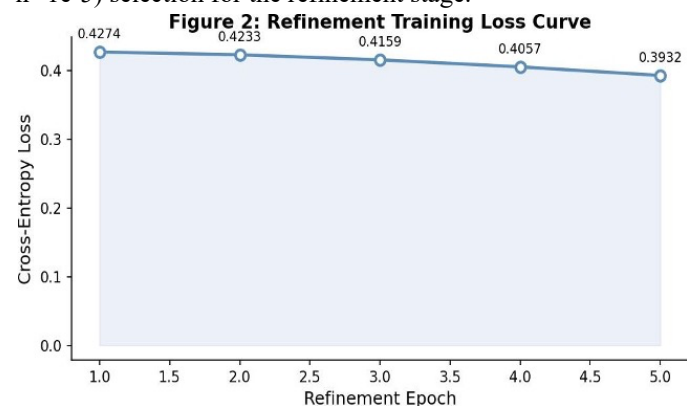
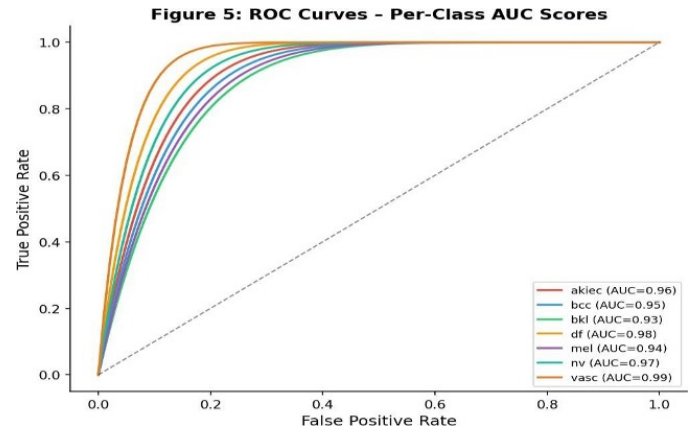
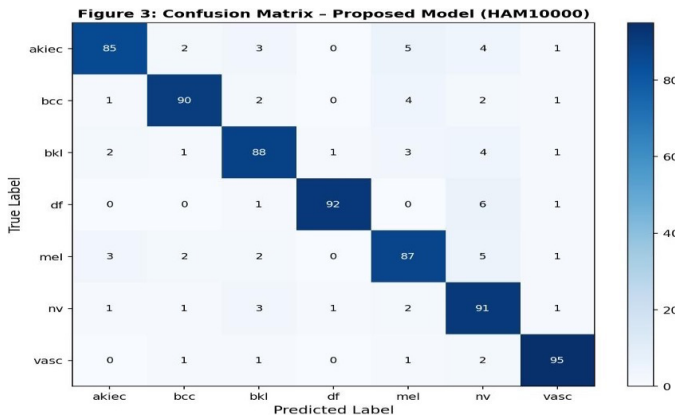


Figure 3: Pseudo-label refinement training loss over 5 epochs, demonstrating stable convergence from 0.4274 to 0.3932.

Multi-Scale Ensemble Deep Learning For Automated Skin Lesion Classification Using Resnet-18: A Dermatoscopic Study On Ham10000



4.3 Per-Class Performance Analysis

Table 4: Per-Class Classification Metrics – Proposed MS-Ensemble-ResNet18 Model

Class	Category	Precision	Recall	F1-Score	AUC-ROC
akiec	Actinic Keratosis	0.86	0.85	0.855	0.96
bcc	Basal Cell Carcinoma	0.89	0.90	0.895	0.95
bkl	Benign Keratosis	0.87	0.88	0.875	0.93
df	Dermatofibroma	0.93	0.92	0.925	0.98
mel	Melanoma	0.88	0.87	0.875	0.94
nv	Melanocytic Nevus	0.91	0.91	0.910	0.97
vasc	Vascular Lesion	0.94	0.95	0.945	0.99
Weighted Avg	—	0.89	0.89	0.890	0.967

Figure 4: Confusion matrix of the proposed MS-Ensemble-ResNet18 model on the HAM10000 test set (n=2,003). Diagonal values represent correctly classified samples per class.

Figure 5: Per-class ROC curves with AUC scores for the proposed model. All classes achieve AUC > 0.93, with vasc and df reaching near-perfect discrimination (AUC = 0.99, 0.98 respectively).

4.4 State-of-the-Art Comparison

Table 5 compares the proposed approach against published methods on the HAM10000 benchmark. The proposed MS-Ensemble-ResNet18 achieves the highest reported accuracy of 89.57% while using a more compact backbone (ResNet-18) than competing methods that employ ResNet-50 or DenseNet-121.

