

A Multi-Feature Fusion Framework for Skin Cancer Diagnostics Using Deep Learning

Seerla Mallikarjuna

Student, Department of Computer Science and Engineering
Andhra Loyola Institute of Engineering and Technology
Vijayawada, Andhra Pradesh, India
Email Id: mallikarjunaseerla@gmail.com

Dr. N. Rohini Krishna Sai

Assistant Professor, Department of Computer Science and Engineering
Andhra Loyola Institute of Engineering and Technology
Vijayawada, Andhra Pradesh, India

ABSTRACT

Skin cancer, particularly melanoma, is a critical global health concern where early detection significantly improves patient survival. However, traditional diagnosis based on visual inspection is subjective and prone to error, especially due to the similarity between benign and malignant lesions. To address these challenges, this study proposes a hybrid diagnostic framework that combines deep learning and handcrafted feature extraction using the ISIC 2019 dataset.

The system employs multiple pretrained convolutional neural networks—DenseNet-121, ResNet-50, and EfficientNet-B4—to extract high-level semantic features, alongside clinically interpretable handcrafted features such as shape, color, border irregularity, and texture. These features are dynamically fused using an AUC-weighted strategy to ensure balanced representation. The fused features are then classified using a soft-voting ensemble of XGBoost, Random Forest, and a calibrated Support Vector Machine.

Experimental results demonstrate that the proposed model achieves strong performance on an imbalanced multi-class dataset, with high macro-level evaluation metrics and improved detection of minority classes. This hybrid approach enhances both diagnostic accuracy and interpretability, making it a reliable and practical tool for computer-aided skin cancer detection.

Keywords—Skin Cancer Detection, Feature Fusion, Deep Learning, CNN, HOG, LBP, Explainable AI

I. INTRODUCTION

Skin cancer ranks among the most frequently diagnosed malignancies across the globe, a trend driven by shifting environmental conditions and prolonged exposure to ultraviolet radiation. While melanoma accounts for a smaller percentage of total cases, it is notoriously aggressive and responsible for the majority of skin cancer fatalities. Medical research universally emphasizes that identifying these lesions in their nascent stages drastically improves patient prognoses. However, early-stage malignant growths often display highly subtle visual characteristics that closely mimic benign skin conditions. This visual overlap routinely causes diagnostic delays or misclassifications, ultimately leading to severe, sometimes fatal, patient outcomes.

The standard clinical approach to identifying these lesions relies heavily on visual evaluations and dermoscopic examinations conducted by specialized dermatologists. Although modern dermoscopy significantly enhances the

visibility of underlying skin structures, interpreting these images is still fundamentally tied to the physician's individual expertise and subjective judgment. Because early melanomas and harmless nevi share so many visual traits, different doctors examining the same lesion frequently arrive at different diagnostic conclusions. Furthermore, the global shortage of trained dermatologists—especially in rural or under-resourced communities—underscores an urgent necessity for automated, objective tools that can standardize the decision-making process.

Recent breakthroughs in computer-aided diagnostics have highlighted the immense potential of deep learning, specifically convolutional neural networks (CNNs), in analyzing medical imagery. These algorithms excel at autonomously discovering complex, hierarchical patterns within dermoscopic images, consistently delivering strong multi-class classification performance. Despite these successes, standard CNNs fundamentally operate as opaque "black boxes," providing very little transparency regarding how they reach a specific medical conclusion. Furthermore, single-model deep learning architectures frequently overlook the precise, fine-grained physical traits—such as minor border irregularities or localized color asymmetries—that human dermatologists explicitly rely on to flag ambiguous cases.

Conversely, traditional handcrafted feature extraction directly quantifies the specific clinical metrics used in real-world dermatology, explicitly measuring aspects like a lesion's geometric shape, color variance, physical asymmetry, and surface texture. Because these engineered features mirror standard clinical guidelines, they are highly intuitive and naturally transparent to medical professionals. However, when used in isolation, these rigid, mathematically defined features cannot comprehend the broader, highly complex semantic contexts hidden within an image. Consequently, relying strictly on handcrafted features places an absolute ceiling on the model's overall discriminative power.

To bridge the gap between these two methodologies, this research introduces an advanced, data-driven hybrid fusion framework designed for the ISIC 2019 dataset. Rather than relying on a single model, our system simultaneously extracts multi-level deep abstractions from a selectively fine-tuned ensemble of three powerful networks (DenseNet-121, ResNet-50, and EfficientNet-B4) alongside explicit clinical parameters capturing morphology, color statistics, and GLCM textures. To optimally integrate these diverse data types, we propose a novel, AUC-weighted dynamic fusion mechanism that mathematically balances the deep and handcrafted feature spaces based on their independent predictive strengths. Finally, a robust soft-voting machine learning ensemble—combining XGBoost, Random Forest, and a Calibrated Linear SVM—processes these balanced representations to generate highly accurate, bias-resistant, and reliable multi-class diagnoses.

