

**A Multimodal Vision-Language Framework for Pulmonary  
Embolism Detection: Integrating 3D CT Analysis with  
Medical LLM Report Generation**

**Sharon Chepkirui Tonui**

Submitted in partial fulfilment of the requirements for the degree of Masters  
of Data Science and Analytics, @iLabAfrica, Strathmore University



**Strathmore Institute of Mathematical Sciences**

**Strathmore University**

**Nairobi, Kenya**

**June 2025**

This dissertation is available for library use through open access on the understanding that it is copy-right material. No quotation from the dissertation may be published without proper acknowledgement.

# Declaration and Approval

## Declaration

I declare that this work has not been previously submitted and approved for the award of a degree by this or any other University. To the best of my knowledge and belief, the dissertation contains no material previously published or written by another person except where due reference is made in the dissertation itself.

© No part of this dissertation may be reproduced without the permission of the author and Strathmore University.

**Student's Name:** Sharon Chepkirui Tonui

**Sign:** 

**Date:** 27/05/2025

## Approval

The dissertation of Sharon Chepkirui Tonui was reviewed and approved for examination by the following:

**Dr. Kennedy Senagi,**  
Lecturer, @iLab Africa, Strathmore University

**Dr. John Olukuru,**  
Lecturer, @iLab Africa, Strathmore University

**Dr. Godfrey Madigu,**  
Dean, Institute of Mathematical Sciences, Strathmore University

**Prof. Bernard Shibwabo,**  
Director of Graduate Studies, Strathmore University

# Abstract

Pulmonary embolism (PE) is a serious and life-threatening condition caused by an artery blockage in the lung. The prevalence of PE varies significantly, ranging from 0.14% to as high as 61.5%, depending on the patient population and medical setting. Mortality rates for those who develop PE are also concerning, with between 40% and 69.5% of patients dying from the condition. Particularly alarming is the case-fatality rate following surgery, where approximately 60% of patients who develop PE do not survive. Computed tomography pulmonary angiography (CTPA) is the standard imaging method for PE detection. The scans generate hundreds of images which need to be manually reviewed, making it time-consuming and prone to overdiagnosis. In Kenya, this issue is further exacerbated by a significant shortage of radiologists with only 1 radiologist for every 270,000 people, far below the recommended 10-12 radiologists per 100,000 people. This deficit delays diagnosis and significantly raises the risk of adverse outcomes, including higher mortality rates for patients. Artificial intelligence (AI) offers a promising solution by automating the image analysis process, reducing diagnostic delays, and supporting radiologists in decision-making. This study presents a deep learning-based pipeline for automated PE detection and radiology report generation. Our system integrates a CT-ViT (Vision Transformer for 3D Medical Image Processing) model to extract features from CTPA scans, followed by the MedTron-7B Large Language Model (LLM), which translates extracted insights into structured radiology reports. Additionally, a Visual Question Answering (VQA) module enhances clinical interpretability by enabling contextual queries on detected abnormalities. Evaluation metrics indicate strong model performance in structured text generation, with peak ROUGE-1 recall reaching 1.0, while BLEU-1 and BLEU-4 scores of 0.35 and 0.22, respectively, highlight challenges in maintaining linguistic coherence. The results indicate the potential of AI-driven diagnostic tools in improving PE detection efficiency, reducing radiologist workload, and enhancing diagnostic accuracy.

**KEY WORDS:** Pulmonary Embolism, PE, CTPA, CTViT, Deep Learning, LLM, Radiology.

# Table of contents

<b>List of figures</b>	<b>ix</b>
<b>List of tables</b>	<b>x</b>
<b>List of Abbreviations</b>	<b>xi</b>
<b>Definition of Terms</b>	<b>xii</b>
<b>Acknowledgement</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background to the Study . . . . .	1
1.1.1 Diagnosis of Pulmonary Embolism . . . . .	2
1.1.2 Artificial Intelligence in PE Detection . . . . .	3
1.1.3 AI Adoption in Kenya . . . . .	4
1.2 Problem Statement . . . . .	5
1.3 Research Objectives . . . . .	6
1.3.1 General Objective . . . . .	6
1.3.2 Specific Objectives . . . . .	6
1.4 Research Questions . . . . .	6
1.5 Scope/Limitations . . . . .	7
1.5.1 Scope . . . . .	7
1.5.2 Limitations . . . . .	7
<b>2 Literature Review</b>	<b>8</b>
2.1 Introduction . . . . .	8

2.1.1	Computer Vision and Advancements in the Field . . . . .	8
2.1.2	Applications in Medical Imaging . . . . .	9
2.2	AI in PE Detection . . . . .	11
2.2.1	Existing Models for PE Detection . . . . .	11
2.2.2	Techniques to Enhance Explainability . . . . .	13
2.3	AI in Report Generation . . . . .	14
2.3.1	Automation of Radiologist Report Generation . . . . .	14
2.3.2	CLIP for Medical Imaging . . . . .	16
2.3.3	Advancements in 3D Medical Vision-Language Pretraining . . . . .	19
2.3.4	African Context and Regional Challenges . . . . .	21
2.4	Research Gaps . . . . .	22
2.5	Conceptual Framework . . . . .	23
<b>3</b>	<b>Methodology</b> . . . . .	<b>25</b>
3.1	Introduction . . . . .	25
3.2	Research Design . . . . .	25
3.3	Data Understanding . . . . .	26
3.3.1	Data Source . . . . .	26
3.3.2	Data Sampling . . . . .	27
3.3.3	Inclusion and Exclusion Criteria . . . . .	28
3.3.4	Data Format . . . . .	29
3.4	Data Preparation . . . . .	29
3.4.1	Image Preprocessing . . . . .	29
3.4.2	Text Preprocessing . . . . .	30
3.4.3	VQA Dataset Creation . . . . .	30
3.5	Modeling . . . . .	31
3.5.1	The CT-CLIP model . . . . .	31
3.5.2	Pre-training . . . . .	32
3.5.3	Fine-tuning . . . . .	33
3.6	Evaluation . . . . .	34

3.7	Deployment . . . . .	35
3.8	Stakeholder Engagement . . . . .	35
<b>4</b>	<b>System Design and Architecture</b>	<b>37</b>
4.1	Introduction . . . . .	37
4.2	System Requirements . . . . .	37
4.2.1	Functional Requirements . . . . .	37
4.2.2	Non-Functional Requirements . . . . .	38
4.3	Overview of System Architecture . . . . .	39
4.4	Frontend Development . . . . .	41
4.4.1	User Interface Design . . . . .	41
4.4.2	Image Upload . . . . .	41
4.4.3	Report Viewer . . . . .	42
4.4.4	Image Viewer . . . . .	42
4.5	Backend Development . . . . .	43
4.6	Deployment Using Docker on Lightning AI . . . . .	44
<b>5</b>	<b>System Implementation and Testing</b>	<b>46</b>
5.1	Introduction . . . . .	46
5.2	User Interface Design . . . . .	46
5.3	CT Scan Upload Process . . . . .	47
5.4	Upload Verification . . . . .	48
5.5	Report View . . . . .	49
5.6	Image View . . . . .	51
5.7	History Panel . . . . .	52
5.8	Testing . . . . .	53
5.8.1	Unit Testing . . . . .	54
5.8.2	Integration Testing . . . . .	54
5.8.3	Usability Testing . . . . .	54
5.8.4	Error Handling . . . . .	55

<b>6</b>	<b>Discussion of Results</b>	<b>56</b>
6.1	Introduction . . . . .	56
6.2	Data Preparation . . . . .	56
6.3	VQA Dataset Creation . . . . .	57
6.4	Image Pre-processing . . . . .	58
6.5	Text Pre-processing . . . . .	58
6.6	Vision Feature Extractor . . . . .	59
6.7	Model Training Pipeline . . . . .	59
6.8	Model Analysis and Performance Metrics . . . . .	62
6.8.1	Training and Validation Loss . . . . .	62
6.8.2	ROUGE-1 Precision and Recall Analysis . . . . .	63
6.8.3	ROUGE-L Precision and Recall . . . . .	65
6.8.4	BLEU Score Comprehensive Analysis . . . . .	66
6.8.5	Clinical Implications . . . . .	66
<b>7</b>	<b>Conclusions, Recommendations and Future Work</b>	<b>68</b>
7.1	Conclusion . . . . .	68
7.2	Recommendations . . . . .	69
7.3	Future Works . . . . .	69
	<b>References</b>	<b>71</b>
	<b>Appendix A Similarity Report</b>	<b>79</b>
	<b>Appendix B Ethical Clearance Confirmation</b>	<b>81</b>
	<b>Appendix C Model Development Code</b>	<b>82</b>
C.1	Loading Pretrained CT-CLIP Model . . . . .	82
C.2	Creating VQA Dataset . . . . .	83
C.3	Fine-Tuning using Meditron . . . . .	85
C.4	Model Code . . . . .	98
C.5	Frontend Code . . . . .	98



# List of figures

Figure 1.1: PE in the segmental artery anterior of the basal segment of the right lower lobe. (Key, n.d.) . . . . .	2
Figure 1.2: PE in the segmental artery posterior of the basal segment of the right lower lobe. (Key, n.d.) . . . . .	2
Figure 1.3: Contrast vs Non-contrast scan. (Key, n.d.) . . . . .	3
Figure 2.1: Conceptual Framework . . . . .	24
Figure 4.1: System Architecture Flow Chart . . . . .	40
Figure 5.1: Welcome screen of the system . . . . .	47
Figure 5.2: Upload results and verification . . . . .	48
Figure 5.3: Invalid format error message . . . . .	49
Figure 5.4: Generated report view . . . . .	50
Figure 5.5: Generated report view . . . . .	50
Figure 5.6: CT Scan Upload Interface . . . . .	52
Figure 5.7: Image viewing interface . . . . .	53
Figure 6.1: Model Training and Validation Loss Progression . . . . .	63
Figure 6.2: Adaptive Learning Rate Schedule . . . . .	63
Figure 6.3: ROUGE-1 Precision and Recall . . . . .	64
Figure 6.4: ROUGE-L Precision and Recall Metrics . . . . .	65
Figure 6.5: BLEU-1 and BLEU-4 Score Comparative Analysis . . . . .	66

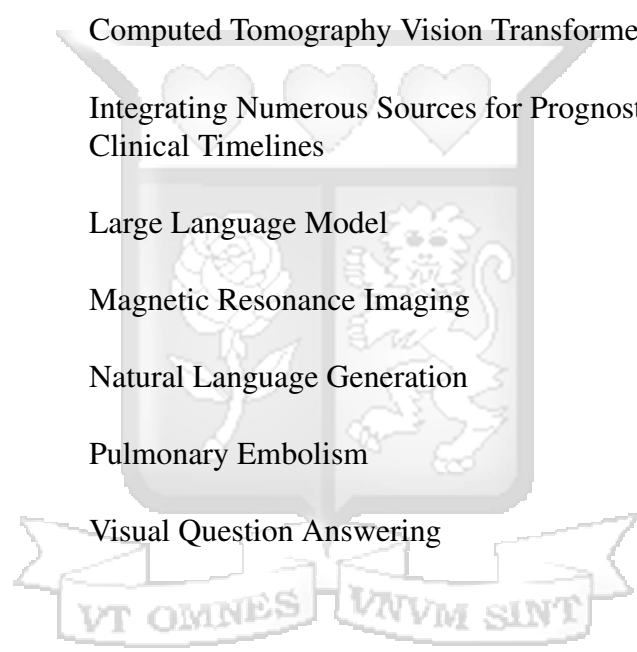
# List of tables

Table 6.1:	3D CT-ViT hyperparameters . . . . .	60
Table 6.2:	CT-CLIP hyperparameters . . . . .	60
Table 6.3:	LoRA-specific parameters for Meditron LLM . . . . .	60
Table 6.4:	General training hyperparameters and justifications . . . . .	61
Table 6.5:	Performance Metric Summary . . . . .	64



# List of Abbreviations

AI	Artificial Intelligence
CLIP	Contrastive Language-Image Pretraining
CT	Computed Tomography
CTPA	Computed Tomography Pulmonary Angiogram
CTViT	Computed Tomography Vision Transformer
INSPECT	Integrating Numerous Sources for Prognostic Evaluation of Clinical Timelines
LLM	Large Language Model
MRI	Magnetic Resonance Imaging
NLG	Natural Language Generation
PE	Pulmonary Embolism
VQA	Visual Question Answering



# Definition of Terms

<b>AI</b>	Computer systems designed to mimic human cognitive functions by employing complex algorithms and advanced computational techniques. AI enables machines to learn, analyze, adapt, and make decisions <a href="#">Russell and Norvig (2016)</a> .
<b>CT</b>	An advanced imaging method that produces 3D views of internal anatomy using multiple X-ray images and computer processing <a href="#">National Institute of Biomedical Imaging and Bioengineering (2023a)</a> .
<b>CTPA</b>	A specialized CT scan that visualizes pulmonary arteries using contrast material to detect pulmonary embolism <a href="#">Essien et al. (2019)</a> .
<b>CLIP</b>	An AI model by OpenAI that learns visual concepts from natural language to describe images accurately <a href="#">Radford et al. (2021)</a> .
<b>CT-ViT</b>	A deep learning model that encodes 3D CT volumes using spatial and causal attention to reconstruct the original image <a href="#">Hamamci et al. (2024b)</a> .
<b>INSPECT</b>	A multimodal dataset combining CT images, radiology reports, and EHR data to aid PE diagnosis and prognosis <a href="#">Huang et al. (2023)</a> .
<b>LLM</b>	AI models trained on massive text datasets to generate human-like language, such as OpenAI's GPT series <a href="#">Brown et al. (2020)</a> .
<b>MRI</b>	A diagnostic tool using magnetic fields and radio waves to create detailed images of soft tissues and organs <a href="#">National Institute of Biomedical Imaging and Bioengineering (2023b)</a> .
<b>NLG</b>	A subfield of AI focused on generating human-like text from structured data <a href="#">Siddharthan (2001)</a> .
<b>PE</b>	A blockage in the pulmonary artery, often due to a clot from the leg, which is potentially life-threatening <a href="#">Merck &amp; Co. (2023)</a> .
<b>VQA</b>	A research area in AI that builds models to answer questions about images by combining vision and language understanding <a href="#">Antol et al. (2015)</a> .

# Acknowledgement

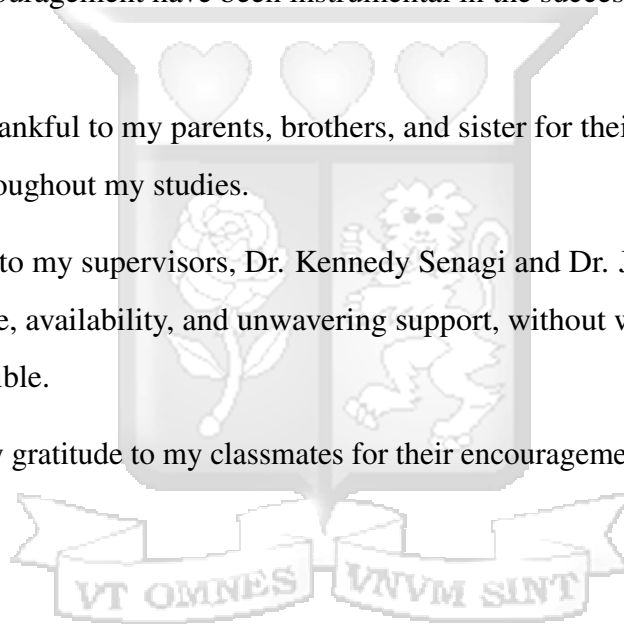
I acknowledge the invaluable support and contributions that made this research possible. First and above all, I am deeply grateful to the Almighty God for His provision, grace, and good health throughout my academic journey.

