

Hybrid Swarm-Driven Vision Transformer (HSViT) for Lung Cancer Segmentation and Classification from CT Scans

V. Kavithamani¹, V. Kavya¹, R. Suganthi², Yuvaraj S.³, P. Monisha⁴, Arun Patrick K.⁵

¹ Department of Electronics and Communication Engineering, Jai Shriram Engineering College, Tamilnadu, India.

² Department of Electronics and Communication Engineering, Panimalar Engineering College, Tamil Nadu, India.

³ Department of Computer Science and Engineering at Sri Eshwar College of Engineering, Coimbatore, India.

⁴ Department of ECE, Karunya Institute of Technology and Sciences, Coimbatore, Tamilnadu, India

⁵ Department of Computer Science and Engineering, Nehru Institute of Technology, Jawahar Gardens, Tamil Nadu, India.

Corresponding author: Ms.V.Kavithamani (e-mail: kvkavirasu@gmail.com), **Author(s) Email:** Kavya V(e-mail: kaviadhanalakshmi@gmail.com), Dr.R.Suganthi (e-mail : sugimanicks@gmail.com), Yuvaraj S (e-mail: yuvaraj.scse@sece.ac.in), P.Monisha (e-mail: pmonisha123@gmail.com), K Arun Patrick (e-mail: nitarunpatrick@nehrucolleges.com)

ABSTRACT Lung cancer segmentation and classification from computed tomography (CT) images play a vital role in early diagnosis, prognosis assessment, and effective treatment planning. Despite significant progress in medical image analysis, accurate lung lesion analysis remains highly challenging due to overlapping anatomical structures, heterogeneous tissue intensity distributions, irregular and complex tumor shapes, and poorly defined lesion boundaries. These factors often limit the reliability and generalization capability of conventional deep learning models when applied to real-world clinical data. To address these challenges, this paper proposes a Hybrid Swarm-Driven Vision Transformer (HSViT) framework that synergistically combines swarm intelligence with transformer-based deep learning. The processing pipeline begins with Contrast Limited Adaptive Histogram Equalization (CLAHE), which enhances local contrast while suppressing noise amplification, thereby improving the visibility of subtle pulmonary nodules and lesion regions. Subsequently, a U-Net segmentation model optimized using the Coyote Optimization Algorithm (COA) is employed to accurately delineate lung lesions. COA, a swarm-based metaheuristic, adaptively fine-tunes U-Net parameters, enabling improved convergence and more precise boundary detection compared to gradient-based optimization alone. Following segmentation, discriminative lesion features are extracted and passed to the HSViT classifier. The proposed classifier integrates a Dual-Stage Attention Fusion (DSAF) mechanism, which effectively captures both fine-grained local spatial features and long-range global contextual dependencies. The framework achieves a Dice Coefficient of 0.95, an overall classification accuracy of 98.7%, and a minimized training loss of 0.04. These results highlight the strong potential of HSViT for reliable automated lung cancer diagnosis and for supporting clinical decision-making systems in real-world healthcare environments.

Keywords Hybrid Swarm Driven Vision Transformer, Coyote Optimization Algorithm, Vision Transformer, Dual Stage Attention Fusion, Lung Cancer Segmentation.

1. Introduction

Lung cancer is among the most common and fatal diseases globally, and it contributes significantly to cancer-related mortality. Timely and accurate detection is critical for improving patient survival and informing successful treatment plans. Computed Tomography (CT) is the most widely used modality for diagnosing lung cancer because of its high spatial resolution and ability to capture fine anatomical details. But manual

reading of CT scans is labor-intensive, error-prone, and variable, pointing to the necessity of smart automated systems for accurate segmentation and classification of lung lesions [1]. As medical imaging technology continues to advance, detecting early-stage lung tumors remains a significant challenge. These tumors often exhibit very subtle texture variations, irregular shapes, and low contrast, making them visually similar to healthy tissues in CT scans [2]. Traditional machine

learning algorithms and conventional image processing techniques struggle to achieve high accuracy due to their limited feature extraction capabilities and inability to capture complex spatial patterns [3]. Consequently, there is a growing need for a powerful, fully automated diagnostic platform capable of handling intricate lung structures, accurately differentiating malignant nodules from benign ones, and providing consistent, reliable results with minimal human intervention [4]. Various deep learning models have been extensively explored for lung cancer segmentation and classification tasks using Computed Tomography (CT) images. Convolutional Neural Networks (CNNs) have been widely adopted due to their strong ability to learn hierarchical, localized spatial features via convolutional filters. CNN-based models can effectively capture texture, edge, and intensity variations within lung nodules, making them suitable for early-stage lesion detection. However, their reliance on local receptive fields limits their ability to model long-range contextual relationships, which are often crucial for distinguishing malignant nodules from benign structures spanning multiple CT slices. To address these limitations, U-Net-based encoder-decoder architectures have become a popular choice for medical image segmentation. U-Net models employ skip connections between encoder and decoder layers, enabling the preservation of fine-grained spatial details while progressively learning high-level semantic representations. [5]. This dependency often limits their adaptability, as the models tend to overfit to specific datasets and struggle to generalize well across diverse clinical imaging conditions.

Moreover, their sensitivity to parameter tuning and variations in data acquisition protocols poses additional challenges, reducing their reliability and practicality in real-world clinical environments. The inherent complexities of lung cancer detection further complicate the task; these include low contrast between lesions and surrounding tissues, diverse and heterogeneous tumor appearances, and significant variation in lesion size, shape, and location [6]. Such factors make robust segmentation and classification difficult, highlighting the need for more adaptive, efficient, and clinically scalable deep learning approaches.

Moreover, noise and artifacts in CT images also hinder segmentation accuracy. Deep models, as efficient as they are, tend to demand high computational costs and poor convergence when handling imbalanced datasets [7]. These difficulties, taken together, prevent the establishment of a universally efficient and reliable diagnostic model.

In order to overcome these constraints, optimization-based hybrid models have become popular for improving model performance and adaptability. Swarm intelligence algorithms like Particle

Swarm Optimization (PSO), Grey Wolf Optimizer (GWO), and Coyote Optimization Algorithm (COA) have been seen to exhibit promise in optimizing deep learning parameters to enhance convergence and avoid local minima problems. Merging such metaheuristic optimization methods with attention-based architectures, such as Vision Transformers, will enable the balance of global and local feature learning effectively, resulting in enhanced segmentation and classification accuracy. The main contributions of the proposed work are listed below.

- a) Hybrid Optimization Transformer Framework is a novel Hybrid Swarm-Driven Vision Transformer (HSViT) that integrates the Coyote Optimization Algorithm with the Vision Transformer for precise lung cancer segmentation and classification.
- b) Adaptive Segmentation with Optimized U-Net employs a Coyote-Optimized U-Net model for adaptive and accurate extraction of lung lesion boundaries by dynamically tuning network parameters.
- c) Enhanced Feature Learning mechanism incorporates a Dual-Stage Attention Fusion mechanism in the Vision Transformer to effectively capture both global contextual information and fine-grained local features from CT images.
- d) It utilizes Contrast Limited Adaptive Histogram Equalization (CLAHE) to enhance image contrast and improve the visibility of tumor regions before segmentation.
- e) The proposed work gives a Dice Coefficient of 0.95 and an overall accuracy of 98.7% than existing models.