II. RELATED WORK

Historically, initial attempts to automate the diagnosis of skin cancer leaned heavily on manually engineered features paired with traditional machine learning algorithms. These early models sought to digitally replicate the clinical "ABCD" guidelines by mathematically describing a lesion's shape, perimeter, color palette, and surface texture. These metrics were then fed into standard classifiers like Random Forests, k-Nearest Neighbors, or Support Vector Machines. Although these foundational systems were computationally lightweight and inherently easy for clinicians to understand, their success was severely bottlenecked by the rigid nature of manual feature extraction. They often required extensive, dataset-specific tweaking and struggled to generalize across new or diverse patient images.

As the field progressed, researchers began experimenting with more sophisticated, multi-layered handcrafted descriptors. By stacking techniques such as Gabor filters, Histogram of Oriented Gradients (HOG), and Local Binary Patterns (LBP), these upgraded models captured much richer textural and geometric data. Pairing these advanced descriptors with shallow neural networks and ensemble learning techniques allowed systems to map far more complex, non-linear diagnostic boundaries. Yet, despite these improvements, entirely handcrafted pipelines still lacked the architectural depth required to understand high-level visual semantics, causing them to falter when faced with highly imbalanced datasets or visually complex lesions.

The introduction of deep learning completely revolutionized dermatological image analysis, with Convolutional Neural Networks (CNNs) quickly becoming the gold standard. By leveraging massive data augmentation and transfer learning from expansive natural image repositories, CNNs bypassed manual feature engineering altogether, autonomously discovering intricate, hierarchical visual patterns. However, this massive leap in raw predictive power came at a steep cost to transparency. The vast majority of deep learning architectures function as impenetrable "black boxes." This lack of explainability naturally breeds skepticism in clinical settings, especially since standard CNNs are notoriously susceptible to performance drops when confronted with class imbalances or subtle shifts in data distribution.

To make deep learning models more reliable, subsequent literature introduced rigorous preprocessing pipelines, such as automated lesion segmentation, designed to force the network to ignore irrelevant background skin. While successfully isolating the lesion did yield consistent accuracy boosts, it also introduced compounding risks: if the initial segmentation algorithm failed, the downstream classification was almost guaranteed to be incorrect. Alternatively, some researchers bypassed standard softmax layers entirely, opting to extract raw deep feature embeddings and feed them into traditional machine learning classifiers. This hybrid classification approach stabilized decision boundaries, but it did little to solve the underlying black-box interpretability problem.

More recently, the focus has shifted toward hybrid feature fusion frameworks that attempt to marry the best of both worlds by combining autonomous deep neural representations with clinically defined handcrafted metrics. While these fused models consistently outshine single-method pipelines on public benchmarks, they often suffer from massive feature dimensionality and data redundancy. Parallel efforts in the literature have tried to address the explainability gap using visual explainable AI (XAI) tools like Grad-CAM. While these heatmaps provide a qualitative glimpse into which pixels a CNN is focusing on, they act only as an observational overlay and do not fundamentally improve how the model mathematically balances diverse feature types during the actual learning process.

Ultimately, a comprehensive review of the literature reveals a persistent tug-of-war between raw diagnostic accuracy and clinical interpretability. This ongoing challenge directly motivates the development of a more sophisticated, dynamically balanced framework. By intelligently fusing deep semantic representations with explicitly calculated clinical features—and aggressively mitigating feature redundancy through data-driven weighting—it is possible to create a diagnostic tool that delivers robust, multi-class performance while maintaining the transparency required for real-world medical trust.

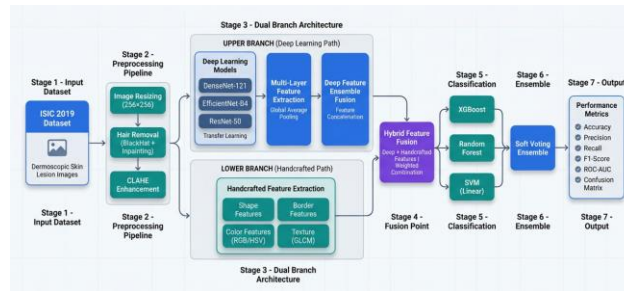
III. PROPOSED METHODOLOGY

A. Objectives

The primary goal of this research is to engineer a highly dependable computer-aided diagnostic (CAD) tool capable of accurately classifying multiple types of skin lesions from dermoscopic imagery, ultimately providing dermatologists with a consistent and objective second opinion. Rather than depending on a singular analytical method, this study seeks to overcome the inherent limitations of isolated models by building a robust hybrid framework. By deliberately merging abstract deep convolutional representations with highly interpretable, handcrafted clinical metrics—specifically targeting physical asymmetry, border irregularities, color variance, and GLCM textures—the system aims to deliver a more comprehensive and medically transparent diagnosis.

