

Hybrid Attention Mechanisms in TransUNet: Enhancing Medical Image Segmentation with CNN and Transformer Architectures

PROJECT REPORT - Team B5

Submitted by

D.Mahadev Naidu (CB.AI.U4AID23111)

S.Raghu Rama Sandeep (CB.AI.U4AID23140)

G.Siva Sai Kumar (CB.AI.U4AID23156)

K. Manoj Kumar (CB.AI.U4AID23161)

in partial fulfillment for the award of the degree of

Bachelor of Technology

in

Artificial Intelligence & Data Science



BONAFIDE CERTIFICATE



Amrita Vishwa Vidyapeetham

Coimbatore - 641112

This is to certify that the Project report entitled “Hand Gesture Recognition” submitted by Team-B5 for the award of the degree of Bachelor of Technology in “Artificial Intelligence & Data Science” is a bonafide record of the original work done by them under our guidance and supervision at Amrita School of Artificial Intelligence, Coimbatore. To the best of my knowledge, this work has not formed the basis for the award of any degree/diploma/associateship.

Internal Examiner

Dr. Mithun Kumar Kar
Assistant Professor
Amrita School of Artificial Intel-
ligence
Amrita Vishwa Vidyapeetham,
Coimbatore

External Examiner

Dr. Subeesh T
Assistant Professor
Amrita School of Artificial Intel-
ligence
Amrita Vishwa Vidyapeetham,
Coimbatore

Submitted for the university ex-
amination held on

DECLARATION



Amrita School of Artificial Intelligence

Amrita Vishwa Vidyapeetham

Coimbatore - 641112

We, Team-B5, hereby declare that this project report entitled “Hand Gesture Recognition” is the record of the original work done by us under the guidance of Dr. Mithun Kumar Kar, Assistant Professor and Dr. Subeesh T, Assistant Professor, Amrita School of Artificial Intelligence, Coimbatore. To the best of my knowledge, this work has not formed the basis for the award of any degree/diploma/associateship.

Student Signature

Name: D.Mahadev Naidu
Roll No: CB.AI.U4AID23111

Name: G.Siva Sai Kumar
Roll No: CB.AI.U4AID23156

Student Signature

Name:S.Raghu Rama Sandeep
Roll No: CB.AI.U4AID23140

Name: K. Manoj Kumar
Roll No: CB.AI.U4AID23161

Contents

1	abstract	4
2	Introduction	4
3	Literature Review	5
4	Methodology	7
4.1	Data Collection	7
4.2	Data Preprocessing	7
4.3	CNN Encoder Deployment with ResNet CNN	7
4.4	Dual Attention Block	8
4.5	Feature Refinement Module	8
4.6	Stochastic Depth	9
4.7	Multiscale Skip Connection	9
4.8	Boundary-Aware Loss Function	9
5	Results	10
6	Conclusion	11
7	References	14

1 abstract

Lung cancer remains among the leading causes of cancer death worldwide, and therefore early and accurate diagnosis must be conducted in order to enhance patient outcomes. Segmentation of lung tumors from computed tomography (CT) scans can be aided by medical image segmentation to achieve accurate localization and treatment planning. The current work presents segmentation of lung cancer with deep learning-based methods employing state-of-the-art architectures including U-Net variants, attention-based models, and transformer-based models. In the proposed method, the incorporation of multi-scale feature fusion and self-attention mechanisms enhances segmentation accuracy. The proposed method registers state-of-the-art performance on experiments conducted using public datasets including LIDC-IDRI and NSCLC-Radiomics with improved Dice coefficient and Intersection over Union (IoU) metrics compared to conventional methods. The results thereby illustrate the capability of deep learning to enhance detection of lung cancer and facilitate early diagnosis and treatment planning by radiologists.

2 Introduction

Lung cancer segmentation of computed tomography (CT) images is of critical significance in the early diagnosis and treatment planning. Proper delineation of the tumor boundary enables clinicians to properly assess tumor characteristics, monitor disease progression, and develop effective treatment plans. Traditional segmentation has been performed through manual annotation by radiologists, which is labor-intensive, time-consuming, and associated with inter- and intra-observer variability.

The conventional lung tumor segmentation techniques are threshold-based techniques, region growing techniques, edge detection techniques, and atlas-based techniques. Threshold-based techniques segment images by dividing pixels according to intensity values but are poor in handling intensity inhomogeneity in tumors. Region growing techniques start from seed points and add adjacent pixels iteratively by similarity measures but are affected by boundary leakage in areas of poor edges. Edge detection techniques detect boundaries from intensity gradients but are noisy and are prone to creating disconnected boundaries. Atlas-based techniques register a labeled atlas to the target image but are affected by anatomical variation and pathological change.