I extend my sincere appreciation to my sponsors, @iLabAfrica, for their unwavering belief in my abilities and commitment to academic excellence. Their generous financial support, resources, and encouragement have been instrumental in the successful completion of this research.

I am profoundly thankful to my parents, brothers, and sister for their constant support and encouragement throughout my studies.

I am also indebted to my supervisors, Dr. Kennedy Senagi and Dr. John Olukuru, for their invaluable guidance, availability, and unwavering support, without which this thesis would not have been possible.

Finally, I extend my gratitude to my classmates for their encouragement and support throughout this journey.



# Chapter 1

## Introduction

### 1.1 Background to the Study

Pulmonary embolism (PE) is a critical medical emergency characterized by blockage of the pulmonary arteries by blood clots, typically originating from deep vein thrombosis in the lower extremities. This condition significantly affects lung function, leading to a varied severity of symptoms and potentially fatal outcomes. PE represents a substantial global health burden, ranking as the third most common cause of cardiovascular death after stroke and myocardial infarction ([Raskob et al., 2014](#); [Wendelboe and Raskob, 2016](#)).

The global impact of PE is evident, with more than 300,000 deaths per year reported in the United States alone ([Raskob et al., 2014](#)). Worldwide, PE is associated with a 20% 90-day mortality rate ([Khan et al., 2022](#)), highlighting its severity and the urgent need for better management strategies. The situation in Africa is even worse, with mortality rates ranging between 40% and 69.5% ([Danwang et al., 2017](#)), significantly higher than global averages. This disparity highlights the unique challenges facing healthcare systems in African countries, including limited resources, delayed diagnosis, and inadequate treatment facilities.

In Kenya, a study conducted at Kenyatta National Hospital (KNH) revealed a mortality rate associated with PE of 28% among 128 patients ([Ogeng'o et al., 2011](#)). Although lower than the continental average, this figure still represents a significant health burden and emphasizes the critical need for improved diagnostic and treatment capabilities in the Kenyan healthcare system.

### 1.1.1 Diagnosis of Pulmonary Embolism

The diagnosis of PE is a complex process that integrates clinical assessment, risk stratification, and imaging studies. Clinical decision support rules, such as the Wells score or the Geneva score, are often used as initial screening tools to assess the probability of PE (Zantonelli et al., 2022). However, these rules are subjective and can be influenced by variability in human interpretation and laboratory output, potentially lowering the clinical pretest probability (Khan et al., 2022).

Computed Tomography Pulmonary Angiography (CTPA) has emerged as the gold standard imaging modality for PE diagnosis, demonstrating high sensitivity (94%) and specificity (98%) (Patel et al., 2020). CTPA provides detailed visualization of the pulmonary arterial tree, with contrast-enhanced scans being significantly more effective than non-contrast scans.



Figure 1.1: PE in the segmental artery anterior of the basal segment of the right lower lobe. (Key, n.d.)

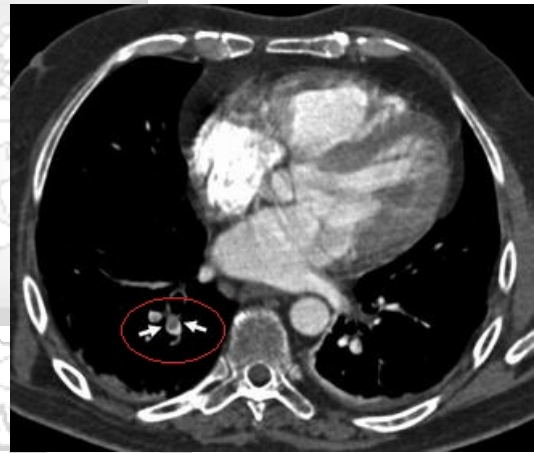


Figure 1.2: PE in the segmental artery posterior of the basal segment of the right lower lobe. (Key, n.d.)

Non-contrast CTPA is occasionally used to detect hyperdense thrombi, but its diagnostic utility is limited as it cannot reliably differentiate between vascular occlusion and normal blood flow (Kaykisiz et al., 2018). In contrast, contrast-enhanced CTPA uses iodinated contrast to opacify the pulmonary arteries, allowing emboli to appear as dark filling defects within the bright vessel lumen (Wittram et al., 2004). This enhances sensitivity and specificity for detecting PE while also providing insights into right heart strain and clot burden.

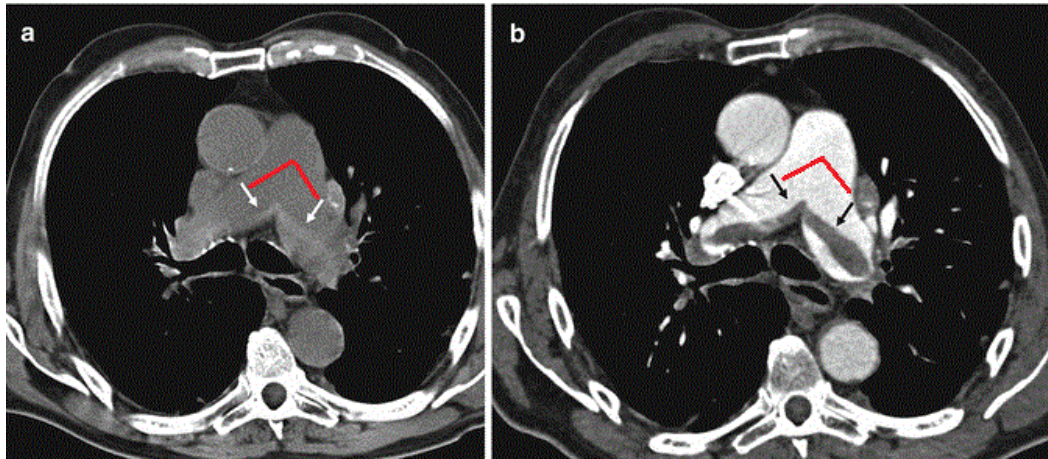


Figure 1.3: Contrast vs Non-contrast scan. (Key, n.d.)

As demonstrated in the provided images, a non-contrast scan (a) may reveal high-attenuation emboli, but a contrast-enhanced scan (b) confirms their presence with greater clarity.

### 1.1.2 Artificial Intelligence in PE Detection

The application of Artificial Intelligence (AI) in medical imaging, particularly in PE detection, has shown promising results. A systematic review and meta-analysis on AI-supported PE detection on CTPA revealed that AI systems achieved a yield rate of 14.6%, compared to a radiologist yield rate of 13.1%, with an overall PE prevalence of 15.8% (Cheikh et al., 2022). This suggests that AI systems can perform at least on par with human radiologists in detecting PE.

A study on AI in pulmonary embolism detection using 11,736 CTPA scans demonstrated major advancements in diagnostic accuracy and efficiency. AI reduced detection time 60-fold and achieved high performance metrics: 91.6% sensitivity, 99.7% specificity, 99.9% negative predictive value, and 80.9% positive predictive value. These results highlight AI's potential to enhance medical imaging interpretation, reduce errors, and accelerate diagnosis. (Topff et al., 2023).

These advancements in AI-enabled PE detection offer the potential to significantly improve the efficiency and accuracy of PE diagnosis. By automating the initial screening of CTPA

scans, AI systems can prioritize cases with a high likelihood of PE, allowing radiologists to focus their expertise on the most critical cases and reducing the time to diagnosis for urgent conditions.

### **1.1.3 AI Adoption in Kenya**

While AI-enabled PE detection is gaining momentum in Europe and America, its implementation in Africa, including Kenya, remains limited. This lag is primarily due to challenges in local healthcare data availability and infrastructure (Colak et al., 2021; Musa et al., 2023). The Kenyan healthcare system faces unique challenges that make the adoption of AI-assisted diagnostic tools both promising and challenging.

Kenya is grappling with an acute shortage of radiologists, with only 1 radiologist for every 1 million people, far below the recommended 10-12 radiologists per 100,000 people (Insights10, 2023). This shortage is particularly pronounced in rural and remote areas. Consequently, many rural hospitals with X-ray departments lack on-site radiologists, necessitating the practice of "teleradiology," where scan images are sent to referral and academic institutions for reporting. While this approach improves access to radiological expertise, it also leads to an overwhelming influx of scans at central institutions, straining the already limited radiologist capacity.

Additionally, many healthcare facilities in Kenya, particularly in rural areas, face significant resource constraints, including limited access to advanced imaging equipment and inadequate IT infrastructure necessary for implementing AI systems. Another potential challenge is that developing AI models require large, diverse, and high-quality datasets, but the limited availability of standardized digital medical imaging data in Kenya poses a significant challenge to the development of locally relevant AI models. These factors collectively hinder the widespread adoption of AI-assisted diagnostic tools in the Kenyan healthcare system.

## 1.2 Problem Statement

Despite its efficacy, the interpretation of CTPA scans presents several challenges. The process is time-consuming. A single study can comprise hundreds of images, each requiring a meticulous review by radiologists before writing a detailed report. This completely manual process is particularly problematic in resource-constrained settings with limited radiological expertise. Additionally, the high sensitivity of CTPA can lead to the detection of small, clinically insignificant emboli, increasing the risk of overdiagnosis and resulting in unnecessary anticoagulation therapy (Weikert et al., 2020). The increasing use of imaging in clinical practice has also led to growing workloads for radiologists, which can result in longer turnaround times for diagnoses and potentially higher error rates due to fatigue or time pressure.

In addition, the absence of a standardized and efficient workflow for interpreting and prioritizing CTPA scans leads to suboptimal resource utilization and potential delays in identifying and treating critical cases. Lastly, the limited implementation of advanced diagnostic tools is a significant barrier. The challenges related to infrastructure, data availability and adaptation to local healthcare contexts hinder the integration of AI technologies into the Kenyan healthcare system.

There is a pressing need for an innovative solution that can improve the accuracy and efficiency of PE detection from CTPA scans while alleviating the burden on radiologists. By developing and implementing an AI-assisted PE detection system, it may be possible to significantly improve the efficiency of radiological services in Kenya and ultimately contribute to better patient outcomes.

## 1.3 Research Objectives

### 1.3.1 General Objective

The main objective of this study is to develop and evaluate a machine learning-based system for the automated detection of PE from CTPA scans, with the goal of reducing radiologist workload and enhancing the efficiency of PE diagnosis and treatment.

### 1.3.2 Specific Objectives

- I. To review existing machine learning models in detecting pulmonary embolism (PE) from CTPA scans.
- II. To design a multimodal deep learning system integrating 3D Vision Transformers and large language models.
- III. To automate the generation of radiologist reports that summarize key findings
- IV. To deploy the system in a web application, making it available for use.

## 1.4 Research Questions

- I. How accurately can a machine learning model detect PE in CTPA scans compared to traditional radiologist evaluations?
- II. Can the automated system generate reliable and clinically useful radiologist reports that summarize key findings such as clot size, location, and severity?
- III. How does the machine learning model's sensitivity, specificity, and overall accuracy compare to current diagnostic methods for detecting PE in CTPA scans?
- IV. To what extent can the use of a machine learning-based system reduce overdiagnosis and diagnostic errors in the evaluation of CTPA scans for PE?

## 1.5 Scope/Limitations

### 1.5.1 Scope

This research focuses on developing and evaluating deep learning architectures PE detection in CTPA imaging. The project encompasses three primary components: implementing 3D CT Vision Transformers for volumetric analysis of CTPA scans, deploying CT-CLIP multimodal models for enhanced feature extraction, and integrating large language models for automated radiological report generation. The scope includes binary classification of PE presence or absence in CTPA studies, localization of embolic clots within pulmonary arterial segments, and generation of structured clinical reports summarizing diagnostic findings. The research evaluates model performance using Natural Language Generation metrics and validates findings against expert radiologist interpretations. The study specifically targets CTPA imaging modality and excludes other pulmonary imaging techniques such as ventilation-perfusion scintigraphy or chest X-rays. Model development utilizes the INSPECT dataset for training, with performance assessment conducted on both annotated validation sets and unannotated clinical cases to evaluate real-world generalizability.

### 1.5.2 Limitations



# Chapter 2

## Literature Review

### 2.1 Introduction

This literature review provides an overview of the models developed to date, the advancements in computer vision, and their application in medical imaging, particularly in PE detection using CTPA scans. A significant interest in applying AI to medical imaging emerged in the early 2000s as deep learning techniques, such as convolutional neural networks (CNNs), began demonstrating superior performance across various image analysis tasks. These breakthroughs have led to the formulation and application of various models, including traditional machine learning models, deep learning architectures, and hybrid approaches. This review examines key studies, compares the effectiveness of different approaches in the context of PE detection, and highlights their ability to provide accurate and reliable diagnostic results.

#### 2.1.1 Computer Vision and Advancements in the Field

Computer vision has undergone a remarkable transformation in recent years, driven largely by advancements in deep learning and neural networks. The field has evolved from simple edge detection and image segmentation to complex tasks such as object recognition, scene understanding, and even image generation. One of the most significant breakthroughs came with the introduction of Convolutional Neural Networks (CNNs). In 2012, [\(Krizhevsky et al., 2012\)](#) demonstrated the power of deep CNNs with their AlexNet architecture, which significantly outperformed previous methods on the ImageNet Large Scale Visual Recognition Challenge. This marked the beginning of the deep learning era in computer vision.

Since then, numerous architectures have been proposed, each pushing the boundaries of what's possible in computer vision. (Simonyan and Zisserman, 2014) introduced in 2014 demonstrated that depth is crucial for good performance. GoogLeNet (Inception) (Szegedy et al., 2015) developed in 2014 introduced the concept of inception modules, allowing for efficient computation and deeper networks. Residual Networks proposed in 2015 (ResNet) (He et al., 2016) allowed for the training of incredibly deep networks through the use of skip connections. DenseNet (Huang et al., 2017) introduced in 2017 proposed an architecture that connects each layer to every other layer in a feed-forward fashion, strengthening feature propagation and encouraging feature reuse. These advancements have not only improved accuracy but also efficiency, allowing for real-time applications on mobile devices and embedded systems. More recently, attention mechanisms and transformer architectures, originally developed for natural language processing, have been successfully applied to computer vision tasks. Vision Transformer (ViT) (Dosovitskiy et al., 2020), showed that a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. Another significant advancement is the development of self-supervised learning techniques. Methods like SimCLR (Chen et al., 2020) and MoCo (He et al., 2020) have shown that visual representations can be learned without requiring large amounts of labeled data, making this approach especially valuable in fields where labeled data is limited or costly to acquire.