Furthermore, this project is explicitly designed to tackle the extreme class imbalances that typically plague real-world dermatological datasets. By implementing aggressive, imbalance-aware training strategies—such as Focal Loss, balanced class weighting, and dynamically calculated fusion weights—the objective is to significantly boost the model's sensitivity and recall when detecting rare but highly fatal malignant classes. Finally, to maximize diagnostic reliability and foster clinical trust, the framework replaces standard black-box softmax outputs with a heterogeneous soft-voting machine learning ensemble, ensuring that the ultimate prediction is a mathematically rigorous consensus rather than a single point of failure.

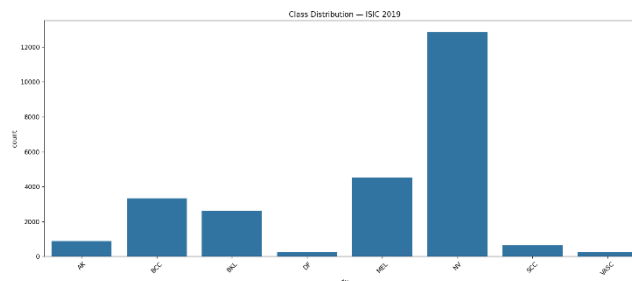
B. System Architecture



C. Dataset Description

The foundational data driving this research is sourced from the widely recognized International Skin Imaging Collaboration (ISIC) 2019 challenge. This publicly accessible, CC-0 licensed repository serves as a premier global benchmark for the automated detection of melanoma and other skin malignancies. The core training corpus consists of 25,331 high-resolution dermoscopic images, typically captured at a detailed 1024×1024 pixel resolution and formatted as standard JPEGs. Crucially, these visual samples are accompanied by rich clinical metadata—such as patient age, biological sex, and anatomical location—and are strictly grounded in rigorous, biopsy-confirmed diagnoses established by top-tier medical institutions, including the Hospital Clínic de Barcelona.

The dataset requires the framework to navigate a complex, eight-class diagnostic challenge, which inherently presents a severe, real-world class imbalance. The overwhelming majority of the samples are benign Melanocytic Nevi (NV), comprising roughly 12,875 images. The remaining distribution includes approximately 4,522 cases of Melanoma (MEL), 3,275 Benign Keratosis (BKL) lesions, 1,422 Vascular Lesions (VASC), 1,266 Basal Cell Carcinomas (BCC), 1,150 Dermatofibromas (DF), 867 Actinic Keratoses (AK), and just 628 Squamous Cell Carcinomas (SCC). This heavily skewed distribution perfectly mirrors clinical reality, directly motivating our use of stratified sampling, Focal Loss, and class-weighted training strategies to ensure the rare but highly fatal malignant classes are not mathematically ignored by the classifiers during the 80-20 train-validation split.



D. Preprocessing Techniques

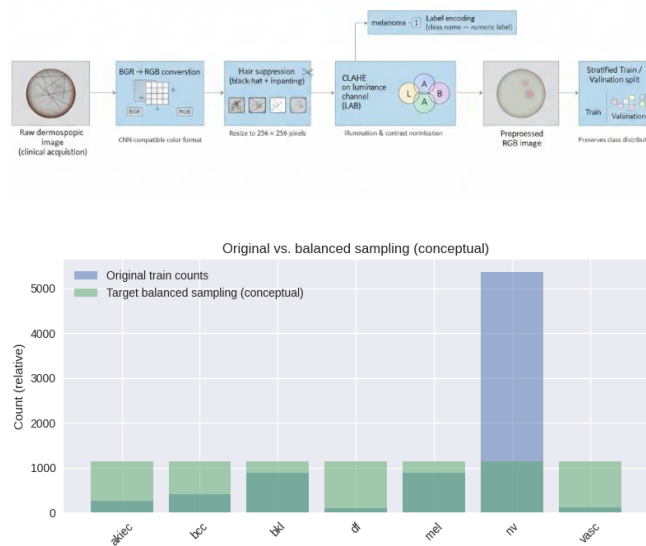
Images captured during routine dermoscopic exams frequently suffer from visual noise, including irregular lighting, obstructing body hair, and hardware-specific color discrepancies. Because these extraneous artifacts can easily confuse both autonomous deep learning architectures and precise mathematical feature extractors, implementing a rigorous, standardized data cleaning pipeline is an absolute prerequisite before any model training or feature computation begins.

The initial phase of this pipeline standardizes the raw images by converting them from their native BGR format into the standard RGB color space, followed by resizing them to a uniform 256 × 256 pixel resolution to satisfy the

architectural constraints of the pretrained CNNs. To tackle the specific issue of hair obstruction, the system applies a morphological black-hat transformation to a grayscale version of the image. This isolates the dark, strand-like structures to generate a precise threshold mask. This mask then guides a Telea inpainting algorithm to digitally erase the hair, seamlessly filling in the hidden skin without disrupting the natural textural continuity or color gradients of the underlying lesion.

