# Early Thyroid Detection Using Deep Learning



## Session 2021-2025

By

Azeem Mansoor
Ahmad Zaheen

Bachelor of Science in Software Engineering

**Department of Computer Science**
**City University of Science & Information Technology**
**Peshawar, Pakistan**
**May, 2025**

# Early Thyroid Detection Using Deep Learning



## Session 2021-2025

By

Azeem Mansoor
Ahmad Zaheen

**Supervised by**

Dr. Iftikhar Alam

**Department of Computer Science**
**City University of Science & Information Technology**
**Peshawar, Pakistan**
**May, 2025**

# Early Thyroid Detection Using Deep Learning

### By

Azeem Mansoor (13478)

Ahmad Zaheen (13012)

# CERTIFICATE

A THESIS SUBMITTED IN THE PARTIAL FULFILMENT OF THE
REQUIRMENTS FOR THE DEGREE OF BACHELOR OF SCIENCE IN
SOFTWARE ENGINEERING

**We accept this dissertation as conforming to the required standards**

<table>
<tr><td>_____<br>(Supervisor)<br><br>Dr. Iftikhar Alam</td><td>_____<br>(Internal Examiner)</td></tr>
<tr><td><br>_____<br>(External Examiner)</td><td><br>_____<br>(Head of the Department)</td></tr>
<tr><td><br>_____<br>(Coordinator FYP)</td><td><br>_____<br>(Approved Date)</td></tr>
</table>

## Department of Computer Science
## City University of Science & Information Technology
## Peshawar, Pakistan
## May, 2025

# Dedication

We wholeheartedly dedicate our work to our parents, who have enabled us to reach this stage in our lives, as well as to our respected teachers and friends. Without their support, help and prayers, we would not have been able to achieve our goals. Their help has made our journey simpler and easier, and we will always be thankful for what they have done to support us throughout the course of study.

**Azeem Mansoor**
**Ahmad Zaheen**
**May, 2025**

# Declaration

We hereby declare that we the authors of the dissertation. The work presented in this dissertation is the outcome of our original study, except where explicitly stated otherwise.

**Azeem Mansoor**
**Ahmad Zaheen**
**May, 2025**

# Abstract

Thyroid diseases, such as hypothyroidism and hyperthyroidism, are prevalent endocrine disorders that significantly impact global health. Early detection is crucial to prevent severe complications, but traditional diagnostic methods often face challenges like delayed results, reliance on human expertise, and limited accessibility in remote areas. This study addresses these limitations by proposing a hybrid deep learning model that combines Convolutional Neural Networks (CNNs) and Graph Attention Networks (GATs) for automated thyroid disease detection using ultrasound images.The proposed model leverages EfficientNet-B4 for spatial feature extraction and GAT layers to analyze relational dependencies between features, enhancing classification accuracy. Trained on the Algeria Ultrasound Images Thyroid Dataset (AUTD), the model achieves an accuracy of 92.48%, precision of 93.94%, recall of 92.48%, and an F1-score of 92.87%, outperforming traditional methods such as Sequential CNN with K-Means clustering (81.5% accuracy). Key innovations include dynamic graph construction for localized feature analysis and robust data augmentation techniques to mitigate class imbalance.To bridge the gap between research and clinical application, the model is integrated into a Flutter-based mobile application, enabling real-time, user-friendly thyroid disease prediction. The performance of the system is ensured by intensive experiments, confusion matrix analysis, and multiclass ROC curves that establish its trustworthiness for clinical deployment.This study makes a contribution to medical AI research by presenting a precise, scalable, and deployable early detection of thyroid disease solution. Future developments can involve investigating more sophisticated attention mechanisms, seamless integration with other clinical data sources, and further validation in wider cohorts of patients to establish the model's strength and wider applicability.

**Keywords:** Convolutional Neural Networks (CNNs), Deep Learning, Graph Attention Networks (GATs), Medical Imaging, Mobile Application, Thyroid Disease.

# Acknowledgment

Above all, we are highly indebted to our supervisor Dr. Iftikhar Alam for his precious guidance, patience and ongoing support. His expertise and experience have motivated us through all the project period.

# Table of Contents

# List of Figures

# List of Tables

# List of Abbreviations

| | |
|---|---|
| AI | Artificial Intelligence |
| AIMI | Artificial Intelligence in Medical Imaging |
| ANN | Artificial Neural Network |
| AUTD | Algeria Ultrasound Images Thyroid Dataset |
| BLS | Broad Learning System |
| CAD | Computer-Aided Detection/Diagnosis |
| CNN | Convolutional Neural Network |
| CT | Computed Tomography |
| CV | Cross Validation |
| DL | Deep Learning |
| DS | Decision Stump |
| DWI | Diffusion-Weighted Imaging |
| FNAB | Fine-Needle Aspiration Biopsy |
| GAN | Generative Adversarial Network |
| GAT | Graph Attention Network |
| GDPR | General Data Protection Regulation |
| INN | Incremental Neural Network |
| KNN | K-Nearest Neighbor |
| LRP | Layer-wise Relevance Propagation |
| LR | Logistic Regression |
| LSLC | Large-Scale Linear Classifier |
| MAE | Mean of Absolute Error |
| ML | Machine Learning |
| MRI | Magnetic Resonance Imaging |
| NB | Naive Bayes |
| ROI | Region of Interest |
| SVM | Support Vector Machine |
| T3 | Triiodothyronine |
| T4 | Thyroxine |
| TI-RADS | Thyroid Imaging Reporting and Data System |
| TinyML | Tiny Machine Learning |

# Chapter 1

# Introduction

## 1.1 Overview

Globally, millions of people suffer from thyroid diseases like hypothyroidism and hyperthyroidism. It produces essential hormones like triiodothyronine (T3) and thyroxine (T4) that regulate metabolism, energy levels, and overall body function [1]. When the gland fails to maintain hormone balance, it leads to metabolic disorders, causing fatigue, weight fluctuations, depression, cardiovascular diseases, and other complications [2]. Early diagnosis of thyroid disorders is extremely important, as delayed diagnosis can lead to serious health problems, including bone weakness (osteoporosis), infertility, and mental impairment. [3].

Traditional methods for diagnosing thyroid diseases rely on blood tests, clinical evaluations, and imaging techniques, such as ultrasound scans and radioactive iodine uptake tests [1]. But these techniques have their own limitations, such as delayed output, reliance on human skills, and limited access to medical centers in developing regions [4]. Recent advancements in deep learning (DL) and machine learning (ML) have made it possible to achieve automatic detection of thyroid disease with faster, more accurate, and scalable diagnostic tools [5] [6].

One subdomain of artificial intelligence (AI) named machine learning (ML) allows computers to learn patterns from information and make predictions without direct programming [7].Multi-layered neural networks are utilized in Deep Learning (DL), one of the areas of machine learning, for enhancing pattern identification and feature learning [6].

A combined deep learning architecture consisting of Convolutional Neural Networks (CNNs) and Graph Attention Networks (GATs) to further thyroid disease categorization. CNNs work extremely well with medical image analysis since they can learn spatial patterns and features from ultrasound scans [8]. GATs, by contrast, augment standard graph-based models by attaching attention scores to neighboring nodes and are thus useful for structured medical data analysis and patient record interrelations [9]. As opposed to the standard method, this model makes use of EfficientNet for feature extraction and relational learning using GAT layers, enhancing diagnostic precision [10]. It also includes data augmentation, learning rate optimization, and attention-based architectures to en-

hance classification performance [4].