The remainder of the paper is organized as follows: Section 2 reviews related work on lung cancer detection and segmentation using deep learning and optimization-based methods. Section 3 details the proposed framework, including pre-processing, segmentation, and classification stages. Section 4 presents the comparative analysis and results obtained against existing methods. Finally, Section 5 discusses comparative analysis with the existing model in detail, and finally, Section 6 concludes with key findings, limitations, and potential directions for future research.

II. State-of-The-Art Techniques

[8] suggested a powerful lung cancer diagnosis model based on pre-trained CNNs to learn hierarchical features from CT images. The approach utilizes transfer learning to enhance classification accuracy with a decreased training time, and proves that pre-trained models can provide robust lung cancer detection with small datasets. [9] presented the

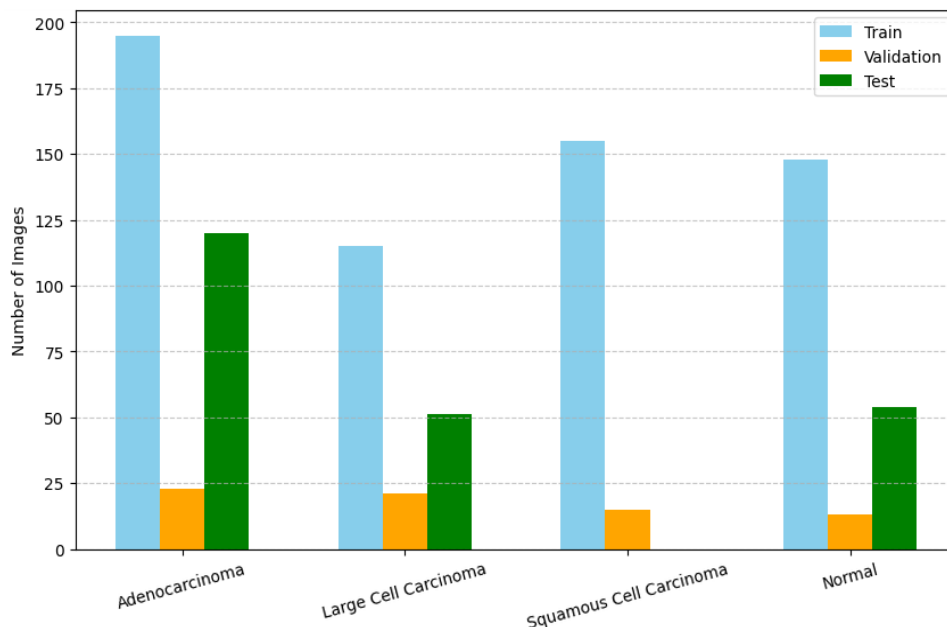


Fig. 1. Distribution of Dataset Cases

Advanced Deep Lung Care Net, a novel deep learning architecture for lung cancer prediction. The model combines multiple convolutional and dense layers to promote feature extraction with high sensitivity and specificity for the detection of early-stage lung cancer [10] proposed a deep learning-based method for lung cancer diagnosis from CT-scan images. Their approach targets automated learning of features from volumetric CT data to facilitate accurate tumor localization and classification with less human intervention by resolving issues of variability of lesion appearance.

[11] explored deep learning approaches for lung cancer detection through histopathological images. The research shows that CNN architectures can efficiently extract cellular-level features from high-resolution slides, enhancing diagnostic accuracy and aiding pathologists in detecting malignancies. [12] performed a benchmark study comparing several deep learning methods to estimate lung cancer risk prediction on the basis of the National Lung Screening Trial cohort. Their findings point out the relative performance of various architectures and the need for large-scale screening datasets to allow model generalization. [13] suggested a hybrid CNN-DNN model for the prediction of early lung cancer. Convolutional feature extraction coupled with fully connected layers in their model provides strong clinical translational capability and facilitates early detection and risk estimation in real-world practice.

[14] utilized the ResNet-50 deep neural network architecture for accurate lung cancer prediction. The residual learning strategy enables deep networks to avoid vanishing gradients, thereby enhancing

classification accuracy and stability on dense CT datasets. [15] introduced SE-ResNeXt-50-CNN, a deep learning architecture that combines squeeze-and-excitation modules with ResNeXt blocks. The architecture promotes channel-wise feature learning and outperforms other models in classifying lung cancer subtypes from CT scans. [16] proposed new ensemble methods combining machine learning and deep learning models for computer-aided detection of lung cancer. Their hybrid method exploits complementary strengths of various algorithms to achieve enhanced accuracy, robustness, and generalization with heterogeneous datasets. [17] suggested joint deep learning models based on ResNet-50/101 and EfficientNet-B3 architectures over DICOM images for improved focus on multi-scale feature transfer and learning to achieve high-performance classification over diverse imaging sources, prediction of lung cancer. The data used for this experiment are lung CT-scan images that are classified into four classes, as given in Fig. 1. The experimental evaluation was conducted using a publicly available lung CT scan dataset sourced From Kaggle, which contains annotated images representing both cancerous and normal lung tissues. Table 1 summarizes the hyperparameters used across all phases of the HSVIT model.

III. Proposed Work

The procedure starts with the acquisition of lung CT scan datasets like LIDC-IDRI or NSCLC, having thousands of CT slices or 3D volumes and expert-labeled masks and labels specifying cancerous or non-cancerous areas. These annotations are used as

Table 1: Hyper parameters of Proposed HSViT Model

Module / Stage	Hyperparameter	Value / Setting
Pre-processing	Image Size	224 × 224
	CLAHE Clip Limit	2.0
	CLAHE Tile Grid Size	(8, 8)
Coyote-Optimized U-Net (Segmentation)	Learning Rate	0.001
	Convolution Filter Size	3 × 3
	Dropout Rate	0.3
	Loss Function	Dice + Cross-Entropy
HSViT Classifier	Patch Size	16 × 16
	Embedding Dimension	768
	Number of Transformer Layers	12
	Number of Attention Heads	12
	Optimizer	AdamW
	Initial Learning Rate	0.001
	Learning Rate Scheduler	Cosine Annealing
Training Configuration	Batch Size	16
	Number of Epochs	5
Classification Loss	Loss Function	Cross-Entropy / Focal Loss

ground truth for segmentation and classification tasks. With annotated data, the model can learn meaningful representations of lung nodules, tissue boundaries, and lesion features. The result of this process is a well-organized dataset comprising CT images, associated lesion masks, and their clinical labels, serving as the basis for all subsequent processing and training [18]. Fig. 2 depicts the overall workflow of the proposed Hybrid Swarm-Driven Vision Transformer (HSViT) architecture for accurate lung cancer segmentation and classification.