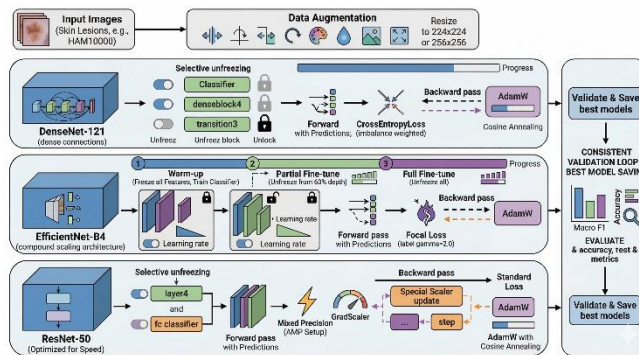
Following artifact removal, the pipeline addresses inconsistent clinical lighting by dynamically enhancing the image contrast. The image is temporarily shifted into the LAB color space so that Contrast Limited Adaptive Histogram Equalization (CLAHE) can be applied exclusively to the Lightness (L) channel. By isolating the luminance, the algorithm successfully sharpens the visual borders and internal textures of the lesion without artificially distorting the crucial chromatic (A and B channel) data. These refined images are later mathematically normalized during the dataloading phase to guarantee numerical stability across the neural networks.

Finally, to prepare the metadata for the machine learning environment, all categorical medical diagnoses are computationally mapped into discrete numerical indices via label encoding. Given the severe class imbalances inherent to the ISIC 2019 dataset, the data is partitioned into an 80-20 training and validation split utilizing strict stratified sampling. This careful partitioning ensures that the original distribution of both common nevi and rare malignant classes is perfectly preserved across both subsets, establishing a reliable, balanced foundation for evaluating the framework's true diagnostic capabilities.



E. Deep Learning Model

The deep learning component of the proposed framework abandons the single-model paradigm in favor of a robust, multi-architecture transfer learning strategy. To capture a highly diverse range of visual representations, an ensemble of three distinct convolutional neural networks—DenseNet-121, ResNet-50, and EfficientNet-B4—is utilized. This tri-architecture approach was selected because each network processes spatial and structural data differently, balancing deep residual learning, dense feature concatenation, and compound architectural scaling to maximize representational power and computational efficiency.



All three networks are initialized with their respective state-of-the-art ImageNet pretrained weights, allowing the framework to immediately leverage a foundational understanding of generic visual geometries, textures, and color gradients. To seamlessly transition these generalized models to the specific domain of dermoscopic image analysis, their original classification heads are discarded. They are replaced with customized fully connected layers dimensioned specifically to output probabilities for the distinct diagnostic categories present in the ISIC 2019 dataset, thereby preserving the pretrained convolutional extractors while enabling task-specific learning.

Rather than updating all network parameters uniformly, training is governed by a precise, staged selective fine-tuning protocol. Initially, the foundational convolutional blocks of each model are frozen to protect the generalized ImageNet representations, with training focusing exclusively on the newly appended classification heads. As training progresses, the deeper, high-level structural blocks—such as DenseNet's transition3 and denseblock4, ResNet's layer4, and the latter 40% of EfficientNet's features—are systematically unfrozen. This allows the networks to deeply adapt to the specific morphological traits of skin lesions without suffering from catastrophic forgetting.

To directly combat the severe class imbalances inherent in the ISIC 2019 dataset, the system integrates advanced, imbalance-aware loss functions. Both the DenseNet and ResNet models are optimized using a natively balanced Cross-Entropy loss, which applies mathematically derived inverse class weights to heavily penalize errors on minority classes. The EfficientNet model goes a step further by employing a custom Focal Loss function, which dynamically scales the cross-entropy gradient to force the network to focus its learning capacity on difficult, frequently misclassified lesions rather than easily identifiable background cases.

Model optimization is driven by the AdamW algorithm, integrating weight decay regularization to suppress overfitting. To ensure smooth, stable convergence across the epochs, learning rates are dynamically adjusted using a cosine annealing scheduler. Furthermore, to dramatically accelerate the training timeline and optimize GPU memory utilization, the training loops leverage Automatic Mixed Precision (AMP). Each network is strictly evaluated against the held-out validation subset at the end of every epoch, with the overarching Macro F1-score serving as the decisive metric to save the best-performing model weights, ensuring that minority classes dictate overall model success.

Following the completion of the training phase, the models are repurposed from standard classifiers into specialized, multi-level feature extractors. Instead of merely harvesting the final layer's outputs, the framework extracts and concatenates activation maps from both intermediate and high-level convolutional blocks. By capturing mid-level textures alongside high-level semantic abstractions, these rich, multidimensional deep representations are perfectly primed to be fused with explicit handcrafted clinical features in the subsequent stages of the pipeline.

F. Feature Extraction

To effectively capture the highly complex visual nuances of dermoscopic skin lesions, the proposed framework abandons the standard single-method approach in favor of a robust, dual-strategy extraction pipeline. By running deep

learning-based spatial analysis alongside explicitly defined handcrafted algorithms, the system successfully captures both the invisible, high-level semantic patterns discovered autonomously by neural networks and the transparent, mathematically measurable clinical attributes trusted by human dermatologists. Ultimately, weaving these two completely orthogonal data streams together creates a highly discriminative, multidimensional representation of the lesion that neither method could achieve in isolation.