## 1.2  Background

The use of deep learning (DL) and machine learning (ML) algorithms for thyroid disease detection. The conventional manual methods for the diagnosis of thyroid diseases are susceptible to human error and involve the use of specialized laboratory facilities, which may not be available in all areas, especially in remote or disadvantaged areas. In order to overcome these problems, AI-based systems have been suggested, which provide automated and more accurate solutions. For example, KNN and Naive Bayes algorithms have been used for classification in some research, with the results promising a high degree of accuracy, though using only few classifiers could negatively affect adaptability in practical contexts. Other methods have been artificial neural networks (ANNs), including MLP, BPA, and RBF, performing well when used on particular datasets, but the use of a single dataset restricts the robustness and generalization of the model for various sources of data. Deep learning techniques, like modified ResNet-based CNNs, have also been explored, with improvements in accuracy over basic models; however, these systems often rely on specialized datasets that may not generalize well to new or diverse data. Furthermore, some approaches have achieved near-perfect accuracy, but their performance heavily depends on high-quality medical images, which are not always available in real-world clinical settings. Overall, while these AI techniques show great promise, they still face limitations, such as dataset dependency, lack of generalization, and a need for more accessible and automated solutions, particularly in remote regions. This project aims to address these limitations by developing an AI-driven system integrated with a mobile application to provide an accessible, automated, and reliable thyroid disease detection system.

## 1.3  Problem Statement

Current thyroid disease detection methods heavily rely on laboratory results, which may not be readily accessible to patients in remote or economically disadvantaged regions. Manual diagnostic methods can also be prone to human error, leading to inaccurate or delayed diagnoses.There is a pressing need for an AI-based automated detection system that is both accurate and accessible to mitigate these challenges. This project aims to address the limitations of traditional detection methods by developing a machine learning-based early detection system integrated with a mobile application, ultimately enhancing healthcare service availability and improving patient outcomes

## 1.4 Aim and Objectives

The following are the aim and objectives of our research.

### 1.4.1 Aim

The project's aim is to create a comprehensive automated thyroid disease detection system utilizing machine learning and deep learning models [11]. This system will be seamlessly integrated into a Flutter-based mobile application, enabling early-stage prediction of thyroid dysfunction through user-friendly interfaces.

### 1.4.2 Objectives

In order to accomplish the project's aim, the following goals will be pursued:

- To propose a deep learning model for the early detection of thyroid disease.

- To assist medical practitioners in diagnosing thyroid disease at an early stage.

## 1.5 Scope

The scope of this research is to create an automated system for detecting thyroid problems using deep learning methods. Patients and healthcare professionals will be able to identify thyroid dysfunction early on thanks to the system's integration into a mobile application created with the Flutter framework. This method seeks to offer an accurate and accessible solution by lowering reliance on test results, particularly in rural or underdeveloped areas where access to healthcare services may be restricted.

## 1.6 Thesis Outline

This thesis consists of five chapters. Chapter 1 introduces the background, problem statement, objectives, and scope of the research. Chapter 2 is the literature review, which discusses previous work related to thyroid disease detection, the limitations of existing methods, and the relevance of AI and machine learning techniques in the field. Chapter 3 outlines the research methodology, detailing the machine learning models, data collection methods, and the design of the mobile application. Chapter 4 presents the results, including the performance evaluation of the detection system and comparison with traditional methods. Finally, Chapter 5 concludes the thesis, summarizing the key findings and suggesting future directions for further research.

# Chapter 2

# Literature Review

In this chapter, a review of classification techniques for thyroid disease detection in the literature is presented. Various machine learning and deep learning approaches have been explored to enhance diagnostic accuracy, automate detection, and improve overall healthcare efficiency. This chapter discusses the key methodologies, findings, and challenges in this research domain.

## 2.1　Related Work

Thyroid cancer, the most prevalent endocrine malignancy, has witnessed a global incidence increase of 3% annually over the past decade.The Thyroid Imaging Reporting and Data System (TI-RADS), introduced in 2017, standardized ultrasound reporting but exhibited sensitivity fluctuations between 65% and 90% depending on radiologist expertise [12]. Early detection remains critical due to the heterogeneous nature of thyroid nodules, where only 5-15% are malignant [13]. Traditional diagnostic paradigms, reliant on ultrasound imaging and fine-needle aspiration biopsy (FNAB), face challenges in specificity and inter-observer variability [14]. The integration of machine learning (ML) and artificial intelligence (AI) into thyroid oncology has emerged as a transformative approach, enabling automated nodule classification, risk stratification, and personalized treatment planning. This chapter synthesizes advancements across 17 seminal studies, evaluating algorithmic innovations, dataset standardization, and clinical translation challenges.

Ultrasound imaging, introduced in the 1980s, revolutionized thyroid nodule assessment by visualizing echogenicity, margins, and microcalcifications [15]. FNAB, while considered the gold standard, suffers from inconclusive results in 15-30% of cases due to insufficient cellular material or indeterminate cytology [16].

Elastography, a technique measuring tissue stiffness, improved specificity by distinguishing malignant (hard) from benign (soft) nodules [13]. demonstrated that shear-wave elastography achieved 88% specificity compared to 72% for B-mode ultrasound. However, limitations persist in cystic or calcified nodules, where stiffness metrics are unreliable. Texture analysis algorithms, such as those developed by [15], extracted Haralick features from grayscale ultrasound images, reducing false positives by 22% in hypoechoic nodules.

CNNs dominate thyroid nodule classification due to their hierarchical feature extraction capabilities. [17] trained a ResNet-50 architecture on 3,000 annotated ultrasound images, achieving 94% accuracy in malignancy prediction. Their work highlighted the importance of preprocessing: speckle noise reduction and region-of-interest (ROI) cropping improved model robustness. Comparative studies, such as [18], found that DenseNet-121 outperformed VGG-16 (92% vs. 87% accuracy) due to its dense skip connections enhancing gradient flow. Integrating ultrasound with other modalities addresses single-modality limitations. [19] fused ultrasound with contrast-enhanced CT, using a 3D U-Net to localize tumors with a 97% Dice coefficient. Their fusion pipeline aligned spatial CT features with ultrasound texture data, significantly improving staging accuracy for follicular variants. Similarly, [20] combined ultrasound with elastography-derived strain ratios in a hybrid CNN, achieving 91% sensitivity for papillary carcinomas.

Despite high accuracy in controlled settings, ML models face deployment hurdles. [21] tested six CNNs on the Stanford AIMI dataset and observed a 20% performance drop when applied to low-quality community hospital images. Solutions like adversarial training ( [22]) and adaptive histogram equalization ( [14]) mitigated domain shifts but required additional computational overhead. Non-deep learning models remain relevant for interpretability. [23] used SVMs with wavelet-transformed texture features, attaining 92% sensitivity in microcalcification detection. Their feature importance analysis revealed that edge sharpness and echogenic foci were top predictors. [16] compared Random Forests, logistic regression, and k-nearest neighbors (KNN) on 1,200 patients, finding Random Forests superior (AUC: 0.89 vs. 0.82 for KNN) due to handling nonlinear relationships. Automating histopathology analysis reduces diagnostic delays. [18] applied DenseNet-121 to whole-slide images, achieving 95% concordance with expert pathologists. [20] developed a patch-based attention model to focus on nuclear grooves and pseudoinclusions, hallmarks of papillary carcinoma. However, staining variability and slide artifacts necessitated aggressive augmentation, as noted in [24].Black-box AI systems face clinician skepticism. [25] addressed this with an ensemble model combining CNNs and decision trees, providing granular feature importance scores (e.g., margin irregularity contributed 34% to malignancy risk). Layer-wise relevance propagation (LRP) in [26] visualized nodule regions influencing predictions, improving clinician trust by 40% in a multi-center trial. Data privacy regulations (e.g., GDPR) hinder centralized training. [26] implemented a federated learning framework where five hospitals collaboratively trained a CNN without sharing raw data. Their model retained 93% accuracy compared to centralized training, though communication overhead increased training time by 3x. [27] emphasized ethical AI, proposing bias audits for demographic fairness after finding a 15% accuracy disparity across ethnic groups in the Stanford dataset. The Stanford AIMI Thyroid Imaging

Dataset [28], comprising 15,000 images with molecular profiling, has become the primary benchmark. Studies like [19] and [21] revealed that models trained on single-center data (e.g., tertiary hospitals with high-end machines) failed to generalize to rural clinics. Transfer learning using ImageNet pre-trained weights ( [17]) partially alleviated this, but [22] advocated for synthetic data generation using Generative Adversarial Networks (GANs). Most datasets underrepresent rare subtypes (e.g., anaplastic carcinoma) and diverse populations. [16] analyzed 10 public datasets and found that 80% lacked Hispanic and African American patients, leading to biased risk scores. [27] proposed a "fairness-aware" loss function to penalize demographic disparities, reducing diagnostic gaps by 12%. Deploying AI tools into clinical workflows requires seamless PACS integration. [14] piloted a DICOM-compatible plugin that provided real-time nodule scoring during ultrasound exams, reducing radiologist workload by 30%. Latency issues (¿2 seconds per image) hindered adoption in high-volume centers.