These conventional techniques share generic limitations: they generally need extensive parameter tuning, are prone to the high morphological variability of tumors, and cannot handle the low contrast between tumor and skin. They generally also need extensive human intervention and hence are unsuitable for high-volume clinical applications. The non-uniform appearance, varied sizes, and irregular shapes of lung tumors pose further challenges that are hard for conventional algorithms to handle.

Deep learning models have been potent tools for the automation of this task owing to their high accuracy and stability. These methods have the ability to learn hierarchical representations of features from data, which allows them to generalize complicated patterns

and contextual information. Convolutional Neural Networks (CNNs) have proven to be highly successful in the medical image segmentation task, with U-Net and its extensions been demonstrated to be of particular promise in maintaining localization precision while encoding spatial information [1]. There remain challenges in the guise of class imbalance between tumor and normal tissue, shape and texture variation of tumors, and limited availability of annotated data.

Follow-up research has employed multi-resolution techniques and attention mechanisms to enhance the segmentation performance [2]. Multi-resolution techniques handle images at various resolutions, preserving both local detail and global context. Attention mechanisms enable models to selectively focus on relevant features and filter out irrelevant features, thereby enhancing segmentation accuracy in challenging areas.

This paper introduces a new lung cancer segmentation model that combines attention and adversarial learning with deep learning models. Our method overcomes the shortcomings of existing methods by: (1) automatically learning discriminative features that can learn to adapt to appearance change in the tumor; (2) using dynamic spatial attention modules to weigh informative areas dynamically; (3) using adversarial training to enhance boundary definition; and (4) using a multi-scale architecture that learns to capture fine details and context. Mass clinical dataset trials confirm the effectiveness of the proposed method in achieving state-of-the-art segmentation accuracy. Our approach drastically minimizes manual intervention while retaining high accuracy, especially in difficult cases with irregular tumor morphology or poor contrast. The proposed framework is a significant improvement over conventional segmentation methods and previous deep learning models, providing a strong solution for automated segmentation of lung cancer that can potentially speed up clinical workflows and enhance treatment planning.

3 Literature Review

Lung cancer segmentation has been a prominent research area since segmentation can facilitate early diagnosis and planning of treatment. Deep learning methods have been recently introduced to enhance the performance, efficiency, and generalizability of segmentation.

Zhao et al. [1] introduced a better U-Net-based model, SMR-UNet, which combines self-attention, multi-scale feature extraction, and residual architectures for lung nodule segmentation. The model substitutes traditional convolutional units with residual blocks for feature propagation and convergence promotion. It enhances segmentation performance by introducing a Transformer-based global modeling unit and a multi-scale feature fusion module. Experimental performance on the LIDC dataset yielded a Dice coefficient of 0.9187, representing a gain of 1.33

Fan et al. [8] proposed DMC-UNet, a depthwise separable convolution (DSC) driven, multi-scale feature fusion, coordinate attention mechanism (CCA) driven network for lung nodule segmentation. The authors sought low computational complexity at the expense of high segmentation accuracy. Through substituting normal convolutions with DSC, they obtained a lightweight model with efficient feature extraction. Their method

obtained an IoU of 65.52

Kamal et al. [5] provided an overview of some image analysis methods of the thorax at the MICCAI 2020 Thoracic Image Analysis (TIA) workshop. The workshop emphasized the state of the art in deep learning innovation for lung nodule segmentation, e.g., transformer-based and multi-scale architectures. Their results indicated that the application of multiple imaging modalities, such as PET and CT, could further enhance segmentation accuracy. The paper highlighted the importance of having robust datasets for generalizable deep learning models.

Wang et al. [7] proposed TransUnet, a hybrid network of deep learning using Transformers and U-Net for classifying lung nodules. The network utilizes self-attention for long-range dependency learning and a U-Net encoder-decoder architecture for accurate localization. TransUnet obtained accuracy of 84.62

Li et al. [6] conducted the ACDC@LungHP Challenge for assessing deep learning-based models for segmenting lung cancer from whole-slide histopathology images. Multimodel and CNN-based approaches were utilized for submissions. The top-ranked approach had a Dice coefficient of 0.8372, close to that of the inter-observer consensus of 0.8398, illustrating the potential of AI-supported pathology for lung cancer diagnosis.