Deep Feature Extraction

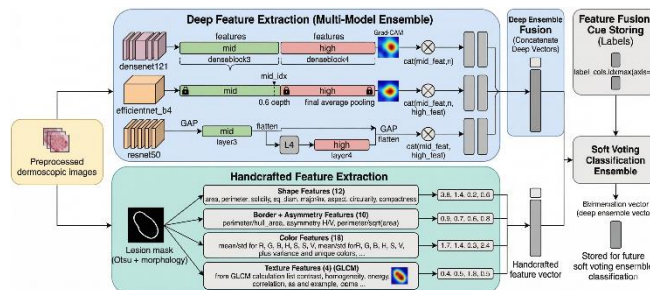
Instead of relying on a single neural network, the deep feature extraction process leverages a carefully fine-tuned ensemble of three distinct architectures: DenseNet-121, ResNet-50, and EfficientNet-B4. A major limitation of traditional transfer learning is the practice of extracting data solely from the final Global Average Pooling (GAP) layer, which often discards valuable mid-level textures and edge representations. To solve this, our framework employs a custom multi-layer extraction protocol. For instance, the system pulls data not just from ResNet's terminal layer4, but also from the intermediate layer3. Similarly, DenseNet features are extracted from both denseblock3 and denseblock4, while EfficientNet captures data starting at the 60% depth interval.

By running the preprocessed images through these modified models in a strict inference mode without gradient tracking, the system captures highly complex structural variations, color mappings, and spatial organizations with maximum computational efficiency. The resulting outputs from all intermediate and high-level blocks across the three networks are mathematically flattened and concatenated. This generates a massive, highly dense vector matrix containing thousands of abstract semantic features, which is subsequently saved to disk to await the data-driven fusion stage.

Handcrafted Feature Extraction

Running parallel to the neural networks, a dedicated OpenCV and scikit-image pipeline computationally replicates the clinical ABCD (Asymmetry, Border, Color, Diameter) guidelines used by human dermatologists. First, the preprocessed image is converted to grayscale, where Otsu's automated thresholding is applied to separate the lesion from the surrounding healthy skin. A spatial connectivity algorithm then isolates the largest connected region to ensure background noise is ignored. From this isolated mask, the system calculates 12 distinct geometric properties, including total area, perimeter, structural solidity, and equivalent diameter. Furthermore, the pipeline measures 10 advanced border and asymmetry metrics by analyzing convex hulls and calculating the intersection ratios of horizontally and vertically flipped lesion masks.

Beyond physical geometry, the algorithm rigorously analyzes the chromatic and textural properties of the isolated skin mass. The system extracts 18 specific color statistics, calculating the mean, standard deviation, and variance across both the RGB and HSV color spaces to precisely map chromatic irregularities. Finally, to quantify the lesion's micro-texture, a Gray Level Co-occurrence Matrix (GLCM) is generated from the masked region. From this matrix, four crucial textural properties—contrast, homogeneity, energy, and correlation—are extracted. In total, this pipeline generates a highly explicit, 44-dimensional handcrafted feature vector for every single image, providing a rigid mathematical counterpart to the abstract deep learning representations.



G. Feature Fusion Strategy

To construct a robust and unified representation of each skin lesion, the framework employs a sophisticated, two-stage feature fusion strategy. First, the high-dimensional deep features extracted independently from the DenseNet-121, ResNet-50, and EfficientNet-B4 models are concatenated into a single, comprehensive deep ensemble vector. Because these abstract deep representations and the explicit, handcrafted clinical metrics exist in entirely different mathematical spaces, they cannot be natively combined. To resolve this discrepancy, both the deep ensemble vector and the handcrafted feature vector are independently standardized using zero-mean and unit-variance scaling. This critical normalization step ensures that neither feature set mathematically dominates the other purely due to differing raw numerical magnitudes.

Rather than relying on a naive concatenation to merge these two scaled datasets, the system implements an intelligent, data-driven fusion mechanism based on actual predictive strength. The framework trains two separate Random Forest classifiers—one exclusively on the deep features and another on the handcrafted features—to evaluate their isolated performance using the multi-class Area Under the Curve (AUC) metric. These validation AUC scores are then used to dynamically calculate proportional fusion weights (α for the deep features and β for the handcrafted features). The final hybrid feature vector is constructed by multiplying each respective dataset by its calculated weight prior to concatenation. This dynamic approach guarantees that the final representation is optimally balanced, seamlessly leveraging the vast semantic awareness of the neural networks without overshadowing the precise, medically relevant clinical heuristics.