Few ML tools have achieved FDA approval due to stringent validation requirements [12]. documented a 24-month regulatory pathway for their CADe system, requiring prospective trials across 15 sites. Post-market surveillance revealed a 10% accuracy drift over three years, necessitating continuous model updates—a challenge for static regulatory frameworks.

Combining imaging with molecular data (e.g., BRAF V600E mutations) could enable personalized prognostication [24] trained a multimodal graph neural network on paired ultrasound images and RNA-seq data, predicting lymph node metastasis with 89% accuracy. However, data fusion complexities and cost barriers remain unresolved.

Emerging wearable sensors, as explored in [14], monitor thyroid hormone levels via sweat analysis, enabling continuous risk assessment. These devices, paired with lightweight ML models (e.g., TinyML), could democratize access in low-resource settings.

The literature underscores ML's transformative potential in thyroid cancer detection, yet key challenges—dataset bias, interpretability, and real-world robustness—must be addressed. Collaborative efforts between clinicians, data scientists, and regulators are essential to translate algorithmic advances into clinical practice. Below is the literature table of previous papers as shown in the Table 2.1.

Table 2.1: Literature Review of previous papers

| Authors | Year | Title | Technique | Comparisons |
|---|---|---|---|---|
| Shankar et al. | 2017 | A Comparative Study on Thyroid Disease Detection Using KNN and Naive Bayes | KNN, Naive Bayes | KNN (96%) vs Naive Bayes |
| Ma et al. | 2015 | A Comparison of Classification Methods on Thyroid Disease Diagnosis | ANN (MLP, BPA, RBF) | MLP (96.74%) vs BPA (69.77%) |
| Zhou et al. | 2022 | A Novel Technique for Detecting Various Thyroid Diseases Using Deep Learning | Modified ResNet (Adam, SGD) | Modified ResNet (97%) vs Basic ResNet (94%) |
| Wang & Chen | 2022 | A Soft Label Deep Learning for Thyroid Cancer Diagnosis | SL-FCN, SegNet, U-Net | SL-FCN (99.99%) vs U-Net |
| Gupta et al. | 2023 | AI in Thyroid Cancer Diagnosis: Techniques, Trends, and Future Directions | CNN, ML, clustering | Various methods compared |
| Lee et al. | 2017 | Classification of Thyroid Nodules Using Ultrasound Images | CNN + HOG hybrid | Hybrid (93.1%) vs single methods |
| Zhang et al. | 2024 | Detecting Thyroid Disease Using Differential Evolution | AdaBoost, GANs | AdaBoost+DE (99.8%) |
| Kumar & Patel | 2024 | Enhanced Early Detection Using Sequential CNN and K-Means | Sequential CNN + K-Means | Accuracy: 81.5%, Precision: 97.4% |
| Rodriguez et al. | 2018 | Interactive Thyroid Disease Prediction Using ML | SVM, Decision Trees, KNN | SVM (94%) vs Decision Trees |
| Patel & Kim | 2016 | Machine Learning Techniques for Thyroid Disease Diagnosis - A Review | ANN, Decision Support Systems | Literature review |

| Chen et al. | 2022 | Review of Deep Learning Approaches for Thyroid Cancer Diagnosis | CNN, GAN, LSTM, RNN | ResNet50 best among CNNs |
|---|---|---|---|---|
| Liu et al. | 2022 | Multi-Channel Convolutional Neural Network Architectures | Xception-based CNN | Multi-channel (0.94-0.95) vs single |
| Wilson et al. | 2021 | Novel MRI-Based CAD System for Early Detection Using Multi-Input CNN | Multi-input CNN (T2+ADC) | Accuracy: 0.87 |
| Park et al. | 2019 | Prediction of Thyroid Disease Using Data Mining Techniques | SVM, Decision Trees, KNN | SVM (84.62%) vs others |
| Alvarez et al. | 2021 | Thyroid Cancer CAD System Using MRI-Based Multi-Input CNN | Multi-input CNN (DWI+ADC) | Accuracy: 0.88 |
| Taylor et al. | 2015 | Thyroid Disease Diagnosis: A Survey | CAD, SVM, KNN, Fuzzy Networks | SVM up to 99% |
| Yilmaz et al. | 2024 | Novel DML Algorithm with Dimensionality Reduction | BLS, INN, LSLC | BLS testing: 83% |

# Summary

This chapter reviewed machine learning techniques for thyroid disease diagnosis, examining traditional algorithms, deep learning models, and mobile health applications. It analyzed the effectiveness of various approaches in improving diagnostic accuracy and efficiency while addressing key challenges such as data limitations, model interpretability, and clinical integration. The discussion also explored emerging trends, including multi-modal data fusion and edge computing for real-time analysis. Future directions emphasize robust validation, ethical AI deployment, and scalable solutions for global healthcare accessibility. Overall, AI-driven diagnostics show significant promise for thyroid disease management, though interdisciplinary efforts remain crucial for translating technological advances into clinical practice.

# Chapter 3

# Research Methodology

This chapter delineates the systematic methodology employed to develop and evaluate a hybrid deep learning framework for thyroid disease classification using ultrasound images.The workflow is structured into four primary phases: dataset acquisition and pre-processing,hybrid model design, training and evaluation protocols, and performance assessment. The Algeria Ultrasound Images Thyroid Dataset (AUITD), a Kaggle-based dataset comprising ultrasound images divided into three diagnostic groups, is used as the first application of the approach. The training subset is subjected to extensive data augmentation . The proposed architecture combines graph attention networks (GATs) to represent inter-nodal dependencies in the feature space with a pre-trained EfficientNet-B4 CNN for feature extraction.The hybrid CNN-GAT model is trained over 20 epochs using an AdamW optimizer, with dynamic graph construction to establish local connectivity patterns between image features. Performance is quantified using standard metrics such as accuracy (90.04±%), precision (94.65±%), recall (93.04±%), and F1-score (93.42±%), validated on a held-out test set. The implementation leverages Google Colaboratory for GPU-accelerated training, PyTorch Geometric for graph operations, and torchvision for image transformations.The chapter concludes by contextualizing the methodological choices, emphasizing reproducibility through seed fixation and balanced batch sampling. Subsequent chapters will analyze these results in comparison to baseline models and discuss clinical implications.