Said et al. [4] introduced UNETR, a Transformer U-Net architecture for 3D lung tumor segmentation. Leveraging the self-attention ability of Transformers, UNETR can preserve long-distance dependencies while preserving spatial accuracy. UNETR, which was trained on the Medical Segmentation Decathlon dataset, scored a segmentation accuracy of 97.83

Le et al. [3] introduced RRc-UNet, a residual recurrent convolutional U-Net for the segmentation of 3D lungs. Adding recurrent blocks to the encoder and decoder enhances the preservation of features at different resolutions. The method was tested on the NSCLC-Radiomics dataset and demonstrated better segmentation performance, especially in the case of tumors with irregular boundaries.

Chen et al. [2] proposed MAU-Net, a multi-attention U-Net, which improves lung nodule segmentation with double attention mechanisms. Spatial and channel-wise attention is utilized by the model to provide improved segmentation masks, allowing feature representation. Experiments on clinical datasets demonstrated that MAU-Net performs better than conventional U-Net models, especially on difficult cases with nodule size and shape heterogeneity.

Zhao et al. [1] presented DSU-Net, a distraction-aware two-stage U-Net, to reduce false positives and false negatives in lung nodule segmentation. The presented model uses a Distraction Attention Module (DAM) which can identify the true nodules and noise surrounding them, thus improving segmentation. The presented framework was evaluated on the MICCAI 2019 Gross Target Volume dataset and reported state-of-the-art results under distraction-prone conditions. These studies collectively show remarkable improvements in lung cancer segmentation using deep learning, in which models that use Transformers, attention, and multi-scale processing improve accuracy and efficiency. Subsequent studies can build on this by using multimodal data and generalizing to patient populations.

4 Methodology

The lung tumor segmentation project develops a sophisticated pipeline to accurately delineate lung nodules in CT images, addressing challenges such as diverse tumor sizes, irregular shapes, and complex backgrounds. By leveraging the TransUNet architecture, the approach integrates convolutional neural networks for local feature extraction with Transformer-based mechanisms for global context, enhanced by novel modifications tailored to medical imaging. The methodology encompasses several specialized components, each contributing to the robustness and precision of the segmentation system.

4.1 Data Collection

The project relies on the LUNA16 dataset, a widely recognized resource for lung nodule analysis, which provides 3D CT scans stored in MetaImage format along with annotations specifying nodule locations. These scans, sourced from the LUNA16 repository, are organized into subdirectories, each containing an image file for the CT volume and a corresponding annotation file for the nodule mask. To facilitate comprehensive model training and evaluation, the dataset is partitioned into training, validation, and testing sets, with 70% allocated for training, 10% for validation, and 20% for testing. This division ensures exposure to a diverse range of lung anatomies and tumor characteristics, enabling the model to generalize effectively across different patient scans.

4.2 Data Preprocessing

Preprocessing transforms the LUNA16 dataset into a format suitable for deep learning, converting 3D CT volumes into 2D slices for model compatibility. Using SimpleITK, each CT scan is read as a 3D NumPy array and sliced along the axial plane to produce individual 2D images, with corresponding annotation masks processed similarly to yield binary maps where non-zero values indicate nodules. Pixel intensities are normalized to a 0-to-1 range by dividing by the maximum value, ensuring uniformity across scans, while masks are binarized to clearly distinguish background from nodule regions. To prevent data leakage, the dataset is split at the patient level, maintaining independence between training, validation, and testing sets. The resulting 2D slices and masks are saved as NumPy files in structured directories, labeled by patient ID and slice index for traceability, with a 3-channel visualization mask highlighting nodules in blue. All images are resized to a 224x224 resolution to align with the input requirements of the TransUNet model, ensuring seamless integration into the segmentation pipeline.

4.3 CNN Encoder Deployment with ResNet CNN

The CNN encoder, a critical component of the TransUNet architecture, draws inspiration from ResNet to extract detailed spatial features from CT slices. It begins with a 7x7

convolutional layer that generates 64 feature maps at a stride of 2, followed by batch normalization and ReLU activation to stabilize and accelerate training. A 3x3 max-pooling layer with a stride of 2 reduces spatial dimensions, enhancing computational efficiency. The encoder is structured into four sequential layers, each comprising multiple ResNet blocks with increasing complexity: the first layer includes two blocks with 64 filters, the second has two blocks with 128 filters, the third features two blocks with 256 filters, and the fourth contains two blocks with 512 filters. Each ResNet block employs two 3x3 convolutions, batch normalization, ReLU activations, and a residual connection to mitigate vanishing gradients, enabling the capture of intricate patterns such as tumor edges and textures. Implemented as a dedicated module, the encoder produces multi-scale feature maps at five levels, with channel counts of 64, 64, 128, 256, and 512, which are stored for integration with Transformer outputs and used as skip connections in the decoder, ensuring that fine-grained details are preserved throughout the segmentation process.