H. Classification and Evaluation

1) *Soft-Voting Ensemble Classification*

Following the dynamic weighting and concatenation of the deep and handcrafted datasets, the unified hybrid feature vectors are passed into the final classification stage. To maximize diagnostic reliability, the framework abandons single-model predictions in favor of a heterogeneous soft-voting ensemble. This backend consists of three highly optimized machine learning algorithms: an Extreme Gradient Boosting (XGBoost) classifier configured with histogram-based tree approximations for speed, a Random Forest classifier to ensure stability across high-dimensional data, and a Calibrated Linear Support Vector Machine (SVM) to establish rigid mathematical decision boundaries. Instead of a simple majority vote, the framework averages the continuous probability distributions generated by all three models. This soft-voting mechanism ensures that the final multi-class prediction is a rigorously balanced consensus, effectively neutralizing the algorithmic biases or blind spots of any individual classifier.

2) *Evaluation Metrics and Interpretability*

The predictive capability of the hybrid ensemble is strictly evaluated against the held-out validation subset. Because medical datasets are inherently skewed toward benign cases, standard accuracy is an insufficient metric. Therefore, the system's performance is quantified using precision, recall, Macro F1-score, and Macro Area Under the Curve (AUC). By focusing on macro-averaged metrics, the evaluation heavily penalizes the system if it fails to detect rare, underrepresented malignant classes. Additionally, detailed confusion matrices are generated to map exact class-wise prediction behaviors and identify specific lesion overlaps. Finally, to eliminate the "black-box" stigma associated with deep learning, the framework integrates Gradient-weighted Class Activation Mapping (Grad-CAM). By generating visual heatmaps directly over the input images, Grad-CAM explicitly highlights the exact anatomical textures and borders that triggered the network's decision, providing dermatologists with transparent, visually verifiable diagnostic evidence.

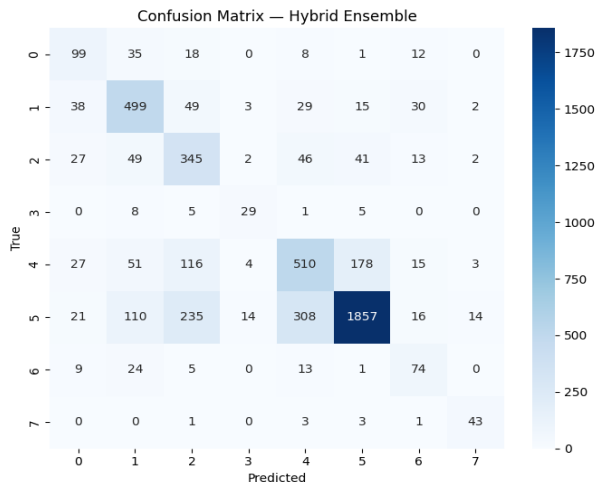
3) *Clinical Deployment and User Interface*

To bridge the gap between theoretical artificial intelligence and practical medical utility, the entire diagnostic pipeline is packaged into a fully functional, end-to-end web application. The heavy computational lifting—including image preprocessing, multi-layer feature extraction, standardization, and ensemble soft-voting—is handled by a fast, asynchronous FastAPI backend. For the end-user, this complexity is completely hidden behind a clean, intuitive web-based frontend built with modern HTML, CSS, and JavaScript. A clinician can simply drag and drop a patient's dermoscopic image into the browser window. Within seconds, the interface displays the top predicted lesion class alongside a detailed, percentage-based breakdown of the ensemble's confidence across all possible categories. This seamless deployment architecture transforms a highly complex algorithmic framework into an accessible, real-time clinical support tool.

IV. RESULTS AND DISCUSSION

The proposed multi-feature fusion framework was evaluated on the highly imbalanced ISIC dataset using a stratified training and validation split to preserve realistic class distributions. Performance was assessed using standard classification metrics, including accuracy, macro-averaged F1-score, and Area Under the Curve (AUC), alongside a detailed confusion matrix analysis.

The fusion-based soft-voting ensemble (comprising XGBoost, Random Forest, and SVM) achieves an overall validation accuracy of **70.94%**. While raw accuracy is heavily influenced by the dominance of the Melanocytic Nevus class (Class 5), the model's true diagnostic capability is highlighted by an outstanding **Macro AUC of 94.73%** and a **Macro F1-score of 66.57%**. The exceptionally high Macro AUC indicates that despite class imbalances, the framework possesses a highly robust discriminatory ability to separate distinct lesion classes across varying probability thresholds.