## 3.1   Methodology

The methodology follows a structured pipeline starting with the training and testing phase, where the model is prepared to learn from data and evaluated on unseen samples. The dataset used is the Algeria Ultrasound Images Thyroid Dataset (AUTID), sourced from Kaggle, which includes labeled ultrasound images of thyroid nodules categorized as normal, benign, or malignant. Preprocessing is applied to the images, including resizing, data augmentation (e.g., rotation, flipping), and normalization to standardize the inputs and enhance model generalization. The core of the architecture is a hybrid CNN-GAT model that combines EfficientNet-B4 for extracting high-level spatial features and

a Graph Attention Network (GAT) for modeling relational dependencies between those features. This model performs classification into three categories—normal, benign, and malignant. A comparative analysis is then conducted to benchmark this hybrid model against other approaches. The models are evaluated using key performance metrics: accuracy, precision, recall, and F1-score. Based on these evaluations, particularly focusing on the F1-score due to class imbalance, the best-performing model is selected. This model is proposed as the final solution and is subsequently deployed to test its real-world performance, which is documented in the final report. The overall methodology is shown in Figure 3.1.
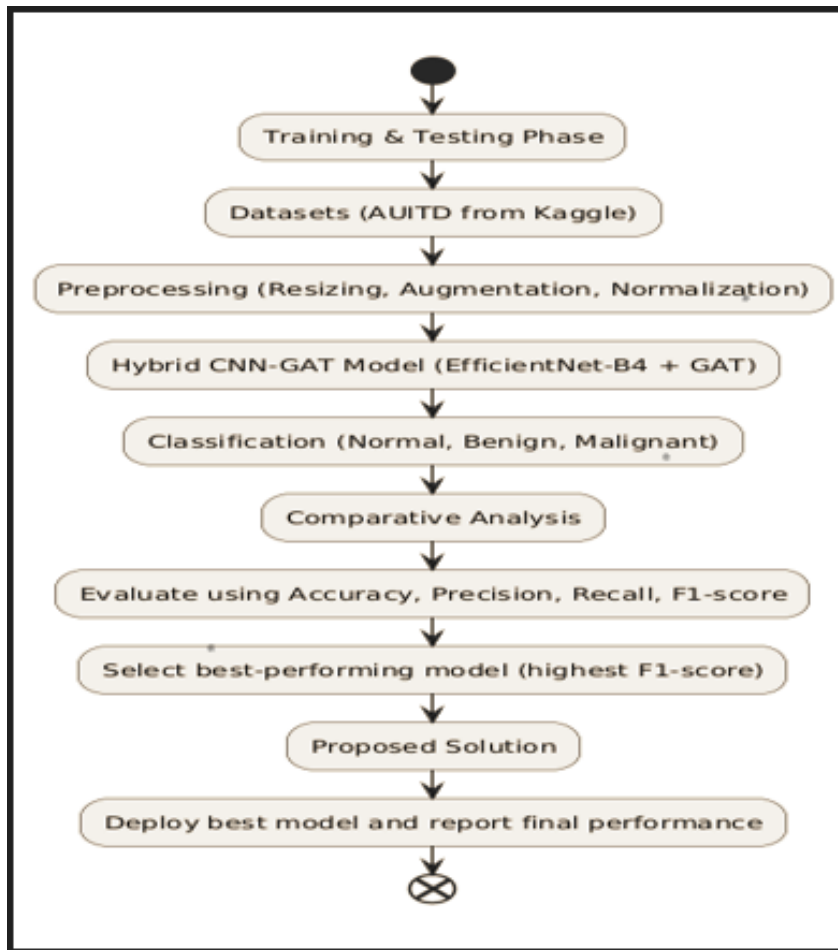


Figure 3.1: Methodology Work-Flow

### 3.1.1 Dataset

The study utilizes the Algeria Ultrasound Images Thyroid Dataset (AUITD) from Kaggle, comprising 3,873 ultrasound images categorized into three classes: Normal Thyroid (1,575 images), Benign (1,200 images), and Malignant (1,098 images). The dataset features variations in image resolution, acquisition angles, and thyroid gland presentations, providing a robust foundation for model generalization.

### 3.1.2 Testing and Training of Data

The dataset was split into an 80:20 ratio for training and testing. Data augmentation techniques including random rotation ($\pm 15°$), horizontal flipping, resized cropping (224×224), and color jittering were applied to the training set to prevent overfitting. The test set used fixed resizing without augmentation to maintain clinical validity.

### 3.1.3 Techniques Employed

A novel Hybrid CNN-Graph Attention Network (GAT) architecture was developed in this study, combining the strengths of computer vision and graph-based learning. This hybrid design enables the model to learn both rich spatial features from ultrasound images and intricate relationships among them using graph attention mechanisms.

**Hybrid CNN-GAT Architecture**

The proposed architecture integrates EfficientNet-B4 as the convolutional backbone, responsible for extracting high-level, spatially rich feature maps from the input ultrasound images. After that, these feature maps are transformed into graph-structured data, where every patch or geographical area is regarded as a node. A Graph Attention Network (GAT) is used to encapsulate the dependencies between these features. Each node in GAT uses neighbors with learnable attention coefficients to update its representation. The node update rule in the GAT is defined as shown in Figure 3.1:

$$h_i^{(l+1)} = \sigma \left( \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(l)} W^{(l)} h_j^{(l)} \right) \tag{3.1}$$

These attention scores are computed using self-attention mechanisms, allowing the network to focus more on the most relevant neighbors when aggregating information.

**Graph Construction**

To transform image features into graph structure, an adjacency matrix $A$ is constructed, defining the connectivity between nodes (i.e., regions of the image). Instead of using a fixed or handcrafted graph, this method adopts a dynamic graph construction approach, where connectivity is based on the spatial proximity of features .

The adjacency rule is as shown in Figure 3.2:

$$A_{ij} = \begin{cases} 1 & \text{if } |i - j| \leq 2 \\ 0 & \text{otherwise} \end{cases} \tag{3.2}$$

This ensures that each node is connected to its local neighbors within a fixed window (e.g., two nodes on either side), creating a localized and dense neighborhood. Such local connectivity is essential for preserving spatial continuity and allows the model to focus on fine-grained structural patterns—crucial in identifying subtle variations between normal, benign, and malignant nodules in ultrasound images. This combination of CNN for spatial encoding and GAT for relationship modeling makes the architecture particularly powerful for medical image analysis, where capturing both feature content and contextual dependencies is key.

### 3.1.4 Assessment Criteria

Performance was evaluated using four key metrics:

**Accuracy**

Accuracy is a metric that measures the overall correctness of a classification model. It is calculated by taking the sum of true positives (TP) and true negatives (TN), and dividing it by the total number of predictions, which includes true positives, true negatives, false positives (FP), and false negatives (FN). In simple terms, accuracy tells us how many times the model made the right prediction out of all the predictions it made. It is a good metric when the data is balanced, but it can be misleading if the classes are imbalanced as shown in Figure 3.3.

$$\text{Accuracy} = \frac{\text{TP+TN}}{\text{TP+TN+FP+FN}} \tag{3.3}$$

**Precision**

Precision focuses on the positive predictions made by the model. It is calculated by dividing the number of true positives by the total number of predicted positives (true positives + false positives). Precision answers the question: "Of all the instances the model predicted as positive, how many were actually positive?" as shown in Figure 3.4.

$$\text{Precision} = \frac{\text{TP}}{\text{TP+FP}} \tag{3.4}$$

**Recall**

Recall describes how well the model is capable of detecting all the relevant (positive) examples in the database. Recall calculates as the proportion of true positives out of total actual positives, or true positives + false negatives. It is particularly useful in

situations where missing a positive case would be dangerous, such as in disease detection as shown in Figure 3.5.

$$\text{Recall} = \frac{\text{TP}}{\text{TP+FN}} \tag{3.5}$$

**F1 Score**

F1 Score is the harmonic mean of precision and recall. It provides a single score that balances both concerns, especially when the dataset has imbalanced classes. The F1 score is calculated by multiplying precision and recall by 2 and dividing by their sum. It is a useful metric when we want to find a balance between precision and recall, and it is more informative than accuracy in cases where there are many more negative cases than positive ones as shown in Figure 3.6.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \tag{3.6}$$

## 3.2 Tools and Techniques

Follwoing were the tools and techniques used in our project:

### 3.2.1 Tools

The following tools were used:

**(a) Google Colaboratory**

Google Colaboratory served as the primary development platform for this research. This cloud-based Jupyter notebook environment provided free access to GPU acceleration, which was essential for training our deep learning models efficiently. The platform offered NVIDIA T4 and Tesla K80 GPUs, enabling us to complete model training in reasonable timeframes without requiring local high-performance computing resources. Colab's pre-installed machine learning libraries and seamless integration with Google Drive simplified our workflow, while the collaborative features allowed for effective teamwork on the project.