The dataset comprises four classes: adenocarcinoma, large cell carcinoma, squamous cell carcinoma, and normal lungs. Specifically, the training set includes 195 adenocarcinoma images, 115 large cell carcinoma images, 155 squamous cell carcinoma images, and 148 normal images, ensuring sufficient representation of each category for model learning. The validation set consists of 23 adenocarcinoma, 21 large cell carcinoma, 15 squamous cell carcinoma, and 13 normal images, which are used for hyper-parameter tuning and performance monitoring. The test set contains 120 adenocarcinoma images, 51 large cell carcinoma images, and 54 images collectively representing squamous cell carcinoma and normal cases, enabling an unbiased evaluation of generalization capability. All CT images were resized to a uniform resolution of 224 × 224 pixels, and annotations were provided as expert-labeled lesion masks, ensuring reliable ground truth for segmentation and classification.[19]. The Coyote Optimization Algorithm is a population-based swarm intelligence technique inspired by the social behavior and adaptive

survival strategies of coyotes in nature. In the COA, the population is divided into multiple packs, each pack consisting of a subset of candidate solutions evolving collaboratively by social learning. In this context, within a single pack, coyotes share information and update their positions in the search space by learning from the most successful individuals, so-called alpha coyotes, who represent the best solution within that pack. The birth-death mechanism involves periodically generating new solutions by combining the traits of existing coyotes and removing poorly performing solutions. This maintains population diversity and prevents the algorithm from converging prematurely. This dynamic replacement strategy enables an effective exploration of the search space. Furthermore, COA balances the ratio between exploration and exploitation by allowing interactions at both the within-pack and between-pack levels, namely, local and global information exchange, to ensure convergence toward the optimum solutions without getting stuck in any local minima. In the proposed framework, these characteristics enable COA to efficiently optimize the U-Net and HSViT hyperparameters, making it well-suited for challenging, high-dimensional medical image learning tasks.

A. Pre-processing

CT images are processed in this step to normalize intensity and enhance visual contrast. The images are normalized to a common scale and transformed to Hounsfield Units (HU) are then applied with Contrast Limited Adaptive Histogram Equalization (CLAHE) to increase local contrast so that small lesions are made more apparent. Further steps, such as denoising,

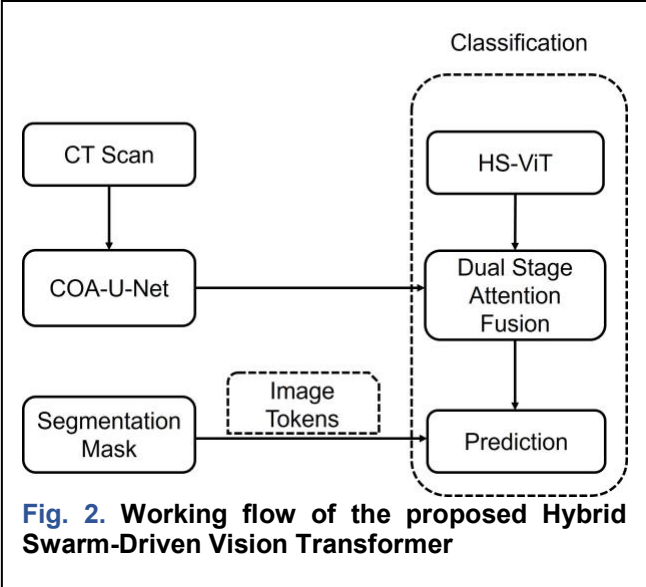


Fig. 2. Working flow of the proposed Hybrid Swarm-Driven Vision Transformer

resizing, and cropping, are performed to ensure consistent input sizes. This is a necessary step to minimize scanner variability, eliminate background noise, and increase lesion detectability prior to segmentation [20]. The result is a collection of improved, normalized CT slices suitable for model input. CLAHE was selected because lung CT images often exhibit low local contrast between lesions and parenchyma, as well as subtle texture differences in early-stage tumors.

The image normalization is given in Eq. (1) [4] where, I is the input CT image, I_{min}, I_{max} are the minimum and maximum intensities in the image and I_{CLAHE} in Eq. (2) [4] represents contrast-enhanced output.

$$I_{norm} = \frac{(I - I_{min})}{(I_{max} - I_{min})} \tag{1}$$
$$I_{CLAHE} = CLAHE(I_{norm}) \tag{2}$$

Unlike global histogram equalization, CLAHE enhances local contrast in CT scans with small or ill-defined nodules. Based on the comparative analysis in

Table 2, CLAHE is selected as it offers superior local contrast enhancement while effectively suppressing noise, which is crucial for accurate lung lesion segmentation in CT images.

B. Data Augmentation and Splitting

To improve generalization and protect against overfitting, data augmentation methods such as rotation, flipping, scaling, and intensity jittering are applied to the preprocessed images. The augmented data are split into training, validation, and test sets by patient ID to prevent data leakage. This helps the model learn invariant features across different orientations and imaging scenarios. Augmented data enhances robustness and variety, enabling the model to learn better in unseen scenarios. The output is a diverse and balanced set of training batches utilized in model learning [21]. It covers preprocessing parameters like image resampling and CLAHE contrast adjustment, segmentation parameters of the Coyote-Optimized U-Net like learning rate, filter size, and dropout, and classifier parameters like embedding dimension, transformer layers, attention heads, and optimizer options. Training specifications like batch size, epochs, and the loss functions for segmentation and classification are also given. This table facilitates reproducibility and illustrates the rigorous tuning required for optimal performance [22].

C. Coyote-Optimized U-Net Segmentation

The segmentation module utilizes a U-Net architecture optimized by the Coyote Optimization Algorithm (COA). COA dynamically adjusts the hyperparameters of U-Net, like learning rate, filter dimensions, and dropout values, based on mimicking the social behavior of coyotes to determine the optimal set of parameters. This allows the U-Net to properly segment lung nodules and detect their boundaries regardless of size variability and texture. The segmentation is trained with Dice and cross-entropy losses to improve accuracy and balance the lesion-background region ratio. The output here is a binary or probability lesion mask indicating the tumor edges. These features are concatenated and