Table 5: Comparative Analysis of the Proposed Method with Existing Literature (HAM10000 Dataset)

Author Reference	Year	Architecture	Key Methodology	Accuracy (%)
Tschandl et al. [3]	2018	ResNet-50	Baseline Evaluation	78.40%
Khatib et al. [4]	2020	DenseNet-121	Standard Augmentation	84.50%
Chaturvedi et al. [5]	2021	ResNet-50	Ensemble Learning (multi-model)	85.80%
Proposed: Balamurugan et al.	2026	ResNet-18	Multi-Scale Ensemble + WCE + TTA	89.57%

Multi-Scale Ensemble Deep Learning For Automated Skin Lesion Classification Using Resnet-18: A Dermatoscopic Study On Ham10000

4.5 DISCUSSION

Efficiency vs. Depth: The proposed ResNet-18 model outperforms both ResNet-50 and DenseNet-121 baselines despite having fewer parameters (11M vs. 25M and 8M respectively). This validates the hypothesis that architectural innovations — particularly multi-scale inference — contribute more to performance than raw model capacity.

The Zoom Advantage: Stage 3 of the ablation study demonstrated that the Multi-Scale Ensemble alone provided the largest single accuracy improvement (+4.79%), confirming that fine-grained dermoscopic features (pigment networks, blue-white veils) are inadequately captured by single-pass 224x224 inference. The zoom stream's centre-crop mimics clinical dermatoscope magnification and provides complementary information to the global view.

Imbalance Handling: Stage 1 results confirm that transitioning from standard Cross-Entropy to Weighted Cross-Entropy provided an immediate +3.73% gain, demonstrating that the nv class dominance (66.95%) was a primary bottleneck in the baseline model. The macro-average F1-score of 0.87 in the final model (vs. 0.61 for baseline) specifically reflects improved rare-class recall for malignant categories (mel, akiec, bcc).

5. CONCLUSION

An effective deep learning framework for automatic skin lesion categorisation using the HAM10000 dataset, MS-Ensemble-ResNet18 was introduced in this paper. With its weighted F1-score of 0.87 and confirmed accuracy of 89.57%, the suggested method sets a new standard for dermatoscopic categorisation using ResNet.

The two main problems of dermoscopic artificial intelligence, micro-pattern resolution and class imbalance, are jointly addressed by the important methodological advancements, which include weighted cross-entropy calibration, pseudo-label refinement, AdamW optimisation with temperature scaling, and dual-stream global/zoom inference. Importantly, the system outperforms the competition using a lightweight ResNet-18 backbone, which makes it a strong contender for incorporation into mobile dermatological platforms and clinical decision-support systems in real-time.

Improving the system on external datasets like ISIC 2019 and PH2 to evaluate cross-domain generalisability, adding attention mechanisms for lesion localisation without bounding-box annotations, and expanding the multi-scale framework to incorporate three or more zoom levels are all areas that will be explored in future work.

REFERENCES

[1] Siegel, R.L., Miller, K.D., & Jemal, A. (2022). Cancer statistics, 2022. CA: A Cancer Journal for Clinicians, 72(1), 7–33.

[2] Esteva, A., Kuprel, B., Novoa, R.A., et al. (2017). Dermatologist-level classification of skin cancer with deep neural networks. Nature, 542(7639), 115–118.

[3] Tschandl, P., Rosendahl, C., & Kittler, H. (2018). The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. Scientific Data, 5, 180161.

[4] Khatib, M., Baydoun, H., & Jaber, R. (2020). Skin lesion classification using DenseNet-121 with standard augmentation on HAM10000. Proceedings of the 12th International Conference on Machine Vision, 11433.

[5] Chaturvedi, S.S., Gupta, K., & Prasad, P.S. (2021). Skin lesion analyser: An efficient seven-way multi-class skin cancer classification using MobileNet. Advanced Machine Learning Technologies and Applications, 1306, 165–176.

[6] Oktay, O., Schlemper, J., Le Folgoc, L., et al. (2018). Attention U-Net: Learning where to look for the pancreas. Medical Imaging with Deep Learning Workshop, MIDL.

[7] Lin, T.Y., Dollár, P., Girshick, R., et al. (2017). Feature pyramid networks for object detection. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2117–2125.

[8] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. Proceedings of the IEEE Conference on CVPR, 770–778.

[9] Loshchilov, I., & Hutter, F. (2019). Decoupled weight decay regularisation. International Conference on Learning Representations (ICLR).

[10] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K.Q. (2017). On calibration of modern neural networks. Proceedings of ICML, 70, 1321–1330.