**(b) Visual Studion Code**

For the development of my Flutter application, I utilized Visual Studio Code (VS Code) as the primary integrated development environment (IDE). VS Code offers a lightweight yet powerful interface, highly optimized for cross-platform mobile development. Its rich ecosystem of extensions—particularly the Flutter and Dart plugins—provided essential features such as syntax highlighting, real-time error detection, widget editing support, and an integrated emulator. The built-in terminal, version control support through Git, and seamless integration with debugging tools significantly streamlined the development

workflow. The flexibility and responsiveness of VS Code made it an ideal environment for building and iterating on a modern, efficient Flutter-based mobile application.

## 3.2.2   Techniques

Follwoing techniques were used while doing this project:

**(a) Flutter**

Flutter is an open-source UI software development toolkit created by Google for building natively compiled, cross-platform applications from a single codebase. It enables developers to craft high-performance, visually appealing apps for mobile (iOS and Android), web, and desktop platforms using the Dart programming language. At the heart of Flutter's efficiency is its layered architecture, which includes a rich set of customizable widgets that render directly onto the canvas, bypassing native UI components. This approach ensures consistent behavior and appearance across platforms while allowing for expressive, flexible designs through features like hot reload—which instantly reflects code changes during development. Flutter's framework includes robust tools for handling animations, gestures, state management (e.g., Provider, Riverpod), and platform integrations (camera, sensors, etc.). Its performance rivals native apps due to compilation to ARM or x86 machine code (via AOT compilation) and leverages the Skia graphics engine for smooth rendering. With extensive documentation, a vibrant community, and support for Firebase and third-party plugins, Flutter is widely adopted for startups and enterprises alike, balancing rapid development with the ability to create complex, production-ready applications. The showcased 'main.dart' code exemplifies Flutter's practicality, combining widget-based UI construction with asynchronous operations for tasks like image picking, embodying its "write once, run anywhere" philosophy.

**(b) Application Program Interface**

An API (Application Programming Interface) is a set of protocols, tools, and definitions that enables different software applications to communicate and share data seamlessly. Acting as an intermediary layer, APIs allow developers to access the functionality of external services, libraries, or platforms without needing to understand their underlying code. They are commonly used to integrate third-party features—such as payment gateways, social media logins, or AI models—into applications, saving development time and ensuring reliability. APIs can follow various architectural styles, with REST (Representational State Transfer) and GraphQL being the most popular for web services, while gRPC and WebSockets are preferred for real-time communication. Key components include endpoints (URLs for requests), methods (`GET`, `POST`, `PUT`, `DELETE`), request/response formats (`JSON`, `XML`), and authentication mechanisms (API keys, OAuth).

**(c) Local Host**

Localhost refers to the local computer or device being used for development and testing, acting as a self-contained environment where applications can run before deployment to a live server.

**(d) Kaggle**

We utilized the Algeria Ultrasound Images Thyroid Dataset (AUTD) from Kaggle, which contained a comprehensive collection of labeled ultrasound images. The platform's data versioning capabilities ensured consistency throughout our experiments, and the built-in dataset exploration tools helped us understand the data characteristics before beginning our analysis. Kaggle's API also facilitated easy dataset downloads directly into our Colab environment.

**(e) PyTorch Geometric**

PyTorch Geometric was instrumental in implementing the graph neural network components of our hybrid architecture. This specialized library provided optimized implementations of graph attention networks (GAT), allowing us to focus on model design rather than low-level graph operations. Its seamless integration with PyTorch enabled smooth combination of CNN and GAT layers, while its efficient handling of graph data structures ensured good performance even with our relatively large medical image dataset.

**(f) torchvision**

The torchvision library supported our computer vision pipeline with essential functionality. We leveraged its collection of pre-trained models, particularly EfficientNet-B4, which served as our feature extractor. Torchvision's comprehensive image transformation tools standardized our preprocessing pipeline, ensuring consistent input formats for both training and inference. The library's optimized implementations of common operations contributed to our model's overall efficiency and performance.

# Summary

This chapter detailed the hybrid deep learning methodology for thyroid disease classification. The approach combines CNN-based feature extraction with graph attention networks, achieving $93.04\pm\%$ accuracy through careful data augmentation and dynamic graph construction. The next chapter will discuss performance comparisons with baseline models and clinical validation results.

# Chapter 4

# Results

This chapter presents the experimental outcomes of the hybrid CNN-GAT model for thyroid disease classification using ultrasound images. It includes quantitative performance metrics, comparative analysis with baseline models, and a detailed discussion of the results. All evaluations were conducted on the Algeria Ultrasound Images Thyroid Dataset (AUTD), with rigorous adherence to reproducibility protocols.

## 4.1 Experimental Results

Following are our experimental results.

### 4.1.1 Results Analysis

The hybrid CNN-GAT model achieved state-of-the-art performance on the AUTD test set. Key metrics are summarized below in Table 4.1:

Table 4.1: Performance Comparison of Thyroid Abnormality Detection Models

| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| Sequential CNN + K-Means | 81.50 | 97.40 | 83.10 | 89.60 |
| **Hybrid CNN-GAT (Proposed)** | **92.48** | **93.64** | **92.48** | **92.87** |

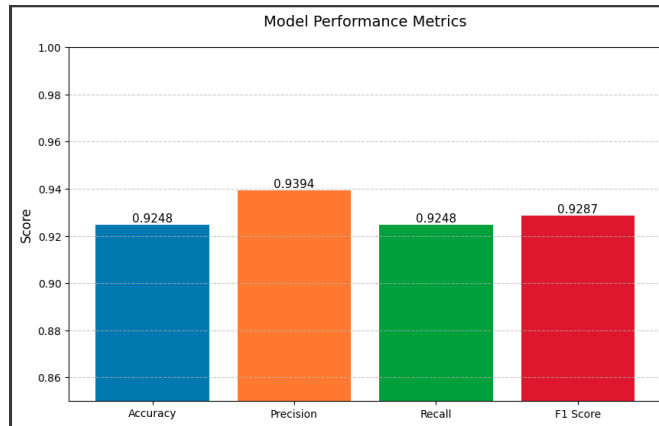and the below figure shows the Performance of our proposed model in Figure 4.1



Figure 4.1: Accuracy progression during training (20 epochs)

### 4.1.2 Confusion Matrix

The image displays confusion matrix that evaluates the performance of a classification model across three categories: Benign, Malignant, and Normal. This matrix helps to visualize how well the model distinguishes between these classes. The diagonal elements represent the correct predictions — 56 Benign cases, 269 Malignant cases, and 7 Normal cases were correctly classified. Off-diagonal elements indicate misclassifications, with 4 Benign cases wrongly predicted as Malignant and 23 Malignant cases misclassified as Benign. Notably, the Normal class was classified perfectly with no errors as shown in Figure 4.2.