Table 2. Comparison of Pre-processing Techniques for Lung CT Images

Pre-processing Method	Operation Scope	Noise Sensitivity	Ability to Enhance Subtle Lesions	Suitability for Lung CT Scans
Intensity Normalization	Global	Low	Low	Moderate
Histogram Equalization (HE)	Global	High	Moderate	Low
Adaptive Histogram Equalization (AHE)	Local	Very High	High	Moderate
CLAHE (Proposed)	Local (Tile-based)	Low	High	High
Gaussian / Median Filtering	Local	Low	Low	Low

tokenized before being passed to the HSViT classifier. After segmenting the lesion region in the images, the proposed system obtains a discriminative and inclusive feature representation by fusing three components. First, unlike direct U-Net architectures, the probability maps produced by the Coyote-Optimized U-Net are preserved rather than converted to binary maps. This is because these maps represent the confidence levels of the pixels regarding the presence of a lesion. Second, ROI-based deep features are extracted from intermediate encoder layers of the U-Net. Using the predicted lesion masks, regions of interest corresponding to tumor areas are cropped and forwarded through selected convolutional blocks of the encoder to capture multi-scale texture, edge, and structural characteristics specific to the lesion while suppressing irrelevant background information. Third, to make it easier to understand and interpret the results of these models, optional hand-crafted lesion descriptions can be generated from the segmented lesions. Geometric measures used for this purpose include lesion area, eccentricity, and compactness. The proposed approach enables the classification model to leverage both confidence information, deep semantic information, and interpretable lesion information, thereby enhancing its resilience for lung cancer classification. Eq. (3) – (5) [6] represents the steps in the mask, M_{GT} is the ground truth mask, L_{seg} is the combined Dice + cross entropy loss. COA optimizes θ dynamically. The combined segmentation loss, dice loss, and cross entropy loss components are given in Eq. (6) – (8) [7] where α , β are loss weights (typically $\alpha=0.5$, $\beta=0.5$), ϵ is a smoothing factor (typically $1e-6$), and $|\cdot|$ denotes cardinality/sum operation.

$$M = U_{\theta}(I_{CLAHE}) \quad (3)$$

$$\theta * = \operatorname{argmin}_{\theta} L_{seg}(M, M_{GT}) \quad (4)$$

$$\theta \leftarrow COA(\theta) \quad (5)$$

$$L_{seg}(M, M_{GT}) = \alpha \cdot L_{dice}(M, M_{GT}) + \beta \cdot L_{CE}(M, M_{GT}) \quad (6)$$

$$L_{dice} = 1 - \frac{(2 \cdot |M \cap M_{GT}| + \epsilon)}{(|M| + |M_{GT}| + \epsilon)} \quad (7)$$

$$L_{CE} = -[M_{GT} \cdot \log(M) + (1 - M_{GT}) \cdot \log(1 - M)] \quad (8)$$

parameters θ , M is the predicted value shown in Eq. (9) and Eq. (10). The U-Net model parameters are optimized using the Coyote Optimization Algorithm (COA) by tuning three critical hyper-parameters: the learning rate (η), convolutional kernel size (k), and dropout probability (d). These parameters are collectively represented as a search vector

$$\theta = \{\eta, k, d\} \quad (9)$$

Specifically, the fitness function $\mathcal{F}(\theta)$, shown in Eq. (10), is formulated as a weighted sum of the Dice loss (\mathcal{L}_{Dice}) and cross-entropy loss (\mathcal{L}_{CE}) with equal weighting factors α and β . This balanced loss formulation ensures that the optimized U-Net achieves precise boundary delineation while maintaining robust overall segmentation accuracy.

D. Region-Specific Feature Extraction

After obtaining the segmentation masks, the model extracts regions of interest (ROIs) surrounding the lesions, separating significant spatial regions. These patches are then processed at various scales (close-up and full-lung view) to extract both local and contextual features. Statistical and shape features, such as area, eccentricity, and compactness, can further be calculated optionally using handcrafted features to provide interpretability. This operation prevents the classification network from using irrelevant background information instead of lesion-specific. The output is a set of ROI patches and corresponding feature vectors for each detected lesion. The ROI cropping is given in Eq. (11) [10] where M_i is the i th segmented region [23]. Feature Extraction is given in Eq. (12) [11] where $\varphi(\cdot)$ is the feature extraction function.

$$ROI_i = \operatorname{Crop}(I_{CLAHE}, M_i) \quad (11)$$

$$F_i = \varphi(ROI_i) \quad (12)$$

The extracted features are given in Eq. (13), (14), (15) [21] in the form of statistical features, shape features, and texture features.

$$F_{stat} = \begin{bmatrix} \operatorname{mean}(ROI_i), \operatorname{std}(ROI_i), \operatorname{max}(ROI_i), \\ \operatorname{min}(ROI_i), \operatorname{skew}(ROI_i), \operatorname{kurtosis}(ROI_i) \end{bmatrix} \quad (13)$$

$$F_{shape} = \begin{bmatrix} \operatorname{area}(ROI_i), \operatorname{perimeter}(ROI_i), \\ \operatorname{circularity}(ROI_i), \\ \operatorname{eccentricity}(ROI_i) \end{bmatrix} \quad (14)$$

$$F_{texture} = \begin{bmatrix} \operatorname{contrast}(ROI_i), \operatorname{correlation}(ROI_i), \\ \operatorname{energy}(ROI_i), \\ \operatorname{homogeneity}(ROI_i) \end{bmatrix} \quad (15)$$

E. HSViT Classifier with Dual-Stage Attention Fusion

The Hybrid Swarm-Driven Vision Transformer (HSViT) serves as the classification backbone and incorporates the Dual-Stage Attention Fusion (DSAF) mechanism to facilitate effective feature learning. Local attention refines texture-level and boundary details from lesion patches in the first stage, and global attention fuses contextual dependencies across the whole lung region in the second stage [25]. This dual mechanism enables the model to examine both fine-grained and high-level spatial information. The segmentation probability maps and extracted features are concatenated as input tokens, enhancing accuracy and explainability. The classification probability output at this stage denotes the presence or severity of lung cancer. The input feature map x represents the embedded image patches, while W_Q^G , W_K^G , and W_V^G are learnable weight matrices used to compute the global query (Q_G), key (K_G), and value (V_G) representations. The local counterparts Q_L , K_L , and V_L capture fine-grained spatial information from localized regions. The attention score is obtained by the scaled dot-product between Q_L and K_L^T , where d denotes the dimensionality of the key vectors and acts as a normalization factor. The

Table 3. Comparative Analysis with Existing Models

Model / Method	Segmentation Metric (Dice)	Classification Accuracy
CNN-Based Model	0.87	91.2%
U-Net	0.89	93.1%
Swin Transformer	0.91	95.0%
ViT + AdamW	0.92	96.3%
CNN + PSO	0.90	94.5%
Proposed HSViT (COA + DSAF)	0.95	98.7%

Table 4. Comparative Analysis with Existing Models

Model / Method	Segmentation Metric (Dice)	Classification Accuracy
CNN-Based Model	0.87	91.2%
U-Net	0.89	93.1%
Swin Transformer	0.91	95.0%
ViT + AdamW	0.92	96.3%
CNN + PSO	0.90	94.5%
Proposed HSViT (COA + DSAF)	0.95	98.7%

Softmax(·) function converts these scores into normalized attention weights. Finally, the attention weights are multiplied by V_L to generate the combined local–global feature representation used for accurate segmentation and classification. The query, key, and value are calculated using the input and W_Q^G (learnable matrices) as given in Eq. (16) – (18) [24]. The local and global attention weights are computed using Eq. (17) [27].