### **2.1.2 Applications in Medical Imaging**

The advancements in computer vision have had a profound impact on medical imaging. These technologies are being applied across various modalities including X-rays, CT scans, MRI, and ultrasound. One of the earliest and most successful applications has been in the detection of diabetic retinopathy. In 2016, (Gulshan et al., 2016) demonstrated an algorithm that could detect diabetic retinopathy in retinal fundus photographs with sensitivity and specificity on par with human experts. This work showed the potential for AI to assist in screening programs, particularly in areas with limited access to ophthalmologists. In breast cancer screening, AI has shown promise in reducing false positives and false negatives. A

study by (McKinney et al., 2020) in 2020 showed that an AI system could reduce false positives by 5.7% and false negatives by 9.4% compared to human experts in mammography interpretation. For brain imaging, AI has been applied to tasks such as tumor segmentation, stroke detection, and Alzheimer's disease prediction. In 2018, (Chilamkurthy et al., 2018) developed a deep learning algorithm that could detect critical findings in head CT scans, including intracranial hemorrhage and midline shift, with high accuracy. Lung nodule detection and classification has been one of the most studied areas in CT imaging . In 2017, (Ardila et al., 2019) demonstrated a deep learning algorithm that could detect lung cancer on low-dose chest CT scans with a performance matching or exceeding that of expert radiologists. The system could predict the risk of lung cancer with an AUC of 0.94, potentially improving the efficiency and accuracy of lung cancer screening programs. In abdominal imaging, AI has been applied to tasks such as liver lesion detection and characterization. A study by (Yasaka et al., 2018) in 2018 showed that a deep learning model could differentiate between three types of liver masses (hepatocellular carcinoma, metastatic tumors, and hemangioma) with an accuracy of 84%. For neuroimaging, AI has been used for tasks such as intracranial hemorrhage detection, stroke diagnosis, and brain tumor segmentation. (Prevedello et al., 2017) developed a deep learning algorithm that could flag head CT scans with critical findings, potentially reducing the time to diagnosis for urgent conditions. In cardiac CT, AI has been applied to tasks such as coronary artery calcium scoring and plaque characterization. A study by (van Assen et al., 2019) in 2019 demonstrated an AI system that could automatically quantify total coronary artery calcium, potentially speeding up cardiovascular risk assessment. AI has also been applied to improve the quality of CT images. Deep learning-based image reconstruction techniques have shown promise in reducing radiation dose while maintaining or even improving image quality. A 2019 study by (Akagi et al., 2019) demonstrated that a deep learning reconstruction algorithm could reduce radiation dose by 92% while preserving diagnostic accuracy in chest CT.

## 2.2 AI in PE Detection

### 2.2.1 Existing Models for PE Detection

Early CNN models often focused on detecting emboli in the main pulmonary arteries, leading to suboptimal performance in identifying smaller embolisms in peripheral arteries (Tajbakhsh et al., 2015, 2019). These initial models often focused on analyzing pre-processed image features, requiring complex segmentation and vessel alignment steps that hindered their efficiency and generalizability. One of the first investigations to use CNNs for PE detection was carried out by (Tajbakhsh et al., 2019). Their model, which was trained on 121 CTPA scans, had an 83% sensitivity in finding individual emboli with a fixed false-positive rate. This ground-breaking study showed that CNNs had the ability to outperform conventional computer vision methods for PE identification. (Liu et al., 2020) assessed a deep learning convolutional neural network (U-Net) trained on 590 patients. Their work established strong correlations between the clot burden detected by the model and traditional scoring methods like the Qanadli and Mastora scores, highlighting the model's potential in enhancing diagnostic accuracy for PE.

Researchers shifted to examining full volumetric CTPA scans rather than just specific slices or areas in recognition of the fact that PE may appear in various places inside the pulmonary vasculature. Models were able to acquire a more thorough understanding of the pulmonary arteries and spot emboli that may be missed by examining separate areas thanks to this all-encompassing strategy. For instance, to automatically identify PE using volumetric CTPA scans, (Huang et al., 2020) created PENet, a 77-layer 3D CNN. PENet was first trained on a sizable video dataset (Kinetics-600) and subsequently adjusted using a CTPA dataset. This showed the value of transfer learning, which involves applying information from other domains to the study of medical images. PENet's AUROC on an internal test set was 0.84, and it was 0.85 on an external dataset, proving its resilience and generalizability. (Ajmera et al., 2022) proposed a 2D segmentation model using a U-Net architecture with an Xception encoder. Their study, involving 251 CTPA scans (55 positive for PE), achieved notable performance metrics with a sensitivity of 0.80, specificity of 0.74, and accuracy of 0.76 at the

scan level. The model's slice-level sensitivity reached 0.93, demonstrating its effectiveness in identifying emboli across CT slices. Building upon existing research, several studies have focused on developing more accurate and efficient models for PE detection. (Pu et al., 2023) presented an innovative approach that detects and segments PEs using deep learning without the need for manual outlining. Their algorithm employs computer vision techniques to identify high-confidence PE regions and subsequently trains a 3D Recurrent Residual U-Net model, achieving better performance metrics than traditional manual methods. The most recent notable development is the use of attention mechanisms in CNN architectures, inspired by human experts' ability to focus on both global and local image features when interpreting CTPA scans. (Bushra et al., 2024) created an Attention-Guided CNN (AG-CNN) for PE detection, which includes an attention mechanism to direct the model's attention to localized lesion areas. This study also used ensemble approaches to combine the results of various object detection models. They were able to obtain cutting-edge performance on the FUMPE dataset (Masoudi and Saadatmand-Tarzjan) by integrating models like YOLOv8, Faster R-CNN, and EfficientDet, achieving an AUROC of 0.927 and a sensitivity of 0.862. Furthermore, incorporating essential arterial context data during training emerged as a crucial strategy for improving the localization of small embolisms, addressing a key limitation of earlier models.

Alongside technological advancements, researchers have recognized the importance of incorporating human-centered design (HCD) principles in developing clinically viable PE detection tools. (Babione et al., 2020) emphasized the crucial role of HCD in creating user-friendly, efficient, and effective clinical decision support systems (CDSS) for PE diagnosis. They highlighted the need for iterative design and evaluation processes, involving clinicians and other stakeholders, to ensure seamless integration into existing workflows and address users' specific needs. While deep learning has shown significant promise in PE detection, translating these advancements into tangible clinical benefits remains an ongoing challenge. Researchers have begun exploring the integration of AI tools into radiologist workflows. For example, (Batra et al., 2023) demonstrated the potential of AI-driven worklist reprioritization to expedite the interpretation of CTPA scans positive for PE. The study found that AI-driven reprioritization significantly reduced report turnaround times (47.6 vs. 59.9 minutes) and wait

times for PE-positive cases (21.4 vs. 33.4 minutes). This decrease in wait time, measured from examination completion to report initiation, indicates that the AI tool effectively assisted radiologists in prioritizing critical cases, enabling faster diagnoses.

### **2.2.2 Techniques to Enhance Explainability**

([Huang et al., 2020](#)) utilized Class Activation Maps (CAMs) to visualize the areas of the CTPA scan that the model focused on when making its predictions. This technique highlights the regions contributing most significantly to the classification decision, providing insights into the model's reasoning. Similarly, [Bushra et al. \(2024\)](#) incorporated CAMs into their Attention-Guided Convolutional Neural Network (AG-CNN) framework. These CAMs provide "explainability" to the model by visually demonstrating the correspondence between the model's attention and the actual lesion areas in the CTPA images. ([Ajmera et al., 2022](#)) used a U-Net architecture not just for classification but also for generating segmentation masks of the PE. These masks visually depict the location and extent of the embolus within the CTPA slice, making the model's prediction readily interpretable for clinicians. They argue that this explainable solution helps minimize interpretation time and ensures peripheral emboli are not overlooked. ([Bushra et al., 2024](#)) adopted a heuristic strategy to refine their detection model's predictions. By integrating classifier outcomes with detection results, they aimed to reduce false positives while preserving sensitivity. This process involved setting a probability threshold (0.018) based on evaluations of the validation dataset to balance precision and sensitivity. This thresholding approach, though heuristic, introduces a level of transparency into the decision-making process.

## 2.3 AI in Report Generation

### 2.3.1 Automation of Radiologist Report Generation

The automation of radiology report generation is an emerging area in medical AI with the potential to significantly reduce radiologists' workload and enhance patient care. This complex task requires models to precisely analyze medical images and generate clear, informative reports. It integrates computer vision for image interpretation with natural language processing (NLP) for text generation.

Early research frequently took inspiration from image captioning architectures, which combined convolutional neural networks (CNNs) for extracting visual features and recurrent neural networks (RNNs) for generating text. (Shin et al., 2016) demonstrated one of the first end-to-end systems for chest X-ray report generation using a combination of CNNs and LSTMs. Their system could generate reports that were rated as comparable to human-written reports in terms of clinical accuracy and fluency. These non-hierarchical RNN models, like the basic CNN-RNN, (Donahue et al., 2015), and (Jing et al., 2018), often struggled to create the longer, more comprehensive reports needed in radiology. To address this, researchers developed hierarchical RNN models that utilized multiple RNNs, typically one to generate sentences and another for words within those sentences. These models frequently incorporated attention mechanisms, which allowed them to focus on relevant parts of both the image and the generated text. Examples include Co-ATTN (Jing et al., 2018), which uses a co-attention mechanism on visual features and predicted tags; (Yuan et al., 2019), which incorporates multiple image views (like front and side) for a more complete report; and (Zhang et al., 2020a), which uses knowledge graphs with information on chest X-ray findings to enhance report generation. A different approach involved retrieval-based models that aimed to retrieve existing reports based on how similar the input image was to reports in a database. These models, like CVSE (Yan et al., 2021) and (Syeda-Mahmood et al., 2020), learned visual-semantic embeddings to assess the similarity between the image and the stored reports. For CTPA specifically, (Zhang et al., 2020b) proposed a system that could automatically generate structured reports for PE studies. Their system not only detected PE

but also described its location and extent, generating text that closely mimicked the style of human radiologists.

The advent of transformers, particularly pre-trained language models like BERT, revolutionized natural language processing, leading to improvements in radiology report generation. Models like RTMIC (Xiong et al., 2019) and RM+MCLN (Leite, 2022) offer benefits over RNNs because they train faster and better utilize the parallel processing power of GPUs. (Alfarghaly et al., 2021) introduced CDGPT2, a conditioned transformer model based on GPT2, which leveraged visual and semantic features to generate complete reports. This model demonstrated superior performance in word-overlap metrics compared to earlier RNN-based methods. Researchers also recognized the importance of incorporating domain knowledge into these models. KAD, proposed by (Zhang et al., 2023) utilizes the Unified Medical Language System (UMLS) and RadGraph to inject medical knowledge into the pre-training process, achieving comparable performance to expert radiologists on disease recognition tasks. Addressing the limitations of traditional retrieval methods based on cosine similarity, (Jeong et al., 2023) presented X-REM, a retrieval-based report generation module that leverages a multimodal encoder trained with supervised contrastive learning to better capture the nuanced relationships between chest X-ray images and radiology reports. This approach significantly improved the accuracy of retrieved reports. The challenge of limited data in medical imaging has also been tackled. (Windsor et al., 2023) explored various techniques for improving vision-language modeling in low-data settings, including domain adaptation of pre-trained models, unimodal self-supervision, and data augmentation. Their findings emphasized the importance of tailoring pre-trained models to the medical domain and leveraging additional supervisory signals to enhance performance in data-constrained scenarios. Another key area of development has been the optimization of latent space geometry in medical vision-language models. M-FLAG, introduced by (Liu et al., 2023), utilizes a frozen language model and an orthogonality loss to harmonize the latent space, leading to significant improvements in downstream tasks like image classification, segmentation, and object detection, while also reducing the number of trainable parameters.

The advancements discussed here demonstrate a clear progression towards more sophisticated, domain-aware, and data-efficient models for automated radiology report generation. However, there is still need to incorporate larger and more diverse datasets, develop robust evaluation metrics beyond word-overlap, and address the ethical and safety concerns associated with clinical deployment.

### 2.3.2 CLIP for Medical Imaging

Contrastive Language-Image Pre-training (CLIP) (Radford et al., 2021) is a powerful technique that learns visual representations from natural language supervision. It works by jointly training a vision encoder and a text encoder to align image and text representations in a shared latent space. This allows CLIP to understand the semantic relationships between images and their corresponding textual descriptions, enabling it to perform various tasks like image classification, image retrieval, and even zero-shot learning.

The application of Contrastive Language-Image Pre-Training (CLIP) in medical imaging has experienced rapid growth in recent years, spurred by its potential to bridge the gap between visual features and human language (Zhao et al., 2024). Early research focused on adapting the CLIP pre-training process to suit the specific characteristics of medical data (Zhao et al., 2024). A key challenge was the limited availability of large-scale annotated medical image-text datasets. Researchers explored using publicly available datasets encompassing diverse modalities like X-ray, CT, MRI, and histology. One notable early study, GLoRIA (Huang et al., 2021), demonstrated the significance of multi-scale feature alignment in CLIP pre-training. By contrasting attention-weighted image regions with corresponding words in paired reports, GLoRIA effectively learned both global and local image-text representations (Zhao et al., 2024)

To overcome the limitations posed by data scarcity, researchers investigated techniques to enhance data efficiency. These techniques included sentence-level text augmentation, multi-stage feature alignment, and clustering based on disease-level semantic information (Zhao et al., 2024). The incorporation of external medical knowledge into CLIP pre-training

was another significant development. Studies like UniBrain (Lei et al., 2023) and KoBo effectively integrated knowledge sources like MedKEBERT and CompGCN, leading to improved performance in tasks like classification and segmentation (Zhang et al., 2020a).

The potential of CLIP in radiology report generation was quickly recognized, as this task presented unique challenges for traditional encoder-decoder architectures. These architectures often struggled to generate reports that accurately reflected the complexity and clinical nuances found in real-world radiology reports (Endo et al., 2021). Researchers proposed innovative retrieval-based approaches to overcome these challenges. One such approach, CXR-RePaiR (Endo et al., 2021), leveraged a pre-trained CLIP model to select the most relevant text snippets from a large corpus of existing reports, effectively circumventing the need to generate text from scratch. CXR-RePaiR achieved superior performance compared to previous methods like R2Gen (Chen et al., 2020) and (Miura et al., 2021) M2 Trans on both internal and external datasets, demonstrating the effectiveness of this retrieval-based paradigm (Endo et al., 2021).

The success of these initial applications led to the development of specialized contrastive learning frameworks designed specifically for medical image-text representation learning. MedCLIP (Wang et al., 2022), for instance, tackled the challenge of limited medical image-text datasets by decoupling images and texts during pre-training. This innovative approach allowed for a combinatorial increase in the number of usable image-text pairs, thereby boosting model generalizability (Wang et al., 2022). MedCLIP further introduced a semantic matching loss function that incorporated medical knowledge, reducing false negatives during training and leading to more robust image representations (Wang et al., 2022). Another framework, PhenotypeCLIP, aimed to learn fine-grained image representations by contrasting images with individual sentences within reports, capturing the nuances of image-text relationships and improving performance in tasks like report generation.