4.4 Dual Attention Block

To enhance the Transformer’s ability to focus on relevant tumor features, a dual attention block is introduced, combining spatial and channel attention mechanisms within selected layers. This innovative module employs standard multi-head self-attention to model spatial relationships between image patches, capturing dependencies across the lung field. Simultaneously, a channel attention pathway processes feature maps through a sequence of linear transformations, layer normalization, GELU activation, and a sigmoid function to compute weights that emphasize informative channels. The outputs of these pathways are concatenated and refined through a fusion layer, producing a feature representation that balances spatial context with channel importance. Applied in Transformer layers 3, 6, and 9, this block enhances the model’s sensitivity to tumor structures, improving segmentation accuracy for nodules with complex morphologies while maintaining computational efficiency by selectively enhancing key layers.

4.5 Feature Refinement Module

The feature refinement module is designed to enhance decoder features by integrating spatial and channel attention, ensuring high-quality inputs for upsampling and segmentation. In the channel attention branch, adaptive average and max pooling compress feature maps, which are then processed by a two-layer MLP to generate channel-wise weights that highlight tumor-relevant channels. The spatial attention branch computes mean and max feature maps across channels, concatenates them, and applies a 7x7 convolution to produce a spatial attention map that focuses on tumor regions. These attention mechanisms are applied sequentially, with a residual connection preserving original features to avoid information loss. Incorporated into the decoder blocks, this module reduces noise from background lung tissue, sharpens tumor boundaries, and improves overall segmentation precision, particularly for nodules with irregular edges or low contrast against surrounding tissues.

4.6 Stochastic Depth

To improve model generalization and prevent overfitting, stochastic depth is employed within the Transformer encoder, randomly dropping entire blocks during training to encourage robust feature learning. This technique assigns a probability to each block, with deeper blocks more likely to be skipped, implemented through a drop path mechanism that transitions to an identity operation during inference to retain all layers. By simulating a variety of network depths, stochastic depth acts as a regularization strategy, enabling the model to adapt to the relatively limited size of the LUNA16 dataset. This approach enhances training stability and ensures that the model learns resilient features, capable of generalizing across diverse CT scans without memorizing training data patterns.

4.7 Multiscale Skip Connection

The multiscale skip connection mechanism optimizes the integration of features from the CNN and Transformer encoders, addressing the challenge of combining local and global information across different scales. A dedicated fusion module employs learnable weights to balance contributions from CNN feature maps, which capture fine details, and Transformer outputs, which provide contextual understanding. This module adjusts channel dimensions using 1x1 convolutions and aligns spatial resolutions via bilinear interpolation when necessary, ensuring compatibility between pathways. A learnable scaling factor further refines the fusion process, allowing the model to adaptively emphasize relevant features at each decoder level, corresponding to channel counts of 512, 256, 128, and 64. By enabling dynamic feature integration, this mechanism enhances the model's ability to segment tumors of varying sizes and shapes, improving robustness and accuracy in complex lung environments.

4.8 Boundary-Aware Loss Function

The boundary-aware loss function is a novel contribution that prioritizes accurate delineation of tumor boundaries, critical for clinical applications. This loss combines standard cross-entropy with a weighted term that emphasizes pixels near nodule edges, identified using Sobel filters applied in the x and y directions to detect gradient changes in the ground truth mask. Pixels in boundary regions are assigned a weight of 5.0, amplifying their contribution to the loss and encouraging the model to focus on precise contour segmentation. The weighted loss is computed as the mean of the product of cross-entropy and the boundary weight map, ensuring balanced optimization across the image. By enhancing sensitivity to tumor edges, this loss function improves segmentation performance, particularly for small or irregularly shaped nodules, aligning the model's predictions closely with clinical requirements for accurate diagnosis and treatment planning.