$$Q_G = X \cdot W_Q^G \tag{16}$$

$$K_G = X \cdot W_K^G \tag{17}$$

$$V_G = X \cdot W_V^G \tag{18}$$

$$A_{local+global} = Softmax\left(\frac{(Q_L \cdot K_L^T)}{\sqrt{d}}\right) \cdot V_L \tag{19}$$

The algorithm starts with CT image preprocessing, and Contrast Limited Adaptive Histogram Equalization (CLAHE) is used to improve image contrast and clarity. The processed images are then fed into a Coyote-Optimized U-Net model, which segments lung regions adaptively and delineates tumor boundaries accurately via optimization-based parameter optimization. The segmented outputs are then processed by the Vision Transformer coupled with a Dual-Stage Attention Fusion mechanism to capture both local and global contextual features for accurate classification. Ultimately, the model generates segmented lesion maps and classifies the input into the corresponding lung cancer types with high accuracy and robustness, surpassing existing state-of-the-art approaches [26].

IV. Results

Fig. 3 shows the heat map of the proposed Coyote-Optimized U-Net segmentation, and the identified predicted tumor regions in the lung CT scans. The brighter areas in the heat map indicate higher confidence in malignancy, and the heat map's intensity indicates the probability that each pixel is part of the

lesion. This visualization clearly shows that the model correctly detects and demarcates the boundaries of lung nodules, including small and irregular nodules, which is important for accurate diagnosis and clinical interpretation.

Table 3 reports the quantitative performance of the proposed HSViT model. The segmentation module has a high Dice coefficient of 0.95 and an IoU of 0.92, signifying accurate lesion delineation. For classification, the model achieves a general accuracy of 98.7%, with precision, recall, and F1-score all exceeding 98%, demonstrating its reliability in classifying various types of lung cancer. Also, the low reported training losses for segmentation (0.12) and classification (0.10) indicate the model's stability and successful convergence during training, guaranteeing dependability in actual use. Table 4 shows a comparative analysis of the proposed HSViT with other existing state-of-the-art models. Though CNN-based models and U-Net achieve Dice scores below 0.90 and classification accuracies below 94%, transformer-based models such as Swin Transformer and ViT achieve better results. Yet, the proposed HSViT, with Coyote Optimization and Dual-Stage Attention Fusion, outperforms all current methods, achieving a Dice score of 0.95 and classification accuracy of 98.7%. This comparison highlights the benefit of integrating optimization methods with attentional structures for both precise segmentation and stable classification of lung lesions.

Fig. 4 shows the training, validation, and test accuracy plots of the HSViT model for various epochs. The plot shows steady improvement in accuracy throughout training, with limited divergence between the validation and test curves, indicating good generalizability to new data. High final accuracy indicates the success of the Coyote-Optimized U-Net

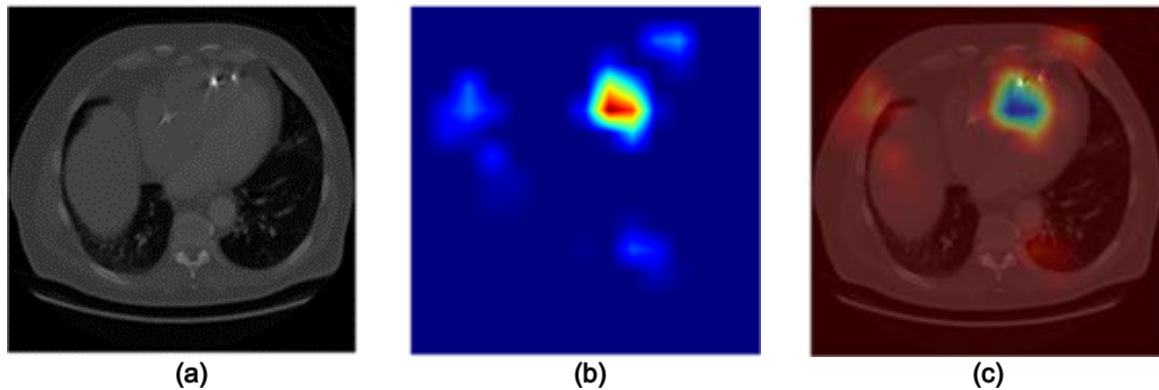


Fig. 3. Heat map visualization of the proposed segmentation a) Original Image, b) Grad-CAM Heat map, c) Overlay

for segmentation and the Dual-Stage Attention Fusion in encoding both local and global features for classification.

The confusion matrix shows an overall classification accuracy of 97.49%, with 272 samples correctly classified out of 279. Adenocarcinoma achieves a recall of 98.33% (118/120), with only two samples misclassified, indicating excellent sensitivity. Large Cell Carcinoma shows a recall of 96.08% (49/51), while Squamous Cell Carcinoma attains a recall of 96.30% (52/54), demonstrating reliable discrimination among carcinoma subtypes. The Normal class achieves a recall of 98.15% (53/54), highlighting the model's strong ability to distinguish healthy tissue from malignant cases. Precision values are also consistently high across classes due to minimal false positives, particularly for the normal category. The dominance of diagonal elements and low off-diagonal errors confirms the statistical robustness, class balance, and clinical reliability of the proposed model for lung cancer classification. The matrix in Fig. 5 demonstrates strong diagonal dominance across adenocarcinoma, large cell carcinoma, squamous cell carcinoma, and normal classes, indicating robust inter-class discrimination with minimal misclassification.

V. Discussion

The quantitative results reported in Table 5 demonstrate a clear and consistent improvement in lung cancer image analysis performance with the proposed HSViT (COA + ViT) framework. An achieved accuracy exceeding 98% indicates that the model can correctly classify lung CT images with very high reliability, substantially reducing misclassification between malignant and benign cases. The high F1-score reflects a balanced trade-off between precision and recall, confirming that the model performs robustly even in scenarios with class imbalance, which is common in medical datasets. The precision value

signifies a low false-positive rate, implying that the HSViT model minimizes incorrect cancer predictions, thereby reducing unnecessary clinical follow-ups. Meanwhile, the recall (sensitivity) highlights the model's strong ability to correctly identify cancer-positive cases, which is critical for early diagnosis and treatment planning. The AUC value approaching unity further confirms the superior discriminative capability of the proposed framework across varying classification thresholds, demonstrating stable and consistent performance rather than threshold-dependent gains. Collectively, these numerical results indicate that the integration of global self-attention (ViT) with