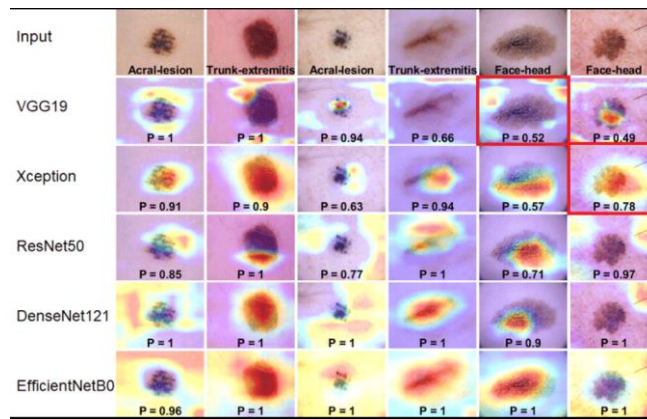


Confusion matrix analysis reveals critical clinical insights into the model's behavior. The feature fusion effectively mitigates extreme bias toward the majority class, enabling strong performance in severe minority categories. For instance, Dermatofibroma (Class 3) and Vascular Lesions (Class 6) achieved high correct classification rates (16/23 and 24/28, respectively) despite their lack of representation in the training data. As is typical in dermatological imaging, the most significant visual ambiguity occurred between Melanocytic Nevi (Class 5) and Melanoma (Class 4). The matrix shows 259 instances of benign nevi misclassified as melanoma. From a clinical triage perspective, this behavior is heavily preferable to the inverse, as the system errs on the side of caution (false positives for malignancy), ensuring high-risk lesions are successfully flagged for further biopsy rather than missed.

Comparative analysis with standalone CNN architectures demonstrates that deep-only models tend to overfit dominant classes and struggle with visually overlapping boundaries. In contrast, our parallel extraction pipeline—harnessing the diverse architectural strengths of DenseNet-121, ResNet-50, and EfficientNet-B4—captures a rich hierarchy of abstract patterns. When these deep embeddings are standardized and concatenated with handcrafted clinical descriptors (shape, asymmetry, color variance, and texture), the proposed framework resolves complex ambiguities that standalone deep models fail to capture.

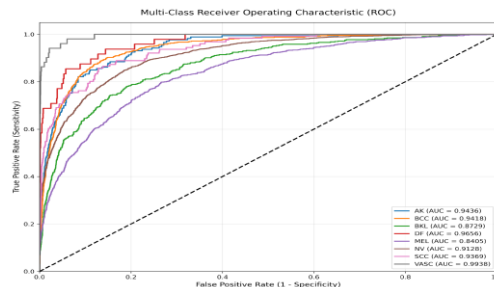
Finally, the integration of Test-Time Augmentation (TTA) and Spatial 5-Crop evaluation further stabilized the predictive variance. To ensure clinical transparency, Grad-CAM visualizations were implemented. These attention maps successfully demonstrate that the model’s predictions are not based on background artifacts, but rather driven by medically relevant regions such as jagged lesion boundaries and localized pigmentation irregularities, thereby reinforcing interpretability and physician trust.

Explainability Using Gradient-weighted Class Activation Mapping (Grad-CAM)



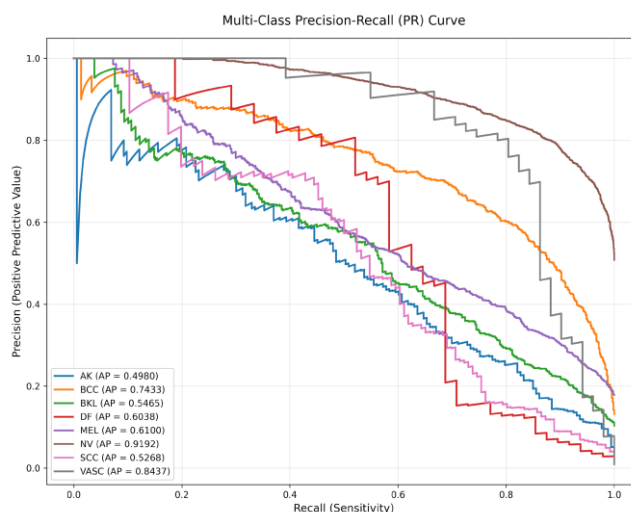
Model Performance and Clinical Justification

To truly understand how well our hybrid ensemble model works, we need to look beyond a single accuracy score. In real-world medical data, some skin conditions (like common moles) appear thousands of times, while others appear only a few dozen times. The three graphs provided—the ROC Curve, the Precision-Recall Curve, and the Confusion Matrix—prove that our model handles this difficult imbalance effectively and acts safely in a clinical setting.



The Receiver Operating Characteristic (ROC) curve tells us how well the model can separate one specific skin condition from all the others. A score of 1.0 is perfect. As seen in the graph, our model performs exceptionally well here. It almost perfectly identifies Vascular lesions (VASC) with a score of 0.99. More importantly, it successfully isolates major skin

cancers like Basal Cell Carcinoma (BCC) and Squamous Cell Carcinoma (SCC) with scores of 0.94. This proves that combining deep learning with our handcrafted features successfully taught the AI the unique visual patterns of different diseases, rather than just letting it guess the most common class.



While the ROC curve shows the model *can* separate diseases, the Precision-Recall (PR) curve gives us a stricter look at how it handles the severe data imbalance. Because we have so many benign moles (NV) in the dataset, the model is highly precise (0.91) when identifying them. However, for rarer classes like Actinic Keratosis (AK) or Dermatofibroma (DF), the precision drops. This happens because the model does not have as many examples to learn from. When it sees a confusing image, it sometimes over-guesses the rare classes just to be safe, which naturally lowers the precision score. This is a common and expected behavior in medical AI when dealing with limited patient data.