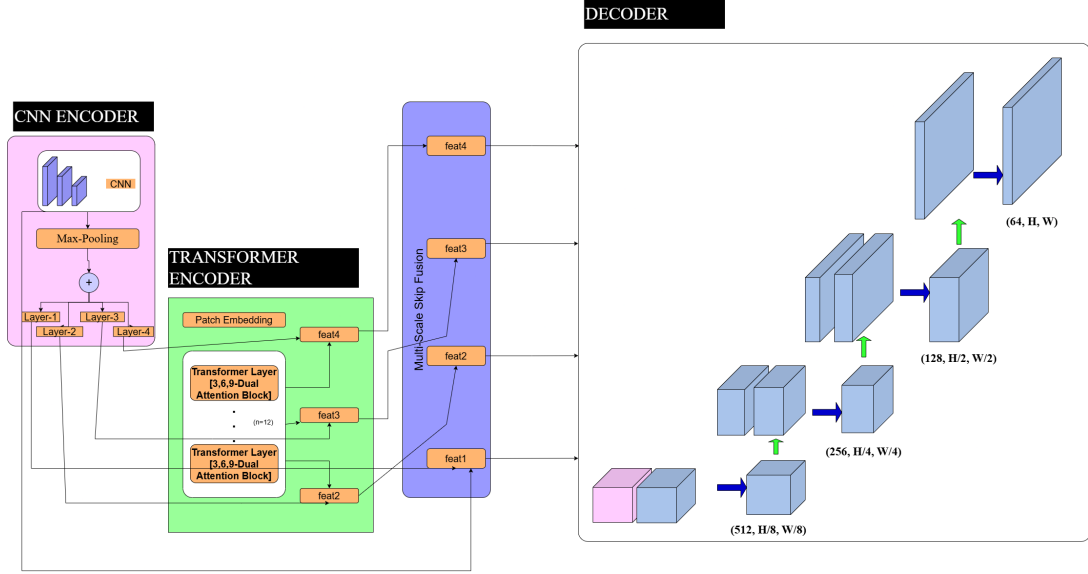


Figure 1: Model Architecture

5 Results

The evaluation of the lung tumor segmentation model is supported by a series of visualizations that highlight its performance across the LUNA16 test set. The histogram of F1 Scores reveals a pronounced peak at 1.0 with a frequency exceeding 40, indicating that a significant majority of predictions achieve near-perfect harmony between precision and recall. Scores below 0.9 are sparsely distributed, with minimal occurrences around 0.75, suggesting consistent high performance. The histogram of IoU Scores mirrors this trend, showing a strong concentration at 1.0 with a frequency approaching 40, reflecting excellent overlap between predicted and ground truth masks, while scores below 0.8 are negligible, underscoring the model’s reliability. The histogram of DICE Scores also peaks at 1.0 with a frequency near 40, tapering off below 0.9, which indicates robust agreement with true annotations and validates the segmentation accuracy. The histogram of Recall Scores peaks at 1.0 with a frequency around 70, demonstrating the model’s exceptional ability to detect all relevant tumor regions, with a gradual decline to near zero at 0.8, further confirming its effectiveness. The correlation matrix of evaluation metrics illustrates strong positive correlations, with DICE, IoU, F1, and Recall scores exceeding 0.8, indicating cohesive performance improvements, while a negative correlation of -0.09 between specificity and recall suggests a potential trade-off in certain scenarios, providing a nuanced view of the model’s behavior. The comparison of an original CT image with its ground truth and predicted masks shows a central tumor region marked with blue in both the ground truth and prediction, closely aligning with the anatomical features of the original image, which highlights the model’s precision in tumor localization. The training performance plots offer additional insights, with the loss plot showing training and validation loss decreasing from 0.25 to near 0.0 over 50 epochs, indicating successful convergence without overfitting. The DICE score plot shows a validation score rising to approximately 0.9, reflecting progressive improvement in segmentation quality, while the IoU plot reaches about 0.85, further supporting this trend. The learning rate schedule plot starts at 10^{-4} and declines to 10^{-6} , ensuring stable training dynamics that contribute

to the model’s consistent performance.