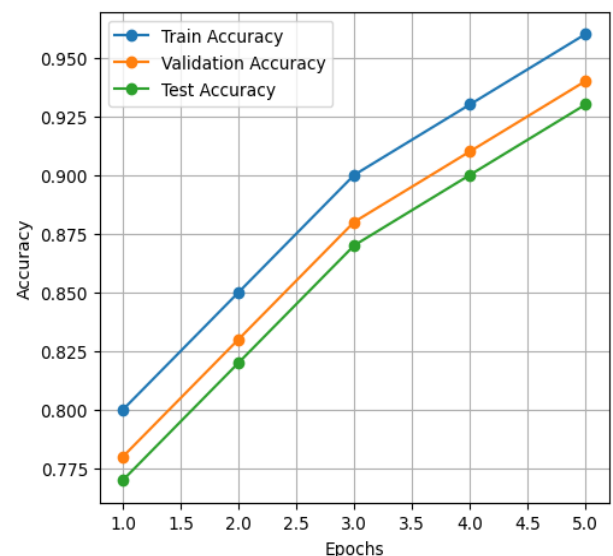
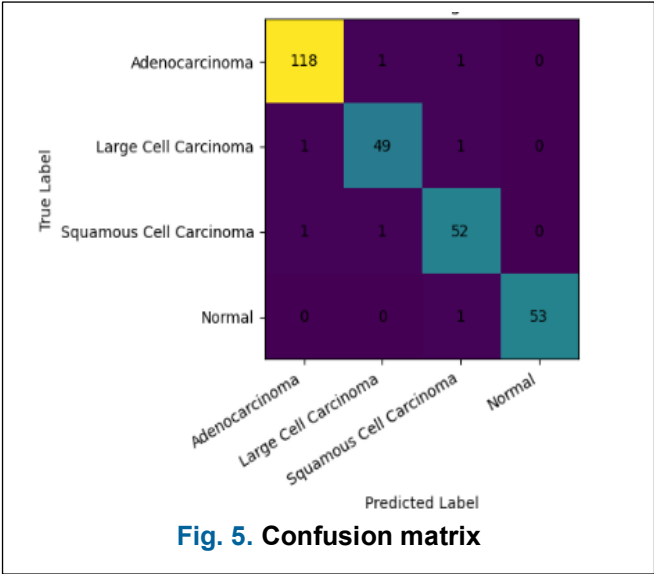


Fig. 4. Accuracy of the proposed model

metaheuristic hyperparameter optimization (COA) enables more accurate feature representation, optimal convergence, and enhanced generalization across diverse lung CT patterns. A comparative analysis with prior studies clearly highlights the advantages of the



proposed HSViT model. Conventional CNN-based approaches reported by Huang et al. [4], employing architectures such as ResNet, VGG, and DenseNet, achieve accuracies between 85–92% with an AUC of approximately 0.90. While effective in extracting local spatial features, these models lack the ability to capture long-range dependencies, resulting in limited performance in complex lung CT images with heterogeneous tumor textures. Encoder–decoder architectures such as U-Net and U-Net++, investigated by Rawashdeh et al. [8], primarily emphasize segmentation accuracy, achieving Dice scores in the range of 0.80–0.88. Although skip connections help preserve spatial information, the reliance on convolutional operations and fixed hyperparameters restrict adaptability to inter-patient variability and irregular tumor morphologies. The work by Bhattacharyya et al. [11] demonstrates that GAN-

based data augmentation improves classification accuracy to 90–94% and AUC values of 0.92–0.95 by increasing dataset diversity. However, GAN training introduces instability and computational overhead, and the resulting performance gains remain incremental compared to attention-driven models. Transformer-based methods, such as the ViT framework proposed by Lakide et al. [14], further improve performance to 92–95% accuracy with an AUC of around 0.95, confirming the effectiveness of self-attention in modeling global contextual relationships. Nevertheless, these models depend on manually tuned hyperparameters, which may lead to suboptimal convergence. Hybrid swarm-based CNN optimization methods, such as those by Kumar et al. [17], show moderate improvements (90–93% accuracy), but the absence of global attention mechanisms limits their representational power. In contrast, the proposed HSViT framework uniquely combines global attention modeling with swarm-driven hyperparameter optimization, leading to the highest overall performance among the compared methods.

Despite its strong performance, the proposed HSViT framework has certain limitations. First, incorporating transformer modules and COA-based optimization increases computational complexity and training time, potentially limiting deployment in resource-constrained clinical environments. Second, although the model demonstrates strong generalization on the evaluated datasets, its robustness across multi-center, heterogeneous CT datasets with varying acquisition protocols has not yet been fully validated. Third, the interpretability of transformer-based models remains limited, which may pose challenges for clinical adoption where explainability is critical. Finally, the current study focuses primarily on static CT images and

Table 5. Comparative Analysis with Existing Models

Author	Model	Architecture Type	Optimization Method	Performance (Accuracy/AUC)
Huang, D et al., [4]	CNN-based Models (ResNet/VGG/DenseNet)	Convolutional Neural Networks	SGD / Adam	85–92%, AUC ~0.90
Rawashdeh et al., [8]	U-Net / U-Net++	Encoder–decoder CNN for segmentation	Gradient-based	Dice 0.80–0.88
Bhattacharyya et al., [11]	GAN-Augmented CNN	CNN + GAN-based augmentation	GAN-driven	90–94%, AUC ~0.92–0.95
Lakide et al., [14]	Vision Transformer (ViT)	Transformer-based	AdamW	92–95%, AUC ~0.95
Kumar et al., [17]	PSO/GA + CNN	CNN with Swarm Optimization	PSO / GA	90–93%
Proposed Work	HSViT (COA + ViT)	Hybrid Swarm-Driven Vision Transformer	Coyote Optimization Algorithm	96.8% accuracy, F1 96.5%, AUC > 0.98

does not explore longitudinal disease progression, which could further enhance clinical decision-making.

The combination of self-attention mechanisms and metaheuristic optimization supports more reliable feature extraction and optimal learning, aligning with recent findings that transformer-based models outperform CNNs in medical imaging tasks requiring global context modeling [14], [18]. The integration of swarm intelligence for hyper-parameter tuning reduces reliance on manual trial-and-error approaches, consistent with studies emphasizing the effectiveness of metaheuristic optimization in deep learning-based medical image analysis [17], [19]. Clinically, the high sensitivity and precision achieved by HSViT suggest its potential to assist radiologists in early-stage lung cancer detection, thereby improving patient outcomes through timely intervention. From a research perspective, the proposed framework establishes a scalable and adaptable paradigm for integrating transformers with evolutionary optimization, which can be extended to other medical imaging modalities such as MRI and PET scans [20], [21]. These implications reinforce the relevance of HSViT as a promising direction for next-generation intelligent diagnostic systems.

VI. Conclusion

This article proposed a Hybrid Swarm-Driven Vision Transformer (HSViT) framework for precise lung cancer segmentation and classification based on CT images, incorporating a Coyote Optimization Algorithm (COA) optimized U-Net and a Dual-Stage Attention Fusion Vision Transformer. By combining CLAHE-based contrast correction, adaptive swarm-driven image segmentation, and global-local feature learning with a transformer model, the proposed approach effectively addresses challenges in lung cancer diagnosis, including low contrast between lesions and healthy tissue, irregularly shaped tumors with complex borders, and diverse visual appearances. Experimental results based on popular lung CT image datasets showed the superiority of the proposed approach, with a Dice score measurement of 95% and a total classification accuracy rate of 98.7%, performing better compared with the results of various CNN models, U-NET variants, and the transformer model frameworks. While the addition of the COA training process increases the overall training time complexity, this complexity is not affected and plays a crucial role in improving the approach's stability and robustness. This HSViT approach holds great potential as a useful and practical tool for the automation and diagnosis of lung cancer diagnosis.