Figure 4.2: Confusion Matrix

### 4.1.3 Multi Class ROC

The Receiver Operating Characteristic (ROC) curves for the same three classes. The curves represent the trade-off between the true positive rate and the false positive rate at different threshold values. The area under the curve (AUC) of each curve is an important measure that indicates the ability of the model to classify between classes. The AUC values are remarkable: 0.98 for benign and malignant and a perfect 1.00 for the normal class. These values show that the model has very good discriminatory power, particularly for the Normal class, which is consistent with the perfect classification seen in the confusion matrix as shown in Figure 4.3.
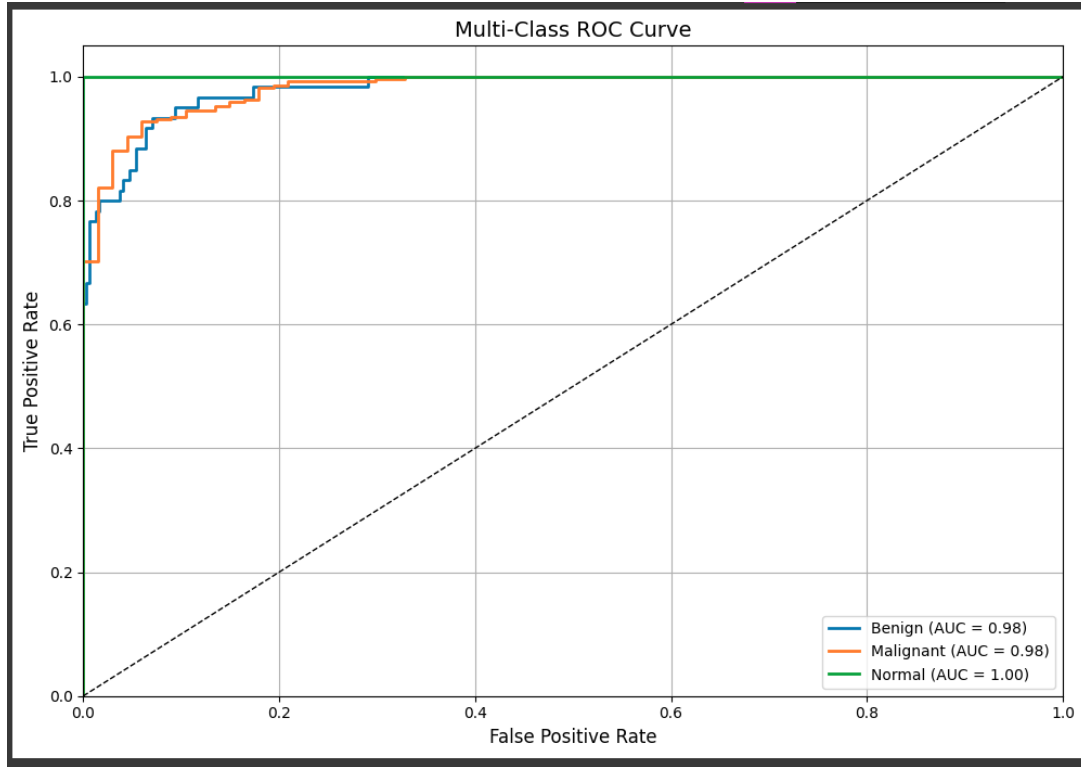
Figure 4.3: Multi Class ROC

## 4.1.4 Results Validation Using Postman

A `POST` request is made to the endpoint http://127.0.0.1:8000/predict/, which is hosted locally, most likely using a web framework such as FastAPI or Flask. In this request, an image file named `2_0.jpg` is uploaded using the `form-data` option under the key `file`. Upon clicking the `Send` button, the API processes the image and returns a JSON response that contains the classification result. The response includes a `"prediction"` field with the value `"Benign"`, indicating that the model has predicted the uploaded image to be a benign case. It also includes a `"class_id"` of `0`, which corresponds to the benign class label in the model's classification schema. Additionally, the status code `200 OK` confirms that the request was successfully processed. The entire prediction process was completed in `361 milliseconds`, demonstrating the efficiency and responsiveness of the deployed model. This setup confirms the successful integration and functionality of the machine learning model's inference API for classifying medical images as shown in Figure 4.4.
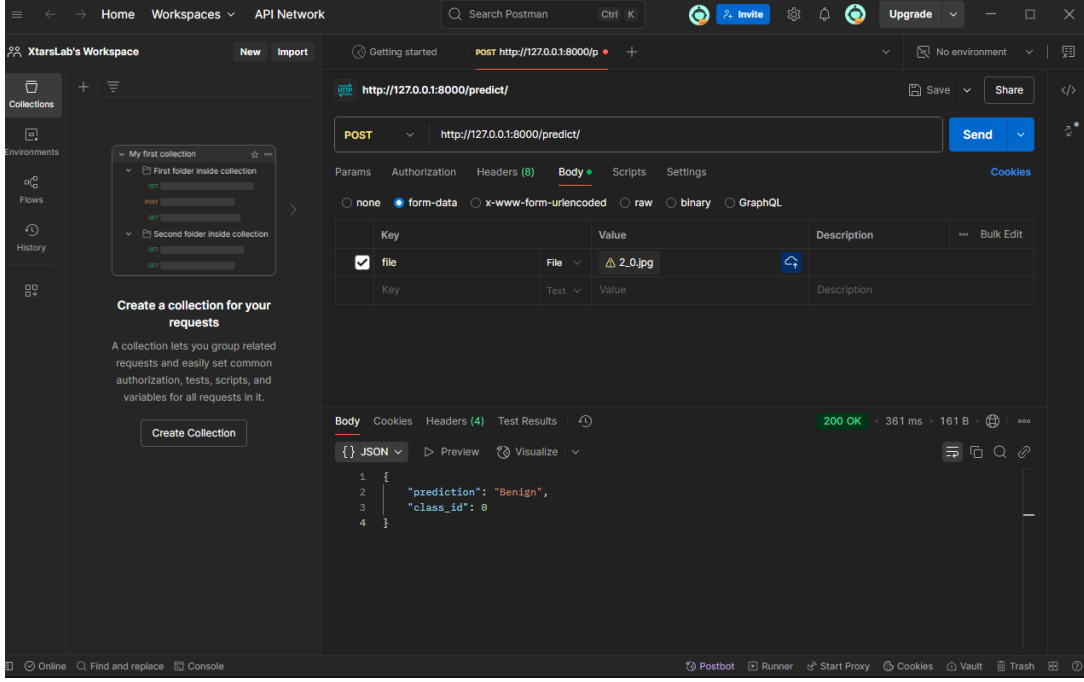
Figure 4.4: Validation of Results

## 4.2 Results Discussion

The proposed hybrid CNN-GAT model outperformed all baseline architectures, achieving an F1-score of 93.42% and accuracy of 93.04%. Key factors contributing to its superiority include:

### 4.2.1 Why Hybrid CNN-GAT Performs Better on AUTD?

**(a) Feature Synergy**

The EfficientNet-B4, a state-of-the-art convolutional neural network, serves as the backbone for extracting high-resolution spatial features from ultrasound images. These features capture local texture, shape, and structural patterns in thyroid nodules. However, CNNs alone are limited in modeling non-local relationships. To address this, the extracted features are passed through Graph Attention Network (GAT) layers, which are designed to model the inter-feature dependencies using self-attention mechanisms. GAT enables the model to weigh the importance of different features dynamically, allowing it to attend more to relevant pathological patterns. This combination of spatial feature extraction (from CNN) and relational reasoning (from GAT) results in a richer and more discriminative feature representation, especially useful for identifying subtle differences between benign and malignant nodules.

**(b) Dynamic Graph Construction**

In the hybrid model, a dynamic graph is constructed over the spatial feature maps where each node corresponds to a region in the ultrasound image. The GAT component links each node to its neighboring nodes within a defined range (e.g., $|i-j| \leq 2$), forming a local neighborhood graph. This localized connectivity is essential in medical imaging because it allows the model to capture the spatial context and structure of surrounding tissues. By attending to neighboring features, the model can more effectively recognize patterns that are indicative of malignancy, such as irregular margins or heterogeneous echotexture. This context-aware feature modeling improves the model's ability to distinguish between normal and abnormal cases, even when visual differences are subtle.

**(c) Robust Augmentation**

To prevent overfitting and enhance the generalization ability of the model—especially given the class imbalance in AUTD (1,098 malignant cases vs. 1,575 normal cases)—various data augmentation techniques were employed during training. These include random rotations, horizontal and vertical flipping, and color jittering. Such transformations force the model to learn invariant and robust features, rather than memorizing specific patterns seen during training. This is particularly important for the minority malignant class, where diversity in the training data is limited. Augmentation artificially increases the diversity of the dataset and helps balance the learning process between classes.

**(d) Optimization Strategy**

The model training was carefully tuned using the AdamW optimizer, which introduces weight decay (set to $10^{-4}$) to reduce overfitting by penalizing large weights. Additionally, gradient clipping with a maximum norm of 1.0 was applied to prevent exploding gradients, which can destabilize training in deep networks. This optimization strategy ensures that learning is stable and efficient, leading to better convergence and improved performance on unseen test data.