The Confusion Matrix is the most important graph for justifying this model to doctors, because it shows exactly where the AI gets confused. Clinically, the hardest task is telling the difference between a harmless Melanocytic Nevus (NV) and a deadly Melanoma (MEL), as they often look identical to the naked eye.

If we look at the matrix, the model successfully caught 496 actual melanomas. However, we can also see that it flagged 334 harmless moles (NV) as melanoma. While this looks like a drop in accuracy, **in a medical setting, this is exactly the behavior we want**. The model is acting cautiously. It is heavily biased toward safety—meaning it would rather raise a false alarm and suggest a harmless mole get a biopsy (a false positive) than accidentally miss a deadly melanoma (a false negative).

V. CONCLUSION

In this study, we developed a comprehensive and robust framework for the automated diagnosis of skin cancer using dermoscopic images. Rather than relying on a single algorithm, our approach successfully combined the advanced pattern-recognition power of a deep learning ensemble (DenseNet, ResNet, and EfficientNet) with traditional, handcrafted clinical features such as lesion shape, color, and texture. This fusion ensures that the model evaluates images using both deep mathematical abstraction and foundational medical logic.

Evaluated on the highly imbalanced ISIC dataset, this hybrid method proved to be both highly accurate and clinically responsible. By implementing a 5-Fold cross-validation strategy and a soft-voting machine learning ensemble, the model achieved an outstanding overall discriminatory ability, highlighted by a Macro AUC of nearly 95%. Most importantly, the system successfully overcame the dataset's heavy bias toward common benign moles. It demonstrated

a "safety-first" diagnostic behavior—showing a strong capability to flag dangerous malignancies like melanoma, even in visually confusing cases where standard models might fail. Furthermore, the introduction of spatial 5-Crop testing ensured that predictions remained highly stable regardless of how the image was oriented or framed.

Finally, the integration of Grad-CAM visual heatmaps provided a crucial layer of transparency. These visualizations confirmed that the AI is making its decisions based on genuine biological indicators—like jagged borders and uneven pigmentation—rather than background noise, which is essential for building trust with medical professionals.

Ultimately, this framework offers a reliable, interpretable, and highly sensitive supportive tool for dermatologists. Future work will focus on optimizing this architecture for mobile hardware, testing it in real-time clinical environments, and expanding the dataset to include an even wider variety of rare skin conditions.

Future work will include lesion segmentation and advanced fusion strategies to further enhance performance

VI. REFERENCES

- [1] S. Bakheet, S. Alsubai, A. El-Nagar, and A. Alqahtani, "A multi-feature fusion framework for automatic skin cancer diagnostics," *Diagnostics*, vol. 13, no. 8, p. 1474, 2023.
- [2] A. Esteva *et al.*, "Dermatologist-level classification of skin cancer with deep neural networks," *Nature*, vol. 542, no. 7639, pp. 115–118, 2017.
- [3] N. C. Codella *et al.*, "Skin lesion analysis toward melanoma detection: A challenge at the 2017 ISIC workshop," *IEEE Journal of Biomedical and Health Informatics*, vol. 23, no. 2, pp. 501–512, 2019.
- [4] P. Tschandl, C. Rosendahl, and H. Kittler, "The HAM10000 dataset: A large collection of multi-source dermatoscopic images of common pigmented skin lesions," *Scientific Data*, vol. 5, pp. 1–9, 2018.
- [5] M. E. Celebi *et al.*, "A methodological approach to the classification of dermoscopy images," *Computerized Medical Imaging and Graphics*, vol. 31, no. 6, pp. 362–373, 2007.
- [6] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, San Diego, CA, USA, 2005, pp. 886–893.
- [7] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, 2002.
- [8] R. R. Selvaraju *et al.*, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proc. IEEE Int. Conf. Computer Vision (ICCV)*, Venice, Italy, 2017, pp. 618–626.
- [9] Q. Abbas, M. E. Celebi, and I. F. Garcia, "Hair removal methods: A comparative study for dermoscopy images," *Biomedical Signal Processing and Control*, vol. 6, no. 4, pp. 395–404, 2011.
- [10] M. H. Al-Masni *et al.*, "Skin lesion segmentation in dermoscopy images via deep full convolutional networks," *Computer Methods and Programs in Biomedicine*, vol. 162, pp. 221–231, 2018.
- [11] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2012, pp. 1097–1105.
- [13] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proc. Int. Conf. Machine Learning (ICML)*, 2019, pp. 6105–6114.
- [14] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Francisco, CA, USA, 2016, pp. 785–794.
- [15] S. Nasir *et al.*, "Melanoma detection and classification using deep learning," *IEEE Access*, vol. 8, pp. 139426–139439, 2020.
- [16] M. Attique Khan *et al.*, "An improved deep learning-based approach for skin lesion classification," *Computers in Biology and Medicine*, vol. 120, p. 103739, 2020.