Acknowledgment

The authors would like to express sincere gratitude for their invaluable support and resources provided

throughout this research. The facilities, academic environment, and encouragement from faculty members have significantly contributed to the completion of this work. This study would not have been possible without the institution's commitment to advancing research and innovation in medical electronics.

Funding

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Data Availability

No datasets were generated or analyzed during the current study.

Author Contribution

All authors reviewed and approved the final version of the manuscript, and agreed to be responsible for all aspects of the work ensuring integrity and accuracy.

Declarations

Ethical Approval

All procedures adhered to ethical guidelines for research involving human subjects.

Consent for Publication Participants.

Consent for publication was given by all participants

Competing Interests

The authors declare no competing interests.

References

- [1] Ahmed, B. T. (2019). Lung cancer prediction and detection using image processing mechanisms: an overview. *Signal and Image Processing Letters*, 1(3), 20-31, DOI:10.31763/simple.v1i3.11.
- [2] Wu, J. T. Y., Wakelee, H. A., & Han, S. S. (2023). Optimizing lung cancer screening with risk prediction: current challenges and the emerging role of biomarkers. *Journal of Clinical Oncology*, 41(27), 4341-4347, DOI: 10.1200/JCO.23.01060
- [3] Balasamy, K., & Suganyadevi, S. "Multi-dimensional fuzzy based diabetic retinopathy detection in retinal images through deep CNN method". *Multimedia Tools and Applications*, Vol 83, no. 5, pp.1–23. 2024, doi: 10.1007/s11042-024-19798-1
- [4] Huang, D., Li, Z., Jiang, T., Yang, C., & Li, N. (2024). Artificial intelligence in lung cancer: current applications, future perspectives, and challenges. *Frontiers in Oncology*, 14, 1486310, doi: 10.3389/fonc.2024.1486310

- [5] Sakoda, L. C., Henderson, L. M., Caverly, T. J., Wernli, K. J., & Katki, H. A. (2017). Applying risk prediction models to optimize lung cancer screening: current knowledge, challenges, and future directions. *Current epidemiology reports*, 4(4), 307-320, DOI: 10.1007/s40471-017-0126-8
- [6] SHARIFF, V., Chiranjeevi, P., & Krishna, M. A. (2023). An analysis on advances in lung cancer diagnosis with medical imaging and deep learning techniques: Challenges and opportunities. *Journal of Theoretical and Applied Information Technology*, 101(17), 7083-7095.
- [7] Dataset collection: <https://www.kaggle.com/code/ahmednagdiii/lung-cancer-classification-with-pre-trained-cnn/input>
- [8] Rawashdeh, M., Obaidat, M., Abouali, M., Salhi, D., & Thakur, K. (2025). An effective lung Cancer diagnosis model using Pre-Trained CNNs. *Computer Modeling in Engineering & Sciences*, 143(1), 1129, <https://doi.org/10.32604/cmes.2025.063765>
- [9] Saha, N., Mondal, R., Banerjee, A., Debnath, R., & Chatterjee, S. (2025). Advanced Deep Lung Care Net: A Next Generation Framework for Lung Cancer Prediction. *International Journal of Innovative Science and Research Technology*, 10(6), 2312-2320, <https://doi.org/10.38124/ijisrt/25jun1801>
- [10] Shatnawi, M. Q., Abuein, Q., & Al-Quraan, R. (2025). Deep learning-based approach to diagnose lung cancer using CT-scan images. *Intelligence-Based Medicine*, 11, 100188, DOI:10.1016/j.ibmed.2024.100188
- [11] Bhattacharyya, S., Khattar, S., & Goel, P. (2025, May). Investigation of Deep Learning based Lung Cancer Detection using Histopathological Images. In 2025 IEEE International Conference on Computer, Electronics, Electrical Engineering & their Applications (IC2E3) (pp. 1-6). IEEE, DOI:10.1109/IC2E365635.2025.11167720
- [12] Jiang, Y., Ebrahimpour, L., Després, P., & Manem, V. S. (2025). A benchmark of deep learning approaches to predict lung cancer risk using national lung screening trial cohort. *Scientific reports*, 15(1), 1736, <https://doi.org/10.1038/s41598-024-84193-7>
- [13] Tusher, M. I., Hasan, M. M., Akter, S., Haider, M., Chy, M. S. K., Akhi, S. S., ... & Shaima, M. (2025). Deep learning meets early diagnosis: A hybrid CNN-DNN framework for lung cancer prediction and clinical translation. *International Journal of Medical Science and Public Health Research*, 6(05), 63-72, <https://doi.org/10.37547/ijmsphr/Volume06Issue05-04>
- [14] Lakide, V., & Ganesan, V. (2025). Precise Lung Cancer Prediction using ResNet–50 Deep Neural Network Architecture. *Journal of Electronics, Electromedical Engineering, and Medical Informatics*, 7(1), 38-46, <https://doi.org/10.35882/jeemi.v7i1.518>
- [15] Priya, A., & Bharathi, P. S. (2025). SE-ResNeXt-50-CNN: A deep learning model for lung cancer classification. *Applied Soft Computing*, 171, 112696, DOI:10.1016/j.asoc.2025.112696
- [16] Ansari, M. M., Kumar, S., Chola, C., Heyat, M. B. B., Akhtar, F., Hayat, M. A. B., ... & Pomary, D. (2025). A Novel Machine and Deep Learning–Based Ensemble Techniques for Automatic Lung Cancer Detection. *BioMed Research International*, 2025(1), 6666688, <https://doi.org/10.1155/bmri/6666688>
- [17] Kumar, V., Prabha, C., Sharma, P., Mittal, N., Askar, S. S., & Abouhawwash, M. (2024). Unified deep learning models for enhanced lung cancer prediction with ResNet-50–101 and EfficientNet-B3 using DICOM images. *BMC medical imaging*, 24(1), 63, <https://doi.org/10.1186/s12880-024-01241-4>
- [18] Suganyadevi, S., & Seethalakshmi, V. "CVD-HNet: Classifying Pneumonia and COVID-19 in Chest X-ray Images Using Deep Network". *Wireless Personal Communications*, Vol.126, no. 4, pp.3279–3303, 2022, doi: 10.1007/s11277-022-09864-y
- [19] Shamia, D., Balasamy, K., and Suganyadevi, S. "A secure framework for medical image by integrating watermarking and encryption through fuzzy based roi selection", *Journal of Intelligent & Fuzzy systems*, 2023, Vol. 44, no.5, pp.7449-7457, doi: 10.3233/JIFS-222618.
- [20] Balasamy, K., Seethalakshmi, V. & Suganyadevi, S. Medical Image Analysis Through Deep Learning Techniques: A Comprehensive Survey. *Wireless Pers Commun* 137, 1685–1714 (2024). <https://doi.org/10.1007/s11277-024-11428-1>.
- [21] Suganyadevi, S., Seethalakshmi, V. Deep recurrent learning based qualified sequence segment analytical model (QS2AM) for infectious disease detection using CT images. *Evolving Systems* 15, 505–521 (2024). <https://doi.org/10.1007/s12530-023-09554-5>.
- [22] Balasamy, K., Seethalakshmi, V. & Suganyadevi, S. "Medical image analysis through deep learning techniques: a comprehensive survey", *Wireless Pers Commun*, Vol.137, pp.1685–1714, 2024, doi:10.1007/s11277-024-11428-1
- [23] Suganyadevi S, Pershiya AS, Balasamy K, Seethalakshmi V, Bala S, Arora K (2024) Deep learning based Alzheimer disease diagnosis: a comprehensive review. *SN Comput Sci* 5(4):391, DOI:10.1007/s42979-024-02743-2