The emphasis on retrieval-based approaches, incorporation of medical knowledge, and exploration of fine-grained representations are indicative of the evolving landscape of CLIP in medical imaging. Despite these advances, the application of CLIP to volumetric imaging modalities like CT scans remains challenging (Zhao et al., 2024). The computational demands

of processing large 3D volumes, the limited availability of large-scale annotated volumetric datasets, and the complexity of establishing semantic correspondences between sentence-level textual features and the information contained in volumetric data pose significant hurdles. The research indicates a trend towards developing more specialized and refined CLIP-based techniques to handle the complexities of volumetric medical imaging data.

### **(a) Limitations**

The foundational CLIP model, while revolutionary for natural image-text understanding, demonstrates severe limitations when applied to medical imaging. CLIP's global-level contrastive learning approach fails to capture the fine-grained, localized features critical for medical diagnosis. Performance evaluations show CLIP achieving only 0.52-0.58 AUC on medical imaging tasks compared to 0.85+ on natural images (Zhao et al., 2024). For 3D medical volumes, CLIP's performance degrades further to 0.45-0.51 AUC due to its inability to process volumetric relationships.

While CXR-RePaiR demonstrates superior performance over generative approaches, several critical constraints emerge for CTPA analysis. The model's performance is fundamentally limited by the quality and diversity of its retrieval corpus, creating bottlenecks when disease presentations differ from the training database. Additionally, unlike generative models that can adapt through fine-tuning, CXR-RePaiR's static retrieval database offers limited flexibility for local clinical practices.

MedCLIP, despite its innovative approach to decoupling images and texts during pre-training, suffers from factual inconsistency and hallucination issues (Wang et al., 2022). The semantic matching loss function, while reducing false negatives during training, cannot fully mitigate the model's tendency to generate plausible but clinically inaccurate associations between images and text. This limitation is particularly concerning in diagnostic applications where factual accuracy is paramount.

GLoRIA's attention-based framework for learning global and local representations improves upon CLIP's global-only approach but introduces computational complexity that limits its

deployment in resource-constrained environments (Huang et al., 2020). The multi-scale feature alignment process requires significant GPU resources, making it impractical for many healthcare facilities in low- and middle-income countries (LMICs). Additionally, GLoRIA's performance heavily depends on the quality and diversity of the training data, which poses challenges when considering the limited availability of annotated medical imaging datasets from diverse populations (Zhao et al., 2024).

### 2.3.3 Advancements in 3D Medical Vision-Language Pretraining

Earlier Medical Vision-Language Pretraining (Med-VLP) approaches mainly focused on 2D medical images, such as chest X-rays, due to limited data availability. However, 3D medical imaging modalities like CT and MRI scans offer richer anatomical details, making them essential for diverse medical applications (Lin et al., 2024; Bai et al., 2023). The transition from 2D to 3D image analysis in Med-VLP introduces challenges in aligning textual descriptions with the inherently sparse representations of 3D images.

CT-GLIP utilizes a grounded learning approach by constructing organ-level image-text pairs, mitigating the complexity of aligning textual descriptions with sparse 3D visual representations (Bai et al., 2024; Lin et al., 2024). This approach breaks down the complex task of full-image alignment into smaller, more manageable units. For instance, instead of aligning a whole CT scan with a comprehensive report, CT-GLIP aligns individual organ images with descriptions specific to those organs. This simplification enables more effective association of visual features with precise diagnostic text, enhancing the model's understanding of both normal and abnormal findings.

M3D-LaMed addresses the computational challenges of 3D medical image analysis by incorporating a specialized 3D spatial pooling perceiver. This perceiver effectively reduces the dimensionality of 3D image embeddings, allowing them to be processed efficiently by the LLM without sacrificing crucial spatial information (Bai et al., 2024). This is particularly important for tasks like segmentation and positioning, where accurate spatial understanding is crucial.

Both CT-GLIP and M3D-LaMed leverage large language models (LLMs) to enhance their understanding and generation capabilities. CT-GLIP uses BioClinicalBERT, a medical expert language model, to generate text embeddings, while M3D-LaMed employs the powerful LLaMA-2-7B model (Bai et al., 2024; Lin et al., 2024). The integration of LLMs allows these models to go beyond simple image-text matching and perform complex tasks like report generation, visual question answering, and referring expression segmentation.

M3D-LaMed introduces a novel approach to referring expression segmentation in 3D medical images by incorporating a promptable segmentation module. This module leverages the LLM's ability to understand natural language descriptions and translate them into prompts that guide a 3D segmentation model (SegVol) (Bai et al., 2024). This allows M3D-LaMed to perform tasks like segmenting a specific organ or lesion based on a textual description, demonstrating the potential of combining LLMs with specialized vision models for complex medical tasks.

Building on the available research, (Hamamci et al., 2024a) introduced CT-RATE, a novel, open-source dataset pairing 3D chest CT scans with corresponding radiology reports, addressing the scarcity of such data in 3D medical imaging. Leveraging CT-RATE, the authors developed CT-CLIP, a self-supervised foundation model for chest CT analysis that outperformed state-of-the-art supervised methods in multi-abnormality detection without requiring manual annotation. CT-CLIP demonstrated versatility through its applications in zero-shot multi-abnormality detection and case retrieval using both image and text queries.

A recent study by (Chen et al., 2024) introduced 3D-CT-GPT, a medical vision-language model designed for generating radiology reports from 3D CT scans, with a focus on chest CTs. This model employs a Visual Question Answering (VQA) framework and integrates a CT ViT for feature extraction, a 3D Average Pooling layer, and a projection layer to improve report generation from 3D CT scans (Chen et al., 2024). The study found that 3D-CT-GPT outperformed existing methods, including CT2Rep (Hamamci et al., 2024b), RadFM (Wu et al., 2023), and M3D (Bai et al., 2024), in both report accuracy and quality. Notably, 3D-CT-GPT achieved higher evaluation scores than M3D on both private and public datasets (Chen et al., 2024).

While 3D-CT-GPT demonstrated superior performance compared to existing methods, a significant limitation for researchers is that the model and its implementation were not made publicly available by the authors (Chen et al., 2024). This lack of public accessibility necessitated the exploration of alternative approaches for this study. CT-CLIP (Hamamci et al., 2024b), being an open-source foundation model specifically designed for chest CT analysis, presented a viable alternative. Although CT-CLIP was originally trained as a classifier rather than a report generator, its self-supervised learning approach and demonstrated versatility in multi-abnormality detection made it suitable for fine-tuning toward report generation tasks. The decision to utilize CT-CLIP over 3D-CT-GPT was therefore pragmatic, balancing the need for state-of-the-art performance with the practical constraints of model availability and accessibility for research purposes.

#### **2.3.4 African Context and Regional Challenges**

A qualitative study of radiographers' perspectives across Africa revealed significant concerns about AI integration, including job security fears and knowledge gaps regarding AI systems (Akudjedu et al., 2023). The study highlighted that while African radiographers recognize the potential benefits of AI in addressing workforce shortages, they also express apprehension about the lack of AI-related training programs and infrastructure support necessary for successful implementation. The radiological equipment landscape in Africa presents another critical challenge. (Mollura et al., 2020) conducted the first detailed analysis of registered diagnostic radiology equipment in a low-income African country, revealing severe shortages of basic imaging equipment, let alone the advanced computational infrastructure required for deploying complex AI models like MedCLIP and 3D-CT-GPT. This infrastructure gap is further compounded by inconsistent power supply, limited internet connectivity, and inadequate data storage capabilities in many healthcare facilities across the continent. Data sovereignty and representation issues also pose significant challenges. Medical imaging datasets used to train models like CT-CLIP predominantly originate from high-income countries with predominantly Caucasian populations, raising concerns about their applicability to African patient populations with different disease presentations, body habitus, and comorbidity profiles

(Akudjedu et al., 2023). The lack of diverse training data can lead to algorithmic bias and reduced diagnostic accuracy when these models are applied in African healthcare settings.

African populations present distinct epidemiological patterns affecting PE presentation. Studies from Kenyatta National Hospital (Ogeng'o et al., 2011) show PE is more frequently associated with HIV-related complications (23% vs. 3% in Western populations) and tuberculosis co-infections (31% vs. 2%), conditions rarely represented in Western training datasets. These demographic gaps contribute to 15-25% accuracy drops when Western-trained models are applied to Kenyan patients.

## 2.4 Research Gaps

The integration of AI into medical imaging has significantly advanced diagnostic capabilities. However, despite these strides, several limitations remain evident in the literature. One prominent limitation in the existing literature is the predominant focus on two-dimensional (2D) imaging modalities for automated report generation. Studies have largely concentrated on chest X-rays or selected slices from CT scans, neglecting the potential of volumetric data in 3D imaging (Hamamci et al., 2024a). This focus can be attributed to the relative simplicity of processing 2D images and the widespread availability of annotated 2D datasets. In contrast, 3D imaging poses significant challenges due to its higher computational demands and the lack of annotated datasets. A single 3D CT scan comprises hundreds of slices, requiring extensive memory and computational power, which can hinder the development and scalability of AI models.

Another limitation lies in the underexploration of VQA systems for 3D CT scans. VQA has shown potential in enabling interactive diagnostic tools, where users can query models for specific findings in medical images. However, current research primarily focuses on VQA systems for 2D imaging modalities. These systems provide answers to user queries but fall short of generating detailed, preliminary reports based on the complex spatial relationships present in 3D scans.

A third limitation is the restricted applicability of current models to contrast-enhanced imaging. While some studies have focused on report generation for 3D chest CT scans, these efforts have primarily utilized non-contrast chest CT volumes. The physiological and radiological differences between non-contrast and contrast-enhanced imaging significantly impact model performance. CTPA scans, for example, involve the administration of contrast dye to visualize pulmonary arteries and detect conditions such as PE. Models developed on non-contrast CT data are unlikely to generalize effectively to contrast-enhanced imaging, as they may fail to detect vascular structures and subtle abnormalities unique to CTPA scans.

This study aims to contribute to the research on report generation from 3D CT scan by focusing on developing a model specifically designed for contrast-enhanced CTPA scans. By tackling this critical gap, this research seeks to advance the field of AI in medical imaging, ensuring that AI systems can better meet the clinical demands of volumetric imaging and contribute more effectively to patient care.

## 2.5 Conceptual Framework

This conceptual framework presents an AI-driven approach to automating radiology report generation for pulmonary embolism (PE) detection using a CTPA-CLIP model.

Inspired by contrastive learning, the model establishes a joint embedding space between CTPA scans and their corresponding radiology reports. The process begins with feature extraction using two deep learning transformers: a 3D visual transformer processes volumetric chest CT images, while a text transformer encodes textual reports into semantic representations. These embeddings are aligned using a contrastive loss function, ensuring that each CTPA scan is paired with its correct report while being distinguished from unrelated ones. This training strategy enhances the model's ability to retrieve the most relevant report or generate structured findings when presented with a new scan.

Once trained, the CTPA-CLIP model will be deployed to assist radiologists in detecting PE and generating structured reports. Given an unseen CTPA scan, the model retrieves

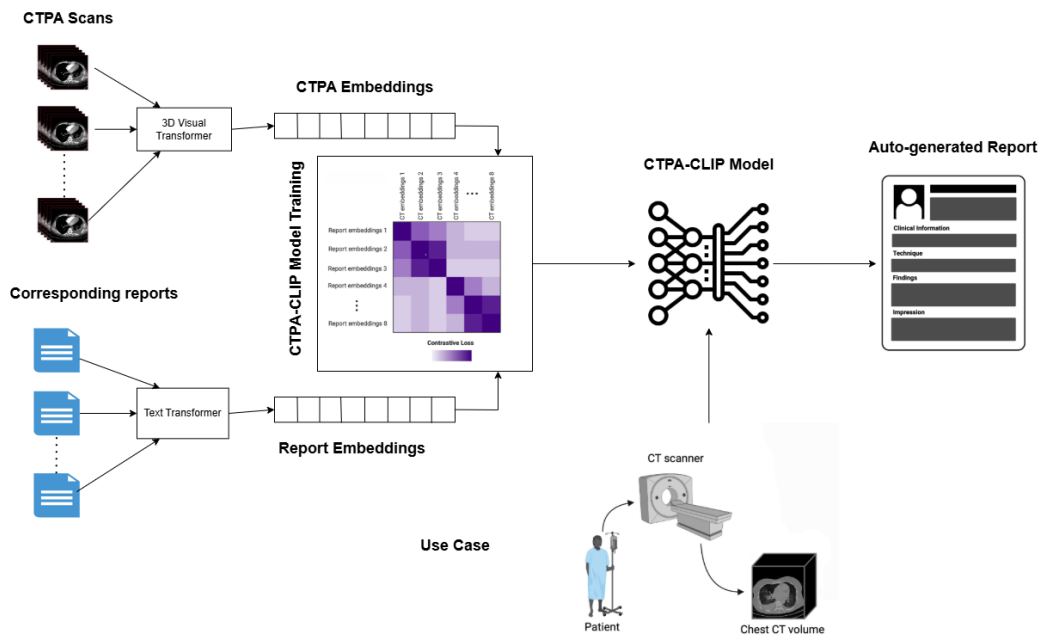


Figure 2.1: Conceptual Framework

a corresponding report or generates one. This automation aims to reduce the diagnostic workload, particularly in high-volume or resource-constrained settings, where radiologists face significant time pressure. The ability to provide first-pass AI-generated reports can expedite the review process, ensuring that radiologists focus on verification and clinical decision-making rather than repetitive documentation. Additionally, the model enhances standardization in radiology reporting, minimizing variability across practitioners and improving diagnostic consistency.

# Chapter 3

## Methodology

### 3.1 Introduction

This chapter outlines the methodology used to analyze CTPA scans, focusing on data understanding, preparation, and modeling techniques. The methodology is guided by the Cross Industry Standard Process for Data Mining (CRISP-DM) framework. The CRISP-DM framework provides a structured approach to data mining and knowledge discovery. The chapter details the data preparation techniques such as windowing and segmentation, and the modeling approaches for PE detection, including state-of-the-art architectures and multi-modal learning strategies. Finally, we evaluate model performance in report generation using NLG metrics.

### 3.2 Research Design

In addressing the critical need for rapid and accurate detection of pulmonary embolism (PE), this study employs a quantitative, experimental, and comparative research design. This approach is well-suited for evaluating the performance of deep learning (DL) models in medical imaging, as it allows for objective measurement and statistical analysis of diagnostic accuracy, sensitivity, and specificity. For instance, a systematic review and meta-analysis by Roberts et al. demonstrated the effectiveness of quantitative methods in assessing DL algorithms across various medical imaging modalities.