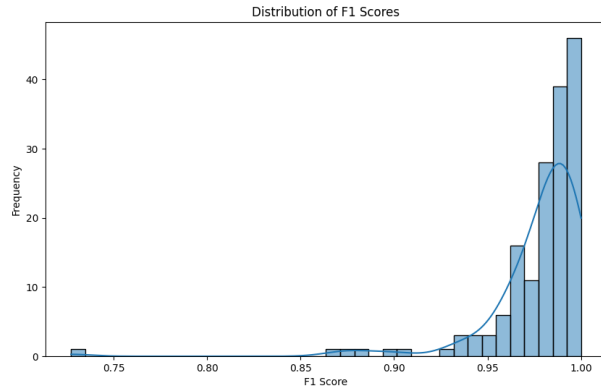


Figure 2: Distribution of F1 Scores

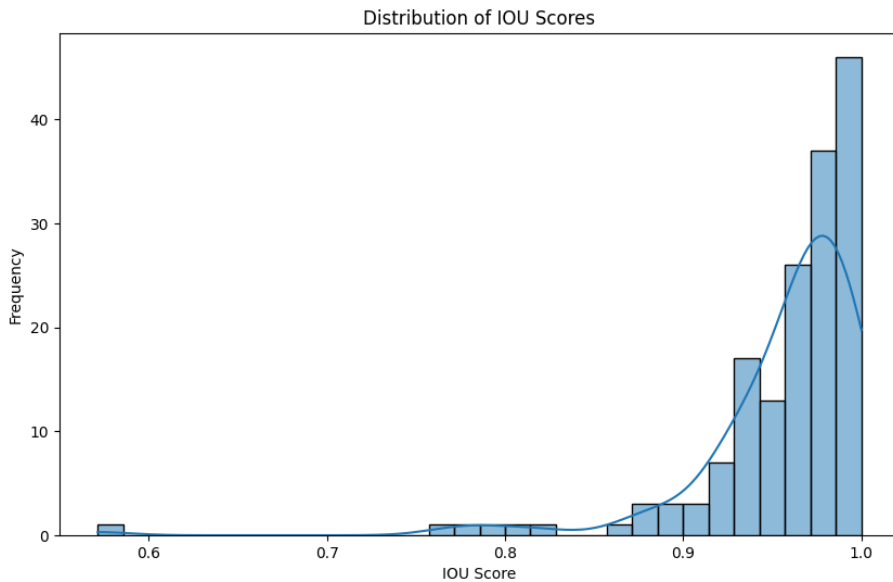


Figure 3: Distribution of IoU Scores

6 Conclusion

The lung tumor segmentation project successfully delivers a TransUNet-based model that integrates convolutional and Transformer-based architectures to achieve precise nodule segmentation in CT images. The methodology leverages the LUNA16 dataset, preprocesses 3D scans into 2D slices, and deploys a sophisticated model enhanced by novel components such as dual attention, feature refinement, stochastic depth, multiscale skip connections, and a boundary-aware loss function. These innovations address critical challenges in medical imaging, including small nodule detection, boundary precision, and generalization across varied scans. The modular design facilitates future enhancements,