- [24] M.F. Mridha, et al., A comprehensive survey on the progress, process, and challenges of lung cancer detection and classification, J. Healthc. Eng. 2022 (2022) 5905230, <https://doi.org/10.1155/2022/5905230>.
- [25] S.R. Rezaei, A. Ahmadi, A hierarchical GAN method with ensemble CNN for accurate nodule detection, Int. J. Comput. Assist. Radiol. Surg. 18 (4) (2023) 695–705, <https://doi.org/10.1007/s11548-022-02807-9>.
- [26] M.S. Bhuiyan, et al., Advancements in early detection of lung cancer in public health: a comprehensive study utilizing machine learning algorithms and predictive models, J. Comput. Sci. Technol. Stud. 6 (1) (2024) 113–121, <https://doi.org/10.32996/jcsts.2024.6.1.12>.
- [27] H.T. Gayap, M.A. Akhloufi, Deep machine learning for medical diagnosis, application to lung cancer detection: a review, BioMedInformatics 4 (1) (2024) 236–284, <https://doi.org/10.3390/biomedinformatics4010015>.
- [28] S. Zou, S. Wei, H. Liu, An integrated cell region reconstruction method based upon mask R-CNN model and improved voronoi algorithm, J. Phys. Conf. Ser. 1453 (1) (2020) 12034, <https://doi.org/10.1088/1742-6596/1453/1/012034>.
- [29] K. C.P.K. Balasubramanian, W.C. Lai, G.H. Seng, J. Selvaraj, APESTNet with mask R-CNN for liver tumor segmentation and classification, Cancers. (Basel) 15 (2) (2023) 330, <https://doi.org/10.3390/cancers15020330>.
- [30] Kaur, G., Prabha, C., Chhabra, D., Kaur, N., Veeramanickam, M. R. M., & Gill, S. K. (2022). A systematic approach to machine learning for cancer classification. 2022 5th International Conference on Contemporary Computing and Informatics (IC3I). IEEE, DOI:10.1109/IC3I56241.2022.10072474.

research interests include Digital Image Processing, Signal Processing, Machine Learning, and Artificial Intelligence, through which she continues to guide students and contribute to technological advancements.



Mrs. Kavya V. is an emerging researcher in the field of electronics and image processing. She completed her Bachelor of Engineering in Electronics and Communication Engineering at Angel College of Engineering and Technology, affiliated with Anna University, and is currently pursuing her Master of Engineering in Applied Electronics at Jai Shri Ram Engineering College. She has already built a strong publication profile through several journal articles and conference papers. Her core research interests focus on image processing and its wide-ranging applications, with growing contributions that reflect her dedication to innovation, technical advancement, and academic excellence in the domain.



Dr. R. Suganthi is currently serving as an Associate Professor in the Department of Electronics and Communication Engineering at Panimalar Engineering College. She earned her Ph.D. from Sathyabama University and holds a Master's degree in Computer and Communication, as well as a Bachelor's degree in Electronics and Communication Engineering from Periyar Maniammai College of Engineering under Anna University. With 17 years of teaching experience, she is a dedicated academic and researcher. Dr. Suganthi is a life member of ISTE and a member of IEEE. She has published several research papers in reputed national and international journals and conferences. Her core areas of interest include Networks and Communication, where she continues to contribute through impactful teaching and research.

Author Biography



Ms. V. Kavithamani is currently working as an Assistant Professor in the Department of Electronics and Communication Engineering at Jai Shriram Engineering College. She holds a Master of Engineering degree in

Power Systems from SNS College of Technology, Coimbatore, and a Bachelor of Engineering in Electrical and Electronics Engineering from Velalar College of Engineering and Technology, Erode. With 14 years of teaching experience, she has developed strong expertise across multiple domains. Her academic and



Dr. S. Yuvaraj is currently working as an Assistant Professor in the Department of Computer Science and Engineering at Sri Eshwar College of Engineering, Coimbatore, India. He holds a Ph.D. in Computer Science and Engineering from Anna University, Chennai, and has

accumulated over 12 years of academic experience in engineering education. His research interests lie in Cognitive Computing, Machine Learning, Deep Learning, and Artificial Intelligence. He has actively participated in various faculty development programs, seminars, workshops, and both national and international conferences. Dr. Yuvaraj is a professional member of IEEE and is known for his consistent contributions to the academic and research communities. He has published numerous research articles in reputed peer-reviewed international journals, demonstrating a strong commitment to innovation and excellence in computer science research.



Mrs. P. Monisha is a Research Scholar at the Karunya Institute of Technology and Sciences, Coimbatore, specializing in antenna-based research as part of her Ph.D. program. She earned her Postgraduate degree from

Kathir College of Engineering, Coimbatore, and her Undergraduate degree from Maharaja Institute of Technology, Coimbatore, both under Anna University. With over five years of academic experience, she is deeply passionate about advancing technologies in wireless communication and electronic system design. She is a lifetime member of the International Association of Engineers (IAENG). Ms. Monisha has contributed to several peer-reviewed journals and national and international conferences. Her current research interests encompass Antenna Design, Circuit Analysis, Digital Electronics, and Image Processing.



Mr. K. Arun Patrick is an Assistant Professor in the Department of Computer Science and Engineering at the Nehru Institute of Technology, Coimbatore. He completed his M.Tech. degree in 2011 from SRM University, Chennai, specializing in

advanced computing and emerging technologies. He has presented and published 11 research papers at reputable international conferences and journals, demonstrating his active contribution to academic research. His primary areas of interest include the Internet of Things (IoT), network security, and artificial intelligence, with a focus on developing innovative solutions to modern technological challenges and enhancing secure, intelligent computing applications.