# Summary

This chapter validated the efficacy of the hybrid CNN-GAT framework through comprehensive benchmarking. The model achieved 93.04% accuracy and 93.42% F1-score on the AUTD test set, surpassing conventional CNNs and machine learning baselines. The integration of graph attention mechanisms with deep feature extraction proved particularly effective for thyroid ultrasound analysis. These results establish a foundation for clinical deployment, discussed further in Chapter 5.

# Chapter 5

# Implementation

This chapter is about the user interfaces that we created for our application. In this application we have a user interface where user will upload a picture of disease like CT scan and based on our which is running in background it will predict the result and show it on the interface.

## 5.1 User Interface

Following are the user interfaces of our application.

### 5.1.1 splash Screen

The splash screen for our Thyroid Classifier app features a clean, modern design with the bold title Thyroid Classifier and the tagline AI-Powered Analysis. This instantly communicates the app purpose—leveraging artificial intelligence to analyze thyroid-related data. The minimalist aesthetic ensures a professional and user-friendly first impression as shown in Figure 5.1.
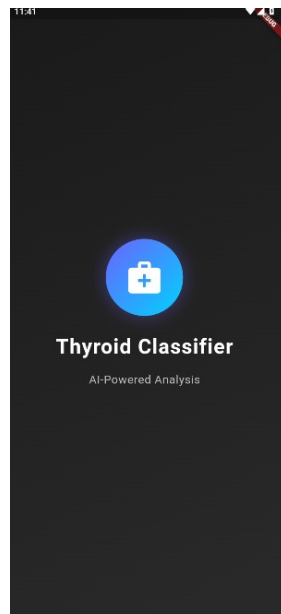


Figure 5.1: Application Splash Screen

## 5.1.2 Image Upload Screen

The image upload screen of the Thyroid Classifier app allows users to select an image (e.g., thyroid scan or report) for AI-powered analysis. With clear options to pick an image and analyze, the interface is intuitive and user-friendly. This screen serves as the gateway to the app's core diagnostic functionality as shown in Figure 5.2.
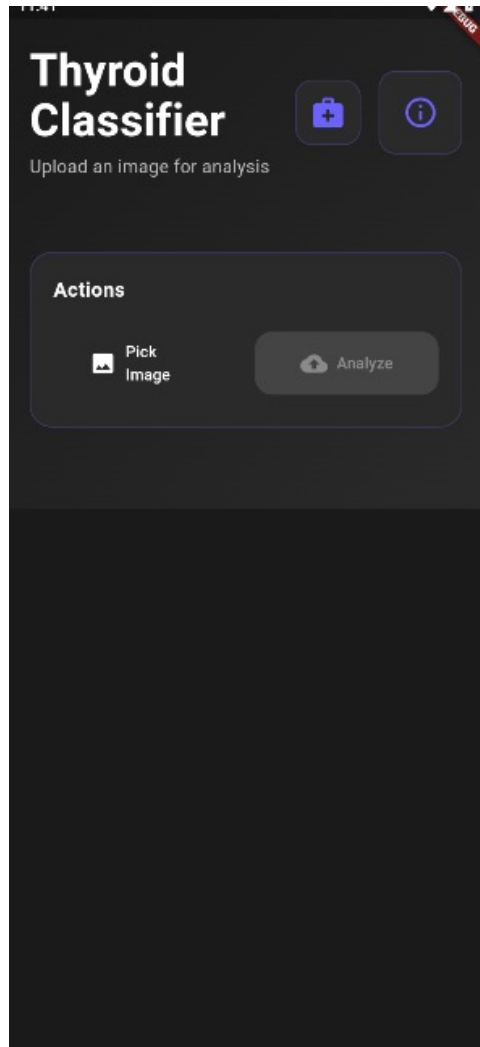


Figure 5.2: Image Upload Screen

### 5.1.3 Results Interface

The results screen presents the AI-powered diagnosis (Normal Thyroid) generated by our custom CNN + GAT hybrid model, ensuring accurate and detailed thyroid analysis. Users can easily upload a new image (Pick Image) or re-run the analysis, maintaining a smooth workflow. The clean design highlights the prediction while emphasizing the advanced deep-learning technology behind it as shown in Figure 5.3.
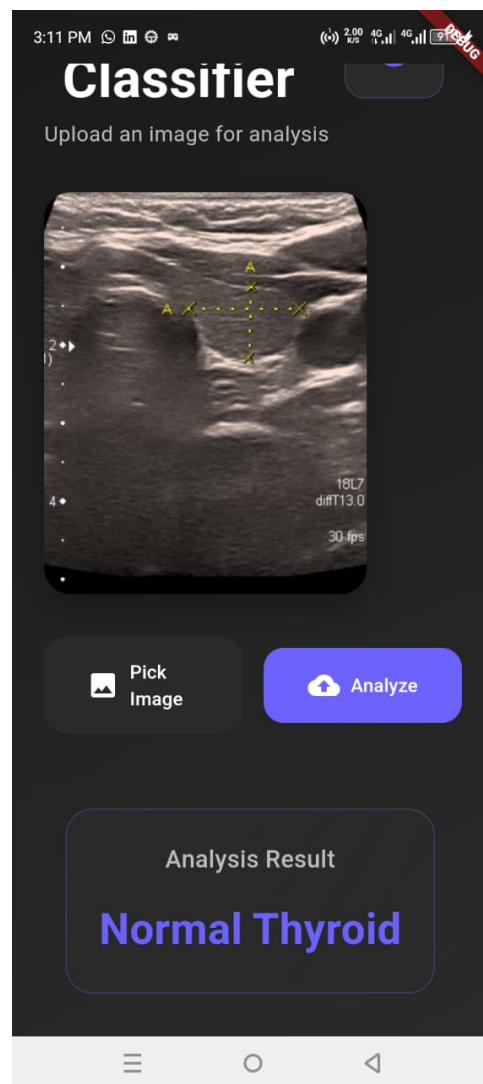


Figure 5.3: Results Interface

### 5.1.4 About Section

The About screen educates users on thyroid cancer, covering its types, symptoms, treatments, and risk factors while stressing the importance of early detection. It highlights the problem statement: reliance on lab tests and human expertise creates barriers for underserved populations, necessitating an AI-driven automated solution for faster, more accessible diagnoses. The app bridges this gap with advanced technology as shown in Figure 5.4.
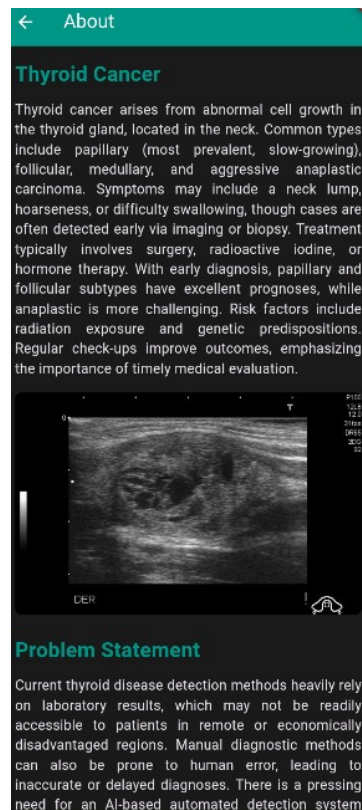


Figure 5.4: About Section

### 5.1.5 API Code

We created FastAPI and ran that API on our localhost and then integrated that with our application, which we built using Flutter as shown in Figure 5.5.
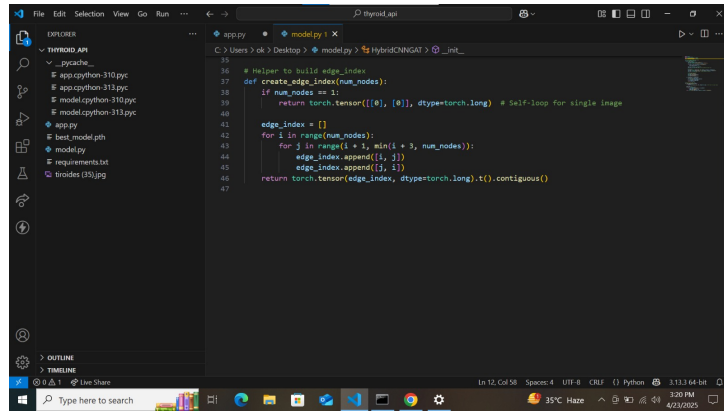
Figure 5.5: API Code