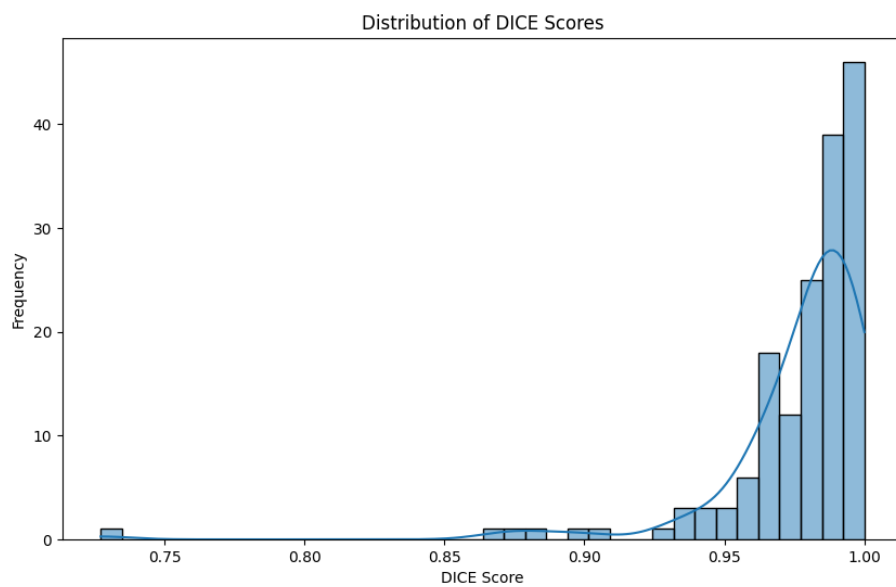


Figure 4: Distribution of DICE Scores

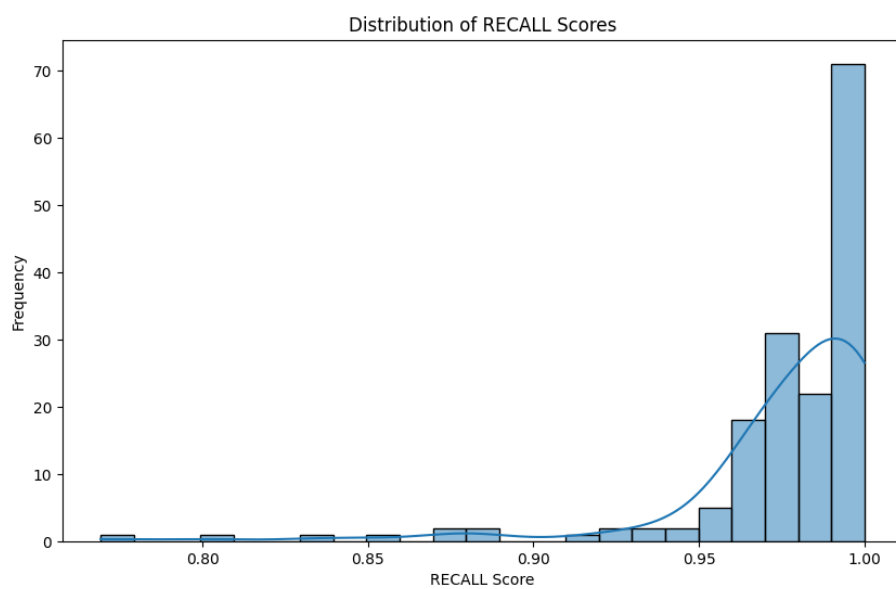


Figure 5: Distribution of Recall Scores

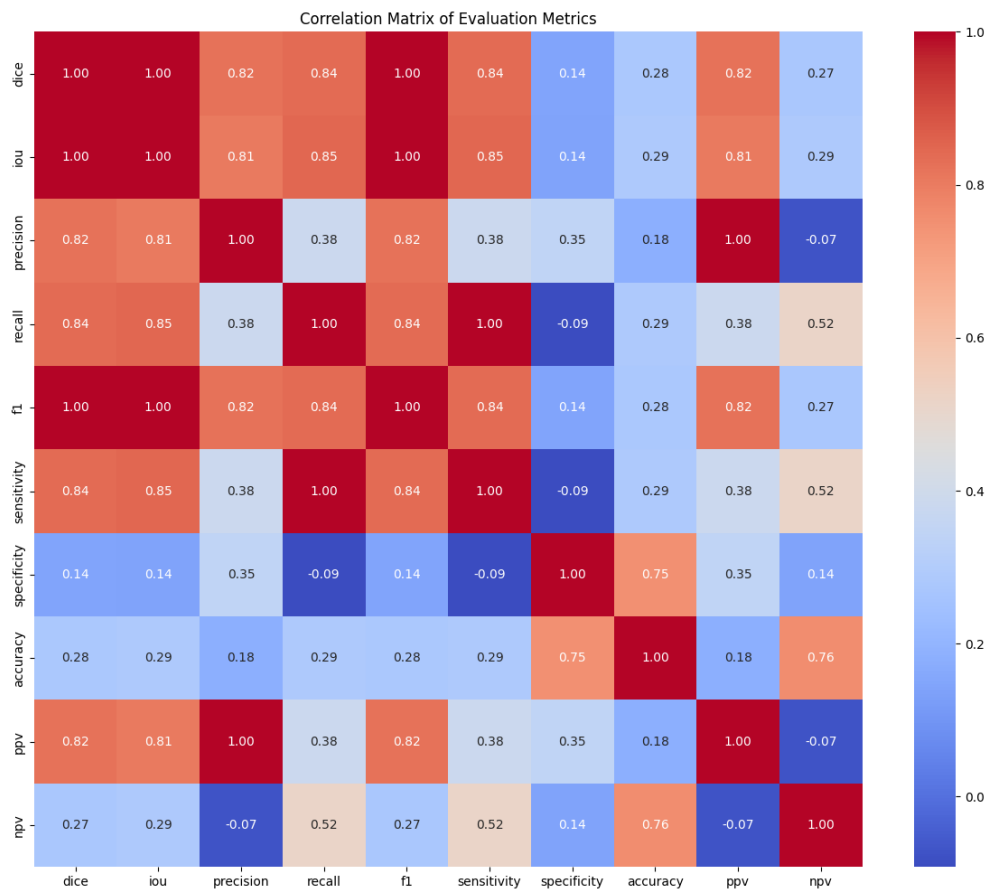


Figure 6: Correlation Matrix of Evaluation Metrics

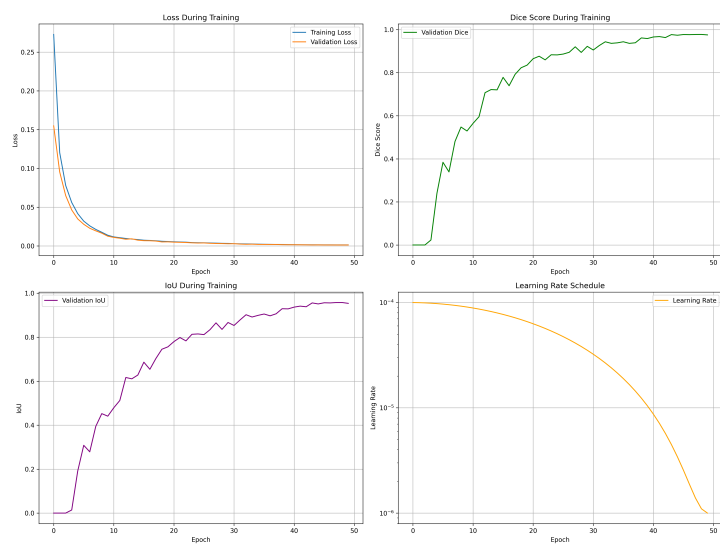


Figure 7: Loss During Training

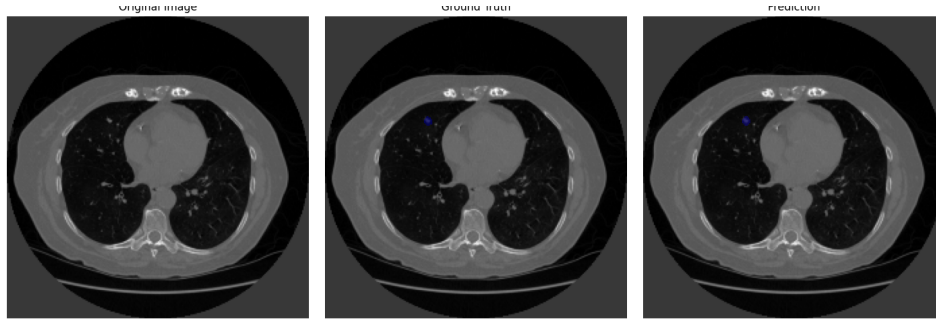


Figure 8: Comparison of Original Image, Ground Truth, and Prediction

such as expanding to other lung abnormalities or integrating into clinical workflows, establishing a robust foundation for advancing computer-aided diagnosis in lung cancer care.

7 References

1. Zhao, Juntong, et al. "DSU-Net: Distraction-Sensitive U-Net for 3D lung tumor segmentation." *Engineering Applications of Artificial Intelligence* 109 (2022): 104649.
2. Chen, Wei, et al. "MAU-Net: Multiple attention 3D U-Net for lung cancer segmentation on CT images." *Procedia Computer Science* 192 (2021): 543-552.
3. Le, Van-Linh, et al. "RRc-UNet 3D for lung tumor segmentation from CT scans of Non-Small Cell Lung Cancer patients." *IEEE/CVF International Conference on Computer Vision*, 2023.