The retrospective component leverages existing CTPA scan datasets for model training and evaluation, an approach that facilitates the development of diagnostic tools without the need

for time-consuming prospective data collection. For instance, [Janizek et al. \(2020\)](#) trained a DL model on retrospective chest X-ray images to assess its generalization across different centers and imaging devices. The experimental aspect involves testing the performance of the developed DL models on new, unseen data to evaluate their generalizability and robustness across diverse clinical settings. [Sengun et al. \(2021\)](#) exemplified this by assessing various DL architectures for liver segmentation, implementing them into a commercial software, and testing their performance on unseen data to determine the most effective architecture. The comparative component involves comparing the performance of the DL models with radiologist interpretations, a crucial step in assessing the potential integration of AI tools into clinical workflows and their impact on diagnostic accuracy. [Aggarwal et al. \(2021\)](#) highlighted the importance of such comparative studies, demonstrating that DL models can achieve diagnostic performance comparable to that of radiologists in certain medical imaging tasks.

### **3.3 Data Understanding**

#### **3.3.1 Data Source**

INSPECT ([Huang et al., 2023](#)) is the largest multimodal dataset from a large cohort of PE patients, along with ground truth labels for multiple outcomes. It was developed to enable reproducible research on strategies for integrating 3D medical imaging and EHR data. The dataset includes CTPA images, structured electronic health record (EHR) data such as demographics, diagnoses, procedures, and vitals, as well as radiology report impressions with ground truth labels for multiple outcomes. The data was collected in a longitudinal study involving 19,438 PE patients. The images are stored in the NIfTI (.nii) formats.

### 3.3.2 Data Sampling

Given the vast size of the INSPECT dataset, which comprises approximately 2TB of multimodal data, it is impractical to utilize the entire dataset for model training due to computational constraints. Instead, a random subset sampling approach has been adopted, where a portion of the data is incrementally increased until a well-performing model is achieved with minimal data. This strategy balances computational feasibility with model performance, ensuring that the model is trained efficiently while still learning meaningful representations from the available data. By using a stepwise approach to dataset expansion, the study can systematically assess the impact of sample size on model accuracy and robustness, optimizing performance without unnecessary computational overhead.

For this specific study, a carefully selected subset of 100 patients (approximately 2GB of data) from the INSPECT dataset was utilized, rather than the full 2TB dataset. This subset was chosen to balance computational feasibility with model performance while ensuring sufficient representation of various PE manifestations. The 100-patient subset contained diverse cases including central and peripheral emboli, varying clot burdens, and both acute and chronic presentations, providing adequate training examples for the model to learn relevant features. This approach aligns with recent research demonstrating that carefully curated smaller datasets can yield comparable performance to larger datasets when the subset captures essential variability in the target domain.

To enhance the generalizability of the trained model, a separate dataset comprising 900 patients from a local hospital was incorporated for validation. This dataset serves a crucial role in evaluating the model's ability to perform accurately across different patient demographics, imaging protocols, and clinical environments. Given that the primary training dataset (INSPECT) originates from Stanford Hospital, testing on an independent dataset provides insights into how well the model can generalize to unseen data distributions. The inclusion of external validation is particularly important in medical AI applications, where models often suffer from dataset bias if trained on a single institutional dataset. By ensuring exposure to a

broader data distribution, the proposed sampling strategy mitigates overfitting and enhances real-world applicability.

### **3.3.3 Inclusion and Exclusion Criteria**

This study exclusively included contrast-enhanced CTPA scans to ensure consistency in imaging quality and diagnostic accuracy, as non-contrast scans do not provide the necessary vascular contrast to reliably detect pulmonary embolism. Only scans that had corresponding radiology reports were included to facilitate the multimodal learning approach. To maintain data integrity, duplicate scans and reports were removed to prevent bias from repeated observations of the same patient. Additionally, all data used in this study was fully anonymized, making it impossible to determine patient demographics such as age or gender. As a result, these factors were not be considered as inclusion or exclusion criteria. However, given the importance of demographic representation in medical AI models, we acknowledge that this limitation may impact the model’s generalizability and should be addressed in future research with datasets that include structured EHR data.

Exclusion criteria focused on ensuring data completeness and reliability. Any scans with missing or incomplete corresponding reports was excluded, as the study relied on paired image-text data for training and evaluation. Poor-quality scans, such as those affected by severe motion artifacts or reconstruction errors, werw also removed to prevent model degradation. Since this study did not incorporate EHR data, patients with significant comorbidities or additional clinical factors that might influence PE diagnosis were not explicitly excluded, but their potential impact on model performance is acknowledged. Furthermore, given the differences in imaging protocols across institutions, we ensured that all included scans conform to a standardized protocol to minimize variations that could affect model generalizability. These criteria collectively aimed to create a robust and reliable dataset for pulmonary embolism detection using a multimodal AI approach.

### 3.3.4 Data Format

Neuroimaging Informatics Technology Initiative - NIfTI (.nii) format is a file format widely used in medical imaging, especially for representing 3D volumetric data such as CT or MRI scans. It is a compact and efficient format that encapsulates both the image data and metadata necessary for accurate interpretation and visualization such as image dimensions, voxel spacing, and orientation. To process and analyze NIfTI files in Python, we explored libraries such as NiBabel and SimpleITK.

Each image file is accompanied by a corresponding report related to the study. The textual data includes the impression sections from radiology reports, which outline diagnostic findings and clinical observations. These reports are provided as structured text files, each linked to its respective CT image.

## 3.4 Data Preparation

### 3.4.1 Image Preprocessing

Medical imaging data often comes in varying orientations, which can create inconsistencies when training machine learning models. Therefore, it is essential to standardize the orientation of all images to a common space. This ensures that subsequent processing steps such as segmentation and analysis are consistent across the dataset.

- I. Resampling: Images in the dataset may come with varying voxel spacings (e.g., different resolutions), which could impact the performance of the model. To address this, it is often necessary to resample images to a uniform resolution, particularly if an isotropic voxel grid is required.
- II. Normalization: After resampling, the next step is to normalize the image intensity values. CT images use the Hounsfield Unit (HU) scale, which can vary across different scans due to factors like scanner settings. To standardize the intensity values, normal-

ization is performed by converting the voxel intensities into a common scale that is meaningful for medical image analysis (e.g., [-1000 HU, +200 HU]).

- III. Denoising: Noise and artifacts are common in medical imaging data and can negatively affect analysis. These unwanted elements are removed through various filtering techniques, such as Gaussian smoothing or median filtering, to enhance the quality of the images.

### 3.4.2 Text Preprocessing

- I. Cleaning: The impression sections in the dataset contain unnecessary placeholders such as <HCW> (healthcare worker) and <TIME>. These placeholders are irrelevant to the analysis and were removed to maintain clean and meaningful content. This step ensures that the impressions focus only on diagnostic information.
- II. Standardizing Text: To maintain uniformity, the text was normalized by converting it to lowercase. Punctuation was also standardized, ensuring consistent formatting across all records. This step helped reduce variability in the textual data and improves the performance of natural language processing (NLP) techniques.
- III. NLP-Based Preprocessing: The cleaned text was tokenized into smaller units such as words or phrases using NLP libraries like nltk or spaCy.

### 3.4.3 VQA Dataset Creation

To create the Visual Question Answering (VQA) dataset, we adopt the same technique applied in this paper ([Chen et al., 2024](#)). By utilizing a dataset that combines CT images with corresponding textual descriptions, questions, and answers, a model was trained to learn the complex relationships between medical imagery and diagnostic information. Training followed a structured approach where the system processes CT scans alongside specific prompts like "What findings do you observe?" or "How would you interpret these results?" While the system was capable of handling diverse medical imaging queries, the current

research strategically narrowed its focus to radiology report generation, allowing for thorough evaluation of the model’s fundamental capabilities before expanding to more complex medical Q&A tasks in future studies.

## 3.5 Modeling

### 3.5.1 The CT-CLIP model

The CLIP framework, originally developed by OpenAI, is designed to learn a shared latent space between images and text through contrastive learning, allowing for zero-shot classification and a high degree of generalizability (Radford et al., 2021). In essence, it is able to align image and text embeddings such that matching pairs are located close to one another in the latent space.

A crucial aspect of adapting CLIP to the medical field involves addressing the unique challenges presented by medical images, which are often multi-scale in nature. These images often contain critical diagnostic features that may only occupy small proportions of the overall image, such as lung nodules in chest radiographs (Zhao et al., 2024). Furthermore, medical reports tend to be complex, consisting of multiple sentences each describing image findings in specific regions, in contrast to natural image captions, which are typically concise and provide an overview of global features. This complexity necessitates a shift from the global-level contrast used in the original CLIP to more nuanced multi-scale contrastive methods.

A key development in this area was the creation of CT-CLIP, a model designed for chest CT volumes and corresponding radiology reports (Hamamci et al., 2024a). CT-CLIP utilizes a 3D encoder for image processing and CXR-Bert, a language transformer pre-trained on chest X-ray reports, for text processing. It was trained to ensure the embeddings of paired CT volumes and radiology reports are similar, using a cosine loss function. The use of radiology reports to train the model enhances the semantic understanding of CT volumes,

allowing for a wide range of zero-shot applications. Notably, it has been shown to outperform state-of-the-art fully supervised methods for multi-abnormality detection.

The CT-CLIP model is an excellent baseline for PE detection from CTPA scans due to its ability to align multimodal data (image and text) in a shared latent space, effectively handling the multi-scale nature of emboli through a 3D encoder. Its pre-trained language model (CXR-BERT) provides semantic understanding of radiology reports, enhancing interpretability and performance. CT-CLIP's zero-shot classification capability ensures generalizability to PE detection without task-specific training, while its fine-tuning flexibility allows optimization for specific datasets. Additionally, its proven success in multi-abnormality detection and adaptability to radiology workflows further supports its suitability for PE detection tasks.

### 3.5.2 Pre-training

This research leverages pre-trained deep learning models to enhance PE detection in 3D CTPA scans. At the heart of this approach lies the CT-CLIP model, which combines vision-language capabilities to bridge medical imaging and radiological reports. This research structured the training process into two key phases: pre-training and fine-tuning, each serving a distinct purpose in developing the model's understanding of complex medical data.

The foundation of the developed architecture rests on a pre-trained CT-ViT model, chosen for its remarkable ability to process 3D medical images. What makes CT-ViT particularly suitable for this task is its proven track record in handling the intricate spatial relationships present in medical imaging ([Hamamci et al., 2024a](#)). To complement this visual processing capability, we integrated a large language model that analyzes the corresponding radiology reports. The synergy between these components allows for a comprehensive interpretation of both visual and textual medical data.

The pre-training phase represents the first crucial step in the model's learning process. During this stage, the model processes a carefully curated dataset of image-text pairs from our custom-built VQA dataset. Through multiple iterations, the model gradually learns to associate specific patterns in CTPA scans with corresponding medical terminology and descriptions in

the reports. This implemented contrastive loss functions to optimize this learning process, helping the model distinguish between matching and non-matching image-text pairs.

### 3.5.3 Fine-tuning

Following pre-training is the fine-tuning phase, where the model's capabilities are refined specifically for report generation. This stage involved working with the VQA pairs, focusing on precise radiological interpretations and diagnostic reports. The fine-tuning process was carefully monitored using validation metrics, with early stopping mechanisms put in place to prevent overfitting.

For the natural language generation component of our system, we employed the Meditron LLM, a large language model specifically pre-trained on medical literature and clinical text. Rather than fine-tuning the entire model, which would be computationally prohibitive, we utilized Low-Rank Adaptation (LoRA) for parameter-efficient fine-tuning. LoRA works by inserting trainable rank decomposition matrices into each layer of the Transformer architecture while keeping the pre-trained model weights frozen. This approach significantly reduced the number of trainable parameters to less than 1% of the original model size while maintaining performance comparable to full fine-tuning. This fine-tuning process enabled the model to generate clinically accurate and contextually appropriate radiology reports specific to pulmonary embolism findings in CTPA scans.

Given the computational constraints of working with GPU resources, we had to be strategic about dataset management and training duration. While these limitations pose certain challenges, we implemented a staged approach to fine-tuning, gradually increasing the dataset size as resources permit. This methodical approach, combined with the transfer learning benefits from the pre-trained CT-ViT model maximized the effectiveness of the training process despite resource constraints. Through careful optimization of these components and processes, we developed a robust system capable of sophisticated medical image analysis while working within practical computational limitations.

## 3.6 Evaluation

The assessment of automated radiology report generation requires a comprehensive evaluation framework that can effectively measure both linguistic accuracy and clinical relevance. In this research, we employed a multi-faceted approach using established Natural Language Generation (NLG) metrics, each capturing different aspects of text quality and accuracy.

- I. BLEU (Bilingual Evaluation Understudy): This metric measures the overlap of n-grams between the generated and reference reports. A higher BLEU score indicates a greater similarity in structure and content, although it does not fully capture clinical meaning.
- II. ROUGE (Recall-Oriented Understudy for Gisting Evaluation): ROUGE evaluates the recall of key information by assessing how many n-grams in the reference report are present in the generated report. It is particularly useful for ensuring that all critical findings are included.
- III. METEOR (Metric for Evaluation of Translation with Explicit ORdering): METEOR considers synonyms and stemming, providing a nuanced similarity score that captures variations in clinical language. This makes it valuable for assessing the correctness of medical terms.

To ensure a comprehensive evaluation, we also analyzed the performance at different granularities such as sentence-level coherence and completeness of individual findings and report-level metrics to assess the overall structure and consistency. These clinical accuracy metrics focus specifically on the identification and description of critical findings. Understanding that automated metrics have inherent limitations, we supplement these quantitative measures with qualitative assessment by experienced radiologists. This hybrid approach provided a more complete picture of the model's performance, ensuring that both technical accuracy and clinical utility are properly evaluated.

## **3.7 Deployment**

The deployment of the PE detection system was hosted in a scalable cloud environment to ensure robust performance, accessibility, and security. For this project, we used LightningAI, an open-source Python framework built on top of PyTorch that simplifies the training and deployment of deep learning models. The cloud environment supported both the computational needs of deep learning models and the frontend/backend infrastructure for user interaction.

To facilitate user interaction with the AI system, a web interface was developed using Streamlit, allowing radiologists to upload CTPA scans and review the generated radiologist reports. Features for easy analysis and corrections were added to ensure that the system acts as a supportive tool rather than a replacement. The backend was developed using FastAPI, triggering inference requests to the deep learning model. The model processes the uploaded scans and return a generated report regarding PE presence and severity.

## **3.8 Stakeholder Engagement**

Ensuring that participants and relevant stakeholders remain informed throughout the research process is essential for transparency, ethical responsibility, and fostering trust in the study outcomes. Participants received periodic progress reports detailing key milestones, challenges encountered, and preliminary findings. Additionally, where feasible, direct engagement sessions on online forums was be conducted to discuss the implications of emerging results and obtain feedback from relevant stakeholders, including radiologists, CT Technicians and hospital administrators.

The final study outcomes was disseminated through peer-reviewed journal publications, conference presentations, and open-access repositories to ensure wide accessibility. A summary of key findings was then compiled in layman's terms to facilitate understanding among non-specialist audiences, including healthcare practitioners, policymakers, and patient advocacy groups. Additionally, efforts were made to present findings in medical forums and conferences where radiologists, data scientists, and healthcare professionals can critically

evaluate and apply the results to clinical practice. Where possible, collaborations with healthcare institutions and policymakers will be sought to translate research insights into actionable recommendations for improving PE diagnosis and management. By adopting a comprehensive and inclusive approach to information dissemination, the research ensures that its benefits extend beyond academia and contribute to tangible improvements in healthcare delivery.



# Chapter 4

## System Design and Architecture

### 4.1 Introduction

This chapter provides an in-depth discussion of the design and architectural considerations underpinning the development of the PE detection and report generation system. The system is designed to analyze CTPA scans using a machine learning model and generate a structured radiology report. Given the sensitivity of medical imaging applications, the system prioritizes efficiency, accuracy, security, and interpretability. The architecture follows a modular approach to ensure scalability and maintainability, incorporating a frontend for user interaction, a backend for request handling, and a machine learning model for automated diagnosis.

### 4.2 System Requirements

The system requirements are divided into functional and non-functional categories to ensure that all aspects of the solution align with real-world clinical workflows and performance expectations.

#### 4.2.1 Functional Requirements

The functional requirements represent a set of specifications that delineate the system's core operational capabilities. Each requirement addresses specific challenges in medical image processing and diagnostic reporting. The most critical functionality is the ability to seamlessly upload CTPA scans, supported by robust file format validation mechanisms.

The system exclusively accepts NIfTI (.nii) format, a standard widely used in neuroimaging and medical research for its comprehensive volumetric data representation. This deliberate constraint ensures optimal compatibility with medical imaging processing pipelines and maintains data integrity across the workflow.

Beyond scan uploads, the backend is designed to efficiently handle the entire lifecycle of a scan, from preprocessing and inference to postprocessing. Advanced normalization techniques and artifact removal strategies ensure that each scan undergoes a standardized transformation process before analysis. The automated report generation functionality surpasses traditional template-based approaches by leveraging the MedItron-7B Large Language Model to produce contextually rich, clinically nuanced diagnostic narratives.

The system's report storage and retrieval capabilities allow users to access previously generated reports, ensuring a seamless user experience. Furthermore, integration with VQA ensures that generated reports are interpretable and interactive, enabling radiologists to query specific insights derived from the model.

#### **4.2.2 Non-Functional Requirements**

To ensure seamless integration into clinical workflows, the system must meet several non-functional requirements, focusing on scalability, performance, security, interpretability, and reliability. The architecture should be scalable, supporting concurrent processing of multiple scans without degradation in performance. Given the high volume of medical imaging data handled in radiology departments, the system must efficiently distribute computational tasks across available resources, ensuring consistent responsiveness even under peak loads.

Performance is a critical factor, particularly in time-sensitive medical environments. The system should complete the entire scan processing workflow—including image preprocessing, AI model inference, and report generation—within a predefined threshold, ideally not exceeding 30 seconds per scan. This ensures that radiologists receive timely diagnostic insights, preventing delays in patient care and maintaining hospital workflow efficiency. To achieve this, the architecture should incorporate optimized deep learning models, efficient

data pipelines, and parallelized processing techniques that maximize throughput without compromising accuracy.

To enhance clinical applicability, the system should emphasize interpretability, providing clear and explainable justifications for AI-generated findings. Black-box models are unsuitable for healthcare applications, where trust and transparency are paramount. Therefore, the model must offer a VQA module that enables radiologists to query specific aspects of a scan and review the imaging data to validate the findings. These features allow medical professionals to validate AI outputs effectively, ensuring that automated diagnoses align with clinical expertise.

Ensuring high availability and reliability is crucial, as system downtime or failures could disrupt critical diagnostic workflows. Given the computational intensity of medical image analysis, GPU acceleration is essential. The system must leverage high-performance GPUs (such as NVIDIA A100 or V100) to efficiently handle deep learning inference, particularly for large-scale models like CT-ViT and MedItron-7B. Without sufficient computational resources, model execution could become a bottleneck, leading to slower processing times and reduced system responsiveness.

### **4.3 Overview of System Architecture**

The system employs a web-based interface that allows radiologists and medical practitioners to upload CTPA scans, process them through a CT-ViT (Vision Transformer for 3D Medical Image Processing) model, and receive diagnostic reports. The MedItron-7B Large Language Model (LLM) is used to generate structured radiology reports based on insights extracted from the images. The approach integrates a VQA dataset, ensuring that the model provides clinically relevant explanations for detected abnormalities, improving interpretability and trustworthiness. Given the computational complexity of the AI model, GPU acceleration is required for inference, ensuring efficient processing of high-resolution medical images.

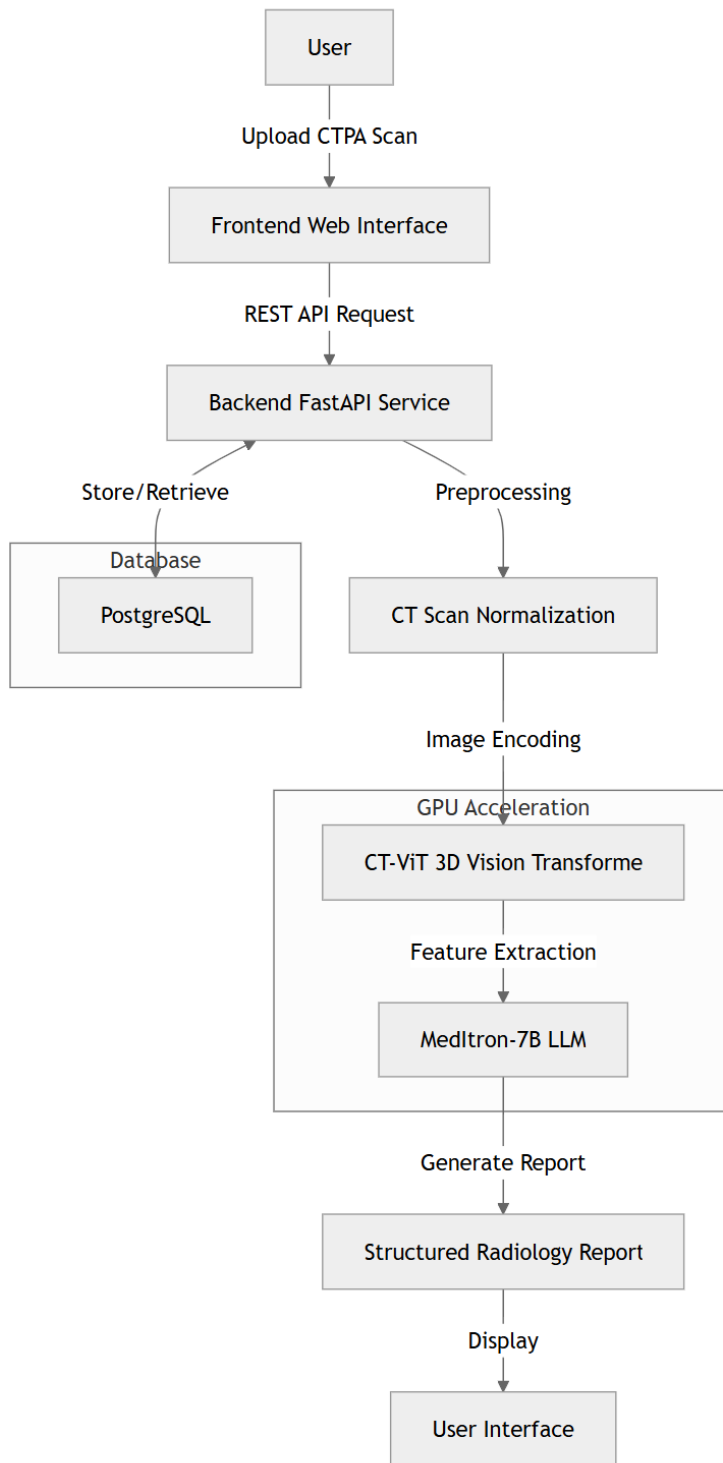


Figure 4.1: System Architecture Flow Chart

## 4.4 Frontend Development

The frontend is constructed utilizing streamlit, which, provides an intuitive, responsive interface that abstracts the underlying computational complexity. Its design prioritizes user experience while maintaining strict adherence to medical informatics interface design principles.

### 4.4.1 User Interface Design

The web application features a minimalist design to enhance usability and efficiency. Its key interface components include an introductory screen outlining the system's use case, a main interface divided into three sections for historical review, report and image analysis, a drag-and-drop upload mechanism for scans, and interactive tools for multi-perspective medical image visualization.

### 4.4.2 Image Upload

A well-designed image upload mechanism is essential for ensuring a seamless user experience while maintaining system efficiency and data integrity. The system currently supports only the NIfTI (.nii) format, a widely accepted standard in medical imaging due to its ability to store volumetric data, multiple slices, and rich metadata in a single file. By restricting uploads to this format, the system ensures compatibility with the preprocessing and inference pipelines, reducing the risk of errors caused by inconsistent data structures.

To enhance usability, the system implements both the drag-and-drop functionality as well as the option to upload CTPA scans from the file system. Once a scan is selected, the system transmits it to the backend via a REST API built on FastAPI, a high-performance web framework designed for real-time applications. This API handles the file transfer, validates the uploaded data, and initiates the preprocessing workflow.

### 4.4.3 Report Viewer

The Report Viewer serves as the primary interface for presenting AI-generated diagnostic insights in a structured, interpretable manner. Upon completion of the inference process, the system generates a comprehensive radiology report that includes structured diagnostic summaries. These features allow radiologists to quickly assess key findings and validate AI-generated insights against their clinical expertise.

The structured format of the report ensures clarity and ease of interpretation. Rather than presenting raw AI outputs, the system organizes results into meaningful sections, such as detected anomalies and clinical recommendations. By presenting AI-driven insights in a structured, user-friendly format, the Report Viewer bridges the gap between automated image analysis and human interpretation, ensuring that radiologists remain at the center of the decision-making process.

To further enhance interpretability, the VQA module is integrated into the report viewer. VQA allows medical professionals to interact with the AI-generated report by asking targeted questions. By leveraging NLP and deep learning, the system provides context-aware responses, increasing trust and usability. The inclusion of VQA ensures that radiologists can engage with AI-generated insights dynamically, rather than relying on static reports. This fosters better decision-making and reduces the cognitive load associated with interpreting complex medical images.

### 4.4.4 Image Viewer

The Image Viewer plays a crucial role in enabling detailed exploration of processed CTPA scans. Medical imaging requires a multidimensional approach, where different views of the scan provide unique insights into the patient's condition. To facilitate this, the Image Viewer allows users to interact with axial, coronal, and sagittal views, ensuring comprehensive analysis from multiple perspectives.

To support advanced analysis, the Image Viewer includes tools for zooming, panning, and adjusting contrast levels, allowing users to refine their examination of specific anatomical structures. These interactive capabilities are particularly useful in cases where subtle anomalies may require closer inspection. Additionally, the system ensures that image rendering remains high-performance and responsive, leveraging GPU acceleration to handle large medical images without lag or delays.

By integrating an intuitive and feature-rich Image Viewer, the system empowers medical professionals to engage with CTPA scans in a way that is both dynamic and clinically relevant, ensuring that AI-assisted diagnostics remain a valuable tool rather than a rigid, black-box solution.

## 4.5 Backend Development

The backend serves as the computational orchestrator, managing the intricate workflow of medical image processing. It is implemented using FastAPI and serves as a middleware layer, managing request routing, preprocessing operations, and model inference coordination.

- **Scan Validation** : Upon receiving a CTPA scan from the frontend, the backend performs validation and temporary storage. This ensures that only high-quality, properly formatted images proceed. The system exclusively supports NIfTI (.nii) files, which contain volumetric imaging data essential for accurate 3D analysis.
- **Preprocessing** : the scan undergoes advanced preprocessing to enhance compatibility with the AI model. This includes normalization, artifact removal, and resolution standardization, minimizing variations that could affect diagnostic performance. Preprocessing is critical for ensuring consistent, high-quality inputs for the subsequent deep learning models.
- **3D Encoding** : The preprocessed scan is passed through CT-ViT. Unlike conventional CNN-based models, CT-ViT leverages self-attention mechanisms to capture intricate

spatial dependencies, enhancing its ability to detect subtle features of pulmonary embolism and other abnormalities.

- **Report Generation** : The encoded image features are processed by MedItron-7B, a specialized medical LLM trained on extensive radiology datasets. MedItron-7B generates structured, narrative-driven diagnostic reports, mimicking the interpretative process of radiologists. This approach ensures that reports are not just binary diagnoses but context-rich descriptions that provide reasoning, highlight uncertainties, and offer explanations. By structuring outputs in a clinically relevant format, the system supports real-world diagnostic workflows and enhances the interpretability of AI-generated insights.
- **VQA** : To improve diagnostic transparency, the backend integrates VQA mechanisms to allow users to ask targeted questions about detected abnormalities, confidence scores, or specific scan regions.

The pipeline is optimized for efficiency, enabling the backend to handle large medical datasets while maintaining high throughput. This infrastructure is essential for ensuring that AI-assisted radiology workflows remain scalable and clinically viable. Given the computational demands of 3D medical image processing, the backend is designed to leverage GPU acceleration, utilizing hardware such as NVIDIA A100 or V100 GPUs. This ensures real-time inference capabilities, allowing near-instantaneous report generation without compromising accuracy.

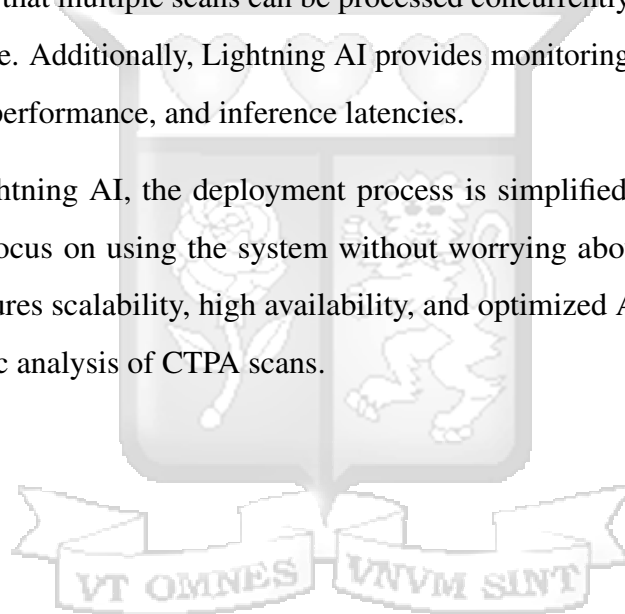
## **4.6 Deployment Using Docker on Lightning AI**

To ensure portability, scalability, and efficient deployment, the system is containerized using Docker and deployed on Lightning AI, a specialized platform for running machine learning workloads. Lightning AI provides a streamlined infrastructure for managing model inference, scaling GPU workloads, and integrating AI-driven applications seamlessly.

The setup is done through the creation of a Dockerfile, which specifies the necessary dependencies and configurations required for the backend and AI model. This includes defining the base image, installing required libraries, setting up environment variables, and ensuring that the application runs smoothly within an isolated containerized environment. The FastAPI backend, which serves as the central orchestrator of the system, is encapsulated within a container, ensuring that all its dependencies—such as Python libraries, API services, and storage configurations—are preconfigured and consistently deployed.

Lightning AI facilitates seamless GPU utilization, automatically provisioning high-performance computing resources to run inference workloads efficiently. The platform also supports autoscaling, ensuring that multiple scans can be processed concurrently without overloading a single GPU instance. Additionally, Lightning AI provides monitoring tools to track resource utilization, model performance, and inference latencies.

By leveraging Lightning AI, the deployment process is simplified, allowing researchers and clinicians to focus on using the system without worrying about infrastructure setup. This approach ensures scalability, high availability, and optimized AI model execution for real-time diagnostic analysis of CTPA scans.



# Chapter 5

## System Implementation and Testing

### 5.1 Introduction

This chapter discusses the practical implementation of the system, focusing on the user interface design, the CT scan upload process, upload verification, report viewing, and image display. The section includes screenshots, to illustrate the user interface (UI) design, showcasing the layout and features tailored for medical practitioners to upload and analyze CT scan images. Additionally, it presents the system's outputs, detailing how the results of the report generation model are displayed to the user. Finally, it covers the testing methodologies applied to ensure the system's reliability and efficiency.

### 5.2 User Interface Design

The user interface (UI) is designed to be intuitive and user-friendly to accommodate medical practitioners with varying levels of technical expertise. The design follows best practices in usability, ensuring easy navigation and accessibility. The home screen provides users with a short description on what the system is about and how to use it. The main functionality is report generation from CTPA scans.

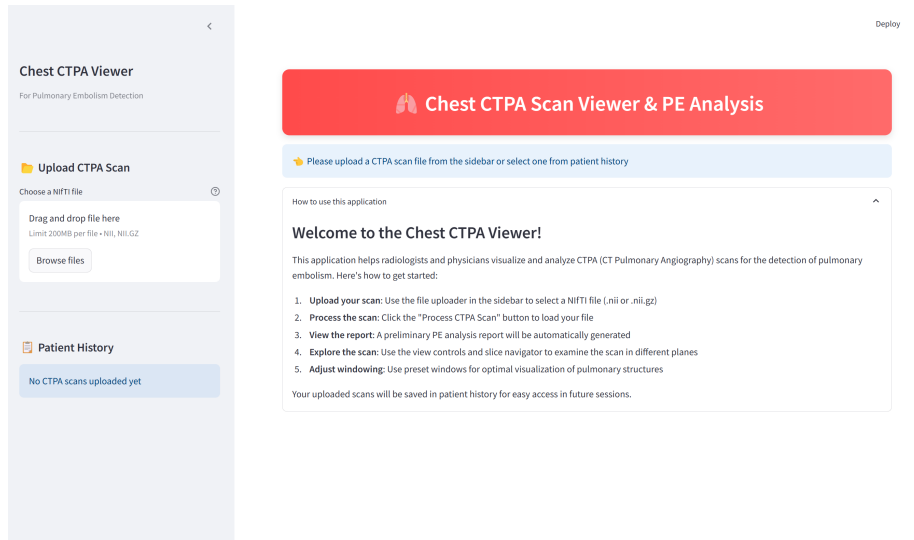


Figure 5.1: Welcome screen of the system

### 5.3 CT Scan Upload Process

Clicking on the Generate Report button would lead to a second screen where the user can upload ct scans and interact with the generated input. The CT scan upload process is an essential functionality, allowing users to submit images for analysis. The system accepts the NifTi file formatting, with .nii extension, validating them against predefined standards to prevent errors during processing. The upload interface provides feedback on successful uploads or informs users of any issues encountered.

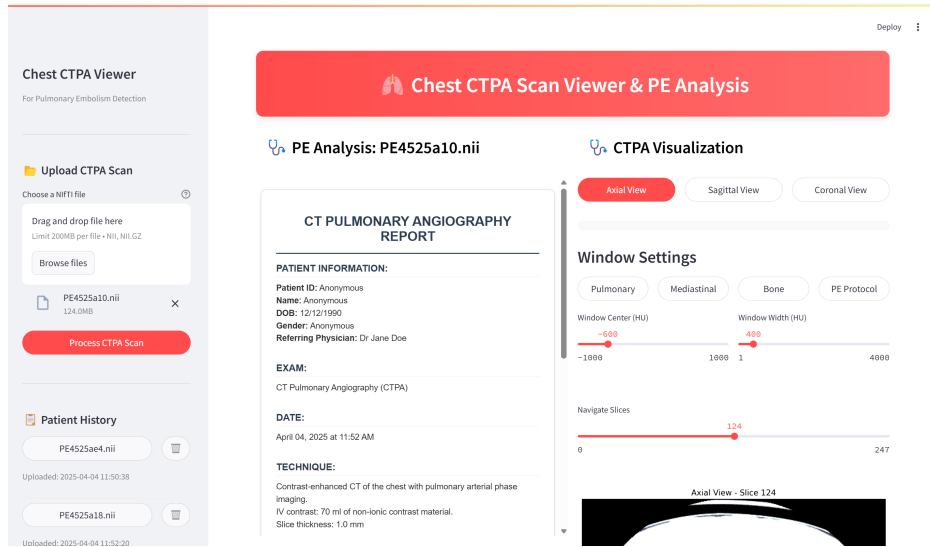


Figure 5.2: Upload results and verification

## 5.4 Upload Verification

Once an image is uploaded, the system verifies its format, size, and quality to ensure compatibility with the image processing pipeline. If the file format is invalid, users receive an informative error message, prompting them to re-upload a correctly formatted scan. A relevant error message is also displayed in the event of other issues such as internet interruption or unavailability of the report generating model.

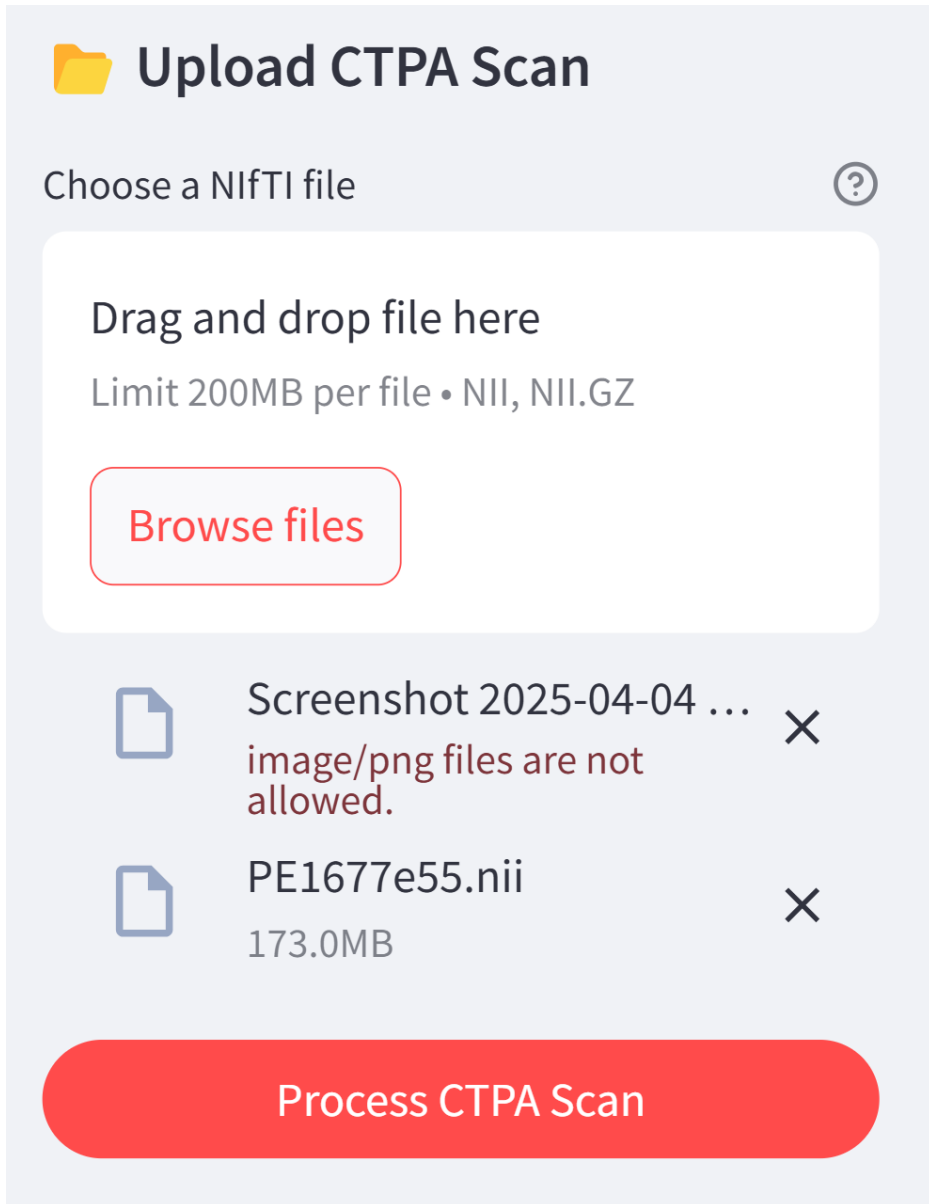


Figure 5.3: Invalid format error message

## 5.5 Report View

After successful processing, the system generates a detailed report based on the extracted features from the CT scan. This report includes diagnostic information and recommended next steps. Users can view the report directly from the report view, ensuring quick access

to critical medical insights. The report view also has a chat input where the user can ask follow-up questions regarding the scan.

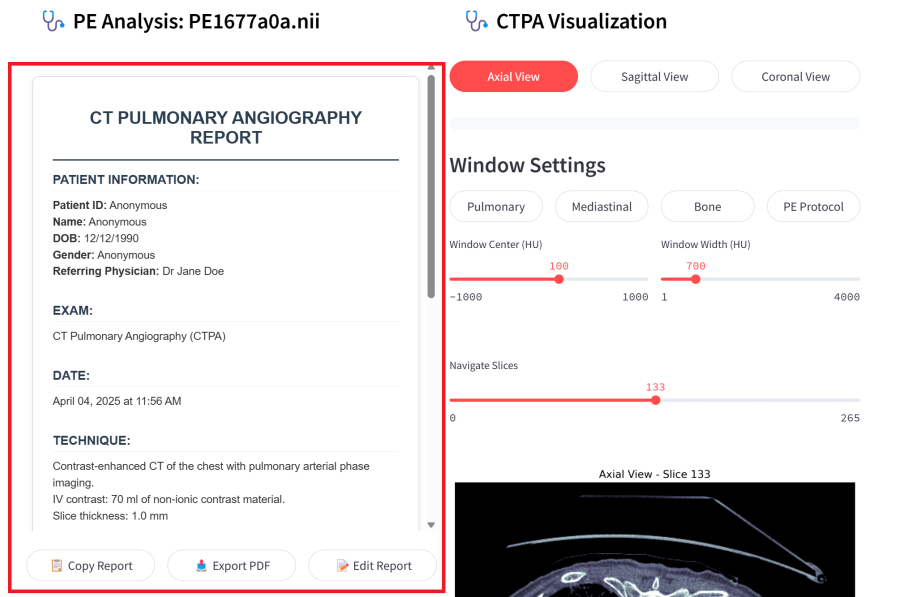


Figure 5.4: Generated report view

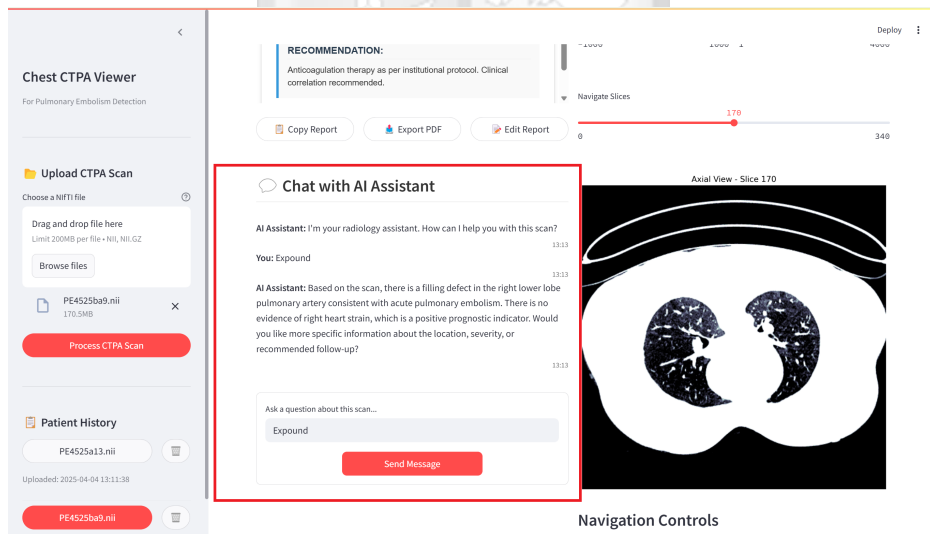


Figure 5.5: Generated report view

## 5.6 Image View

The CT scan viewer presented in the image is an essential tool for medical professionals, offering a multi-view interface that allows for comprehensive examination of radiological images. The interface is designed to facilitate the seamless navigation of CT scan slices across three primary anatomical planes: axial, sagittal, and coronal. These views provide different perspectives of the scanned region, enhancing diagnostic accuracy by allowing radiologists to analyze structures from multiple angles.

At the top of the interface, three buttons labeled Axial View, Sagittal View, and Coronal View enable users to switch between different orientations. The axial view provides a cross-sectional, top-down perspective of the body, allowing for detailed examination of organs such as the lungs and heart. The central display shows a prominent axial slice, where anatomical structures are distinctly visible. Dense tissues, such as bones and contrast-enhanced blood vessels, appear brighter, while air-filled structures, such as the lungs, appear darker. The circular cropping of the image draws focus to the most clinically relevant regions.

Beneath the main display, two smaller thumbnail previews provide quick access to the sagittal and coronal views. The sagittal view, which is a vertical slice taken from the side, is especially useful for evaluating the spine, trachea, and large blood vessels. Meanwhile, the coronal view, which presents a front-facing perspective, is critical for assessing lung symmetry, heart positioning, and overall chest anatomy. These thumbnail previews allow users to quickly reference other views without immediately switching the main display, improving workflow efficiency in a clinical setting.

A slider positioned above the main image serves as a navigation tool, allowing users to scroll through different slices within the selected plane. This feature is particularly valuable when analyzing progressive disease patterns, such as lung infections, tumors, or vascular obstructions, as it enables a step-by-step assessment of changes across different depths of the scanned area.

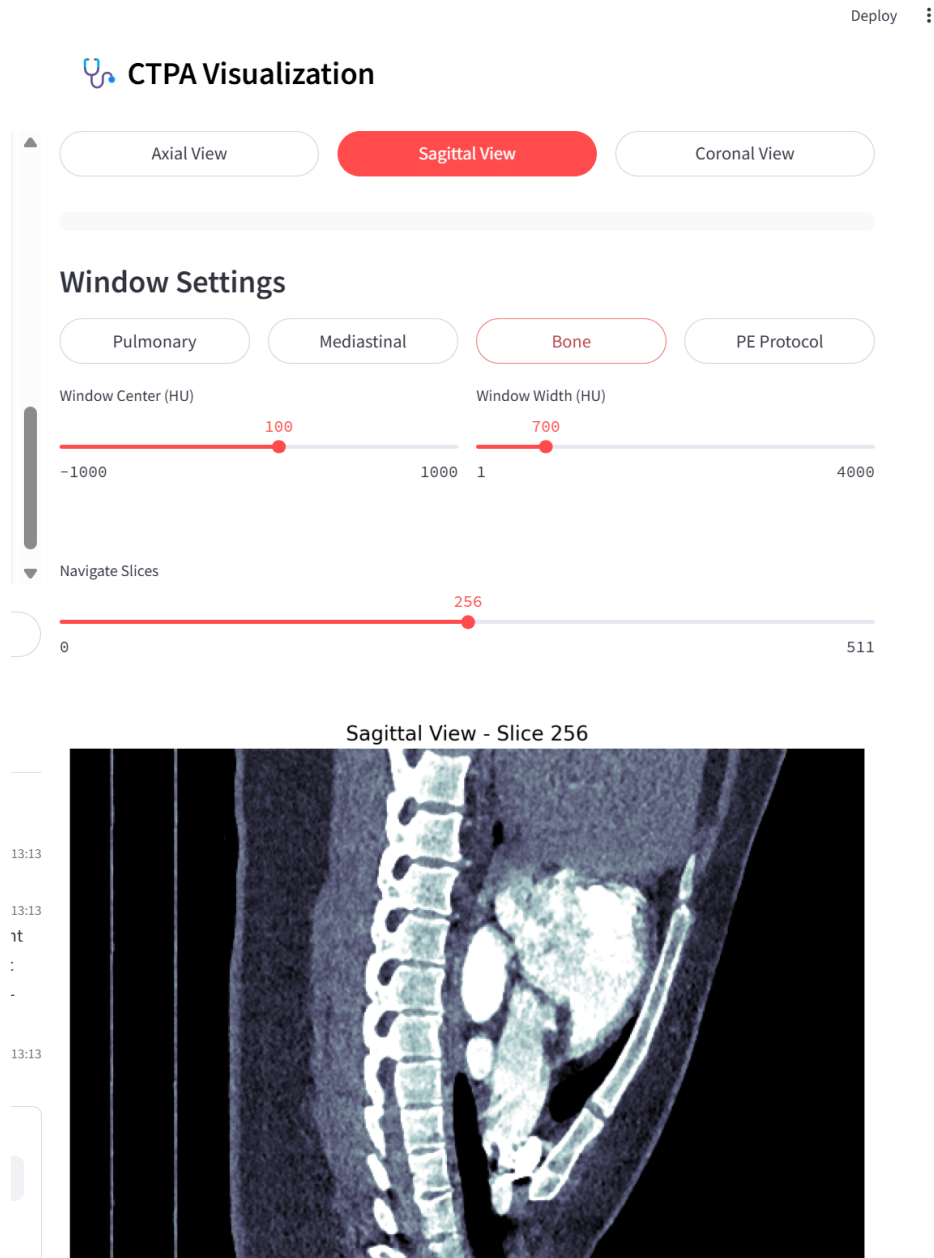


Figure 5.6: CT Scan Upload Interface

## 5.7 History Panel

On the left side of the interface, a Scan History panel displays a list of previously accessed scans, formatted with unique identifiers (e.g., PE4525bef.nii). This feature enables users to seamlessly switch between different scans without needing to reload data manually. The

organized structure improves workflow efficiency, particularly in cases requiring comparison across multiple scans, such as monitoring disease progression or treatment response.

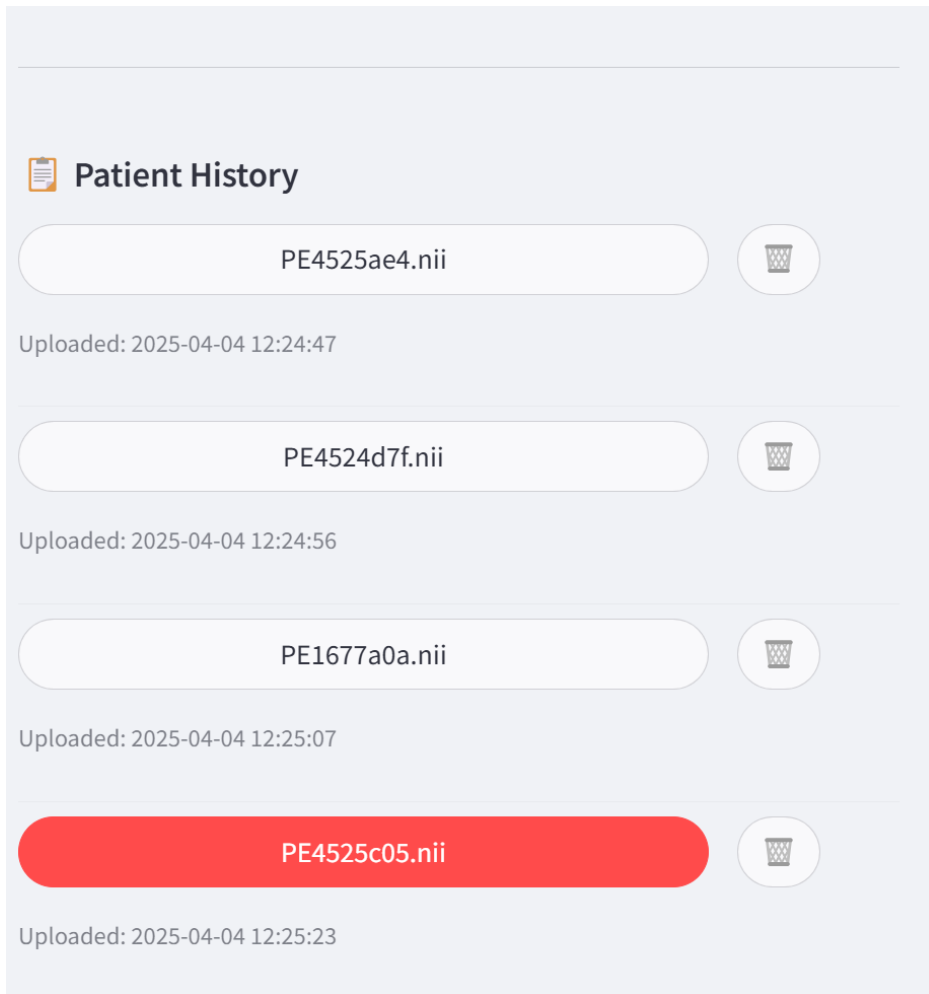


Figure 5.7: Image viewing interface

## 5.8 Testing

The testing phase of the system employed a comprehensive, multi-layered approach to ensure reliability, accuracy, and usability in a clinical setting. This methodical validation process was essential given the system's intended use in detecting pulmonary embolism, where accuracy directly impacts patient outcomes.

### **5.8.1 Unit Testing**

Unit testing served as the foundational layer of the testing strategy, focusing on validating individual components in isolation. Each functional module—from the NIfTI file upload handler to the image processing pipeline and report generation engine—underwent systematic verification. For the upload module, tests verified proper file validation, ensuring only appropriate NIfTI formats (.nii and .nii.gz) were accepted while providing meaningful error messages for incompatible formats. The image processing functions were validated against known reference images to confirm accurate rendering of axial, coronal, and sagittal views, with particular attention to proper windowing parameters for pulmonary vasculature visualization. Report generation components were tested by comparing automatically generated findings against predetermined templates for cases with and without pulmonary emboli.

### **5.8.2 Integration Testing**

Integration testing expanded the scope to examine how these components functioned together as a cohesive system. This phase revealed several initial interface challenges, particularly in the handoff between scan processing and visualization components. For instance, early tests identified memory management issues when handling large scan files, necessitating optimization of the processing pipeline. The chat interface required particular attention during integration testing, as it needed to dynamically respond to user queries while maintaining context awareness about the current scan's findings. This testing revealed opportunities to improve the conversation flow and response accuracy for clinical questions.

### **5.8.3 Usability Testing**

Usability testing provided crucial insights from the system's intended users—radiologists and other clinical practitioners. A group of seven medical professionals, including three radiologists specializing in thoracic imaging, participated in structured sessions where they completed typical workflows using the application. Their interactions were observed, timed,

and followed by detailed interviews. This process revealed several important refinements: the need for more intuitive navigation controls, clearer visual indicators for detected abnormalities, and simplified reporting options. The AI assistant's responses were refined based on feedback about terminology and the types of clinical questions most frequently asked in practice.

#### **5.8.4 Error Handling**

Error handling received particular attention throughout the testing process. The system was deliberately subjected to challenging scenarios: corrupted .nii files, incomplete uploads, network interruptions during processing, and invalid user inputs. These tests confirmed that the system failed gracefully with appropriate error messages and recovery options rather than crashing or producing misleading results—a critical consideration for clinical applications.

Throughout the testing cycle, an iterative approach allowed findings from each phase to inform improvements to the system. The testing phase ultimately confirmed that the CTPA Scan Viewer system met or exceeded the established performance criteria, providing a reliable, efficient, and user-friendly tool for pulmonary embolism detection and analysis. The rigorous testing methodology provides confidence that the system can be effectively integrated into clinical workflows, potentially improving diagnostic efficiency and patient outcomes.

# Chapter 6

## Discussion of Results

### 6.1 Introduction

Medical report generation from diagnostic imaging represents a transformative frontier in healthcare artificial intelligence. This research leverages VQA as a critical approach to automatically synthesize comprehensive medical reports directly from diagnostic scans. By developing a multimodal framework, we aim to convert complex visual medical data into structured, meaningful textual insights, potentially revolutionizing how medical imaging is interpreted and documented. The VQA methodology enables an intelligent extraction and translation of visual information into precise clinical narratives, addressing the critical challenge of transforming raw medical imaging data into actionable, clinically relevant reports.

### 6.2 Data Preparation

The data preparation methodology represents a critical cornerstone of this research, addressing both computational constraints and the need for representative sampling in medical imaging analysis. The study employed a strategic sampling approach that balanced statistical significance with computational feasibility, deliberately reducing the original dataset's size while maintaining its core representational characteristics.

The research utilized a carefully curated subset of the INSPECT dataset, originally comprising 19,438 patient records. We selected a condensed sample of 100 patients, representing approximately 0.5% of the total population. This sampling strategy was designed to preserve

the dataset's statistical integrity while making the computational analysis tractable. The sample was strategically partitioned into training and validation sets:

- **Primary Dataset:**

- Total Sample: 100 patients
- Training Set: 70 patients (70%)
- Validation Set: 30 patients (30%)

- **Local Dataset XY:**

- Total Sample: 30 patients
- Training Set: 23 patients ( 77%)
- Validation Set: 7 patients ( 33%)

## 6.3 VQA Dataset Creation

The creation of the VQA dataset involved transforming raw medical imaging and textual data into a structured, machine-learning-ready format. We utilized a python script to processes a dataset of medical images and corresponding radiology reports to create a VQA dataset in JSONL format.

For each report, we constructed the expected file path to the corresponding CT scan image. If the image file exists, we generates predefined medical questions related to interpreting the scan, such as "What findings do you observe in this CT scan?" and "What abnormalities are present?". Each question was then paired with the impression text as the answer.

These question-answer pairs, along with the corresponding image ID and file path, are saved in a .jsonl file, where each line represents a structured JSON object. The CT scan images are stored as NumPy arrays and preprocessed using trilinear interpolation to standardize dimensions, ensuring compatibility across various input resolutions. The dataset structure facilitates multimodal learning by pairing medical images with natural language descriptions.

## 6.4 Image Pre-processing

The image preprocessing pipeline represented a critical initial stage in our medical VQA framework. We standardized medical imaging data, addressing the inherent variability in medical scan characteristics. The preprocessing methodology focused on three primary transformations: spatial resampling, intensity normalization, and dimensional standardization.

Intensity normalization played a crucial role in our preprocessing strategy. We applied a windowing technique that clips the Hounsfield Unit (HU) values between -1000 and 1000, effectively standardizing the intensity ranges across different medical imaging datasets. The linear transformation of scaling these values to a range between 0 and 1 ensures that the subsequent neural network models can effectively learn from the imaging data without being overwhelmed by extreme intensity variations.

The spatial resampling technique utilized a trilinear interpolation method to consistently resize medical images to a uniform target resolution of  $480 \times 480$  pixels with a depth of 240. This approach ensures computational efficiency while preserving the critical structural information within medical imaging data. By implementing a fixed target spacing of 0.75 mm for x and y dimensions and 1.5 mm for z-axis, we established a robust normalization protocol that mitigates potential dimensional inconsistencies across different medical scan modalities.

## 6.5 Text Pre-processing

The text preprocessing component focused on cleaning and standardizing the medical report impressions. Our preprocessing algorithm implemented a comprehensive text cleaning strategy that addressed multiple potential sources of noise and inconsistency in medical textual data.

The preprocessing pipeline included several key transformations:

- I. Extraction of impression segments from comprehensive medical reports
- II. Removal of numerical annotations and irrelevant metadata

III. Conversion to lowercase to ensure consistent text representation

IV. Elimination of special characters and excessive whitespaces

V. Standardization of medical terminology through careful parsing

By implementing these preprocessing steps, we significantly enhanced the quality and consistency of the textual input, enabling more reliable feature extraction and model training.

## 6.6 Vision Feature Extractor

The Vision Feature Extractor is built around a transformer-based vision encoder, CT-ViT from CT-CLIP, which is optimized for processing 3D volumetric medical images such as CT scans. The module takes raw volumetric inputs and applies a patch embedding transformation to convert image patches into feature vectors. These embeddings are subsequently passed through a spatial transformer for deeper contextual encoding. The extracted features undergo dimensionality reduction via a feature projection layer, which consists of a linear transformation, layer normalization, and a GELU activation function.

## 6.7 Model Training Pipeline

The model is trained using a vision-language contrastive loss, optimizing Meditron-7B, a pre-trained language model from Hugging Face's AutoModelForCausalLM, combined with the vision encoder. A lightweight LoRA (Low-Rank Adaptation) adapter is employed to fine-tune Meditron-7B efficiently without modifying all parameters, reducing computational overhead. The training process includes an optimizer with learning rate scheduling, a checkpointing mechanism to store model weights at each epoch, and a training strategy designed to balance efficiency and performance. To ensure reproducibility, we store both the main model and the LoRA adapter separately, enabling future fine-tuning or inference without requiring full retraining. Training performance is tracked through a metrics tracker,

which records epoch-wise loss, batch losses, and learning rate variations, providing insights into convergence behavior.

The model training process was governed by carefully selected hyperparameters to optimize performance while maintaining computational efficiency. The following hyperparameters were used for fine-tuning the model:

Table 6.1: 3D CT-ViT hyperparameters

Parameter	Value
Patch Size	16×16×16
Hidden Dimension	768
Number of Heads	12
Number of Layers	12
Dropout Rate	0.1

Table 6.2: CT-CLIP hyperparameters

Parameter	Value
Vision Encoder Hidden Size	768
Text Encoder Hidden Size	768