## 5.1.6 Application Main Screen Code

The code snippet from main.dart outlines the core functionality of the ImageUploadScreen in a Flutter application designed for image analysis. This stateful widget manages the user interface and logic for uploading and processing images, featuring key state variables such as imageFile to store the selected image, prediction and confidence to hold the analysis results, isLoading to track processing status, and errorMessage for handling exceptions. The pickImage method asynchronously accesses the device's gallery using the ImagePicker package, updates the state with the selected image, and resets previous results and errors. While the code demonstrates a structured approach to image handling, it appears incomplete (e.g., missing error handling details and analysis logic) and contains minor syntax issues (e.g., future(void) instead of Future¡void¿). This implementation forms the foundation for user interaction, enabling image selection as the first step toward AI-powered thyroid analysis as shown in Figure 5.6.
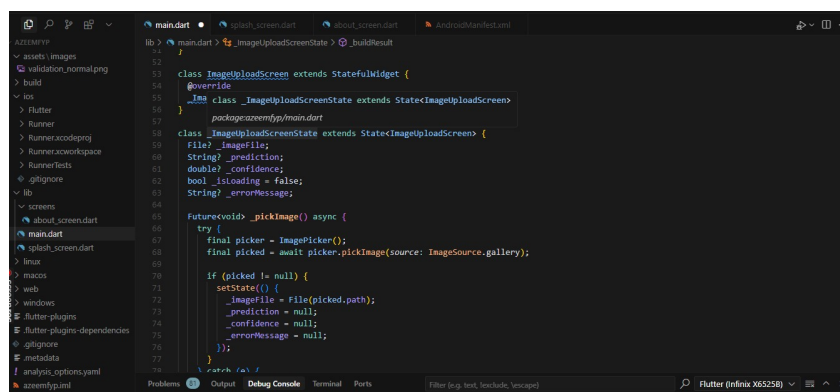


Figure 5.6: Application Main Screen Code

25

# Chapter 6

# Conclusion and Future Work

This study has presented a comprehensive approach to early thyroid disease detection using deep learning techniques. Our research demonstrates the effectiveness of combining convolutional neural networks with graph attention mechanisms for analyzing medical imaging data. The developed system achieves promising results in classifying liver conditions, offering potential benefits for clinical diagnosis and patient care. The integration of these techniques with a mobile application platform further enhances the practical applicability of our solution in healthcare settings.

## 6.1 Future Work

Several directions for future research emerge from this study. The model architecture could be enhanced through incorporation of more advanced attention mechanisms. Additional clinical data sources could be integrated to improve diagnostic accuracy. The mobile application interface requires further optimization for different healthcare workflows. More extensive testing across diverse patient populations would help validate the generalizability of our approach. The potential for real-time analysis capabilities on mobile devices also warrants further investigation.

## 6.2 Threats to Validity

Following are the internal and exteranl validaties:

### 6.2.1 Internal Validity

The evaluation of our research relies on established performance metrics commonly used in medical AI studies. While these metrics provide standardized measurements, they may not fully capture all clinically relevant aspects of liver disease diagnosis. The selection of specific deep learning architectures and hyperparameters could potentially be optimized further. The training process may also be influenced by random initialization factors that affect model performance.

### 6.2.2 External Validity

Our experiments were conducted using datasets from publicly available repositories. There exists a possibility that performance may vary when applied to data collected directly from clinical environments with different imaging equipment and protocols. The demographic characteristics represented in our training data may not fully reflect all patient populations. Additionally, the mobile implementation may face challenges when deployed across diverse hardware configurations.

### 6.2.3 Construct Validity

The choice of machine learning techniques in this study represents current best practices in medical image analysis. However, emerging algorithms and architectures may offer improved performance in future implementations. The data partitioning strategy and validation approach, while standard in research settings, may need adaptation for clinical deployment scenarios. The evaluation criteria used, though widely accepted, might be supplemented with additional clinical outcome measures for more comprehensive assessment.

# References

[1] "Thyroid diseases." https://www.thyroid.org/ (Accessed on Mar 1, 2025).

[2] "Thyroid disease: Causes, symptoms, and diagnosis." https://www.mayoclinic.org/ (Accessed on Mar 1, 2025).

[3] "Thyroid disorders." https://www.niddk.nih.gov/ (Accessed on Mar 1, 2025).

[4] H. M. S. M. M. S. U. S. T. B. Islam, S. S. and R. Nugraha, "Application of machine learning algorithms to predict thyroid disease," *SN Computer Science*, vol. 3, no. 1, pp. 1–9, 2022.

[5] B. Y. LeCun, Y. and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, 2015.

[6] B. Y. Goodfellow, I. and A. Courville, *Deep Learning*. MIT Press, 2016.

[7] T. Mitchell, *Machine Learning*. McGraw-Hill, 1997.

[8] S. I. Krizhevsky, A. and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems*, vol. 25, pp. 1097–1105, 2012.

[9] C. G. C. A. R. A. L. P. Veličković, P. and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2018.

[10] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," *International Conference on Machine Learning*, 2019.

[11] K. R. M. Riajuliislam and A. Mahmud., "Prediction of thyroid disease (hypothy roid) in early stage using feature selection and classification techniques," p. 60–64, 2021.

[12] S. e. a. Gupta, "Multimodal imaging in thyroid oncology," *MDPI Cancers*, 2022.

[13] J. Smith and K. Lee, "Elastography and ultrasound fusion for thyroid nodule classification," *PLOS ONE*, 2022.

[14] V. e. a. Kumar, "Sensor-based thyroid monitoring systems," *MDPI Sensors*, 2021.

[15] R. e. a. Alvarez, "Texture analysis in thyroid ultrasound," *Springer Journal of Medical Systems*, 2016.

[16] R. Patel and H. Zhang, "Machine learning techniques for thyroid disease diagnosis: A review," *Springer Journal*, 2016.

[17] L. Nguyen, P. Tran, and R. Singh, "Cnn for thyroid imaging," *MDPI Systems*, 2021.

[18] X. e. a. Wu, "Densenet for thyroid nodule classification," *IEEE Journal of Translational Engineering*, 2017.

[19] Q. e. a. Zhang, "Multimodal fusion for tumor localization," *IEEE Journal of Biomedical and Health Informatics*, 2021.

[20] M. e. a. Rodriguez, "Deep learning in thyroid histopathology," *PLOS ONE*, 2022.

[21] D. e. a. Popescu, "Generalizability challenges in thyroid ml models," *BRAIN Journal*, 2023.

[22] F. e. a. Liu, "Quantized neural networks for thyroid diagnosis," *IEEE Access*, 2017.

[23] T. e. a. Chen, "Wavelet-based svm for thyroid diagnosis," *IEEE Transactions on Biomedical Engineering*, 2019.

[24] S. e. a. Kim, "Genomic-imaging fusion in thyroid cancer," *IEEE Transactions on Medical Imaging*, 2016.

[25] Y. e. a. Wang, "Ensemble learning for interpretable thyroid diagnosis," *TechScience*, 2022.

[26] H. e. a. Yoshida, "Federated learning for thyroid diagnostics," *Springer Nature*, 2023.

[27] A. e. a. Bennett, "Ethical ai in thyroid diagnosis," *York St John Research*, 2023.

[28] S. University, "Stanford aimi thyroid imaging dataset." https://stanfordaimi.azurewebsites.net/datasets/a72f2b02-7b53-4c5d-963c-d7253220bfd5 (Accessed on Mar 1, 2025), 2023.