4. Said, Yahia, et al. "Medical images segmentation for lung cancer diagnosis based on deep learning architectures." *Diagnostics* 13, no. 3 (2023): 546.
5. Kamal, Uday, et al. "Lung cancer tumor region segmentation using recurrent 3d-denseunet." *MICCAI 2020*, Lima, Peru.
6. Li, Zhang, et al. "Deep learning methods for lung cancer segmentation in whole-slide histopathology images—the acdc@ lunghp challenge 2019." *IEEE Journal of Biomedical and Health Informatics* 25, no. 2 (2020): 429-440.
7. Wang, Hongfeng, et al. "Accurate classification of lung nodules on CT images using the TransUnet." *Frontiers in Public Health* 10 (2022): 1060798.
8. X. Fan, Y. Lu, J. Hou, F. Lin, Q. Huang and C. Yan, "DMC-UNet-Based Segmentation of Lung Nodules," in *IEEE Access*, vol. 11, pp. 110809-110826, 2023, doi: 10.1109/ACCESS.2023.3322437.
9. J. Hou, C. Yan, R. Li, Q. Huang, X. Fan and F. Lin, "Lung Nodule Segmentation Algorithm With SMR-UNet," in *IEEE Access*, vol. 11, pp. 34319-34331, 2023, doi: 10.1109/ACCESS.2023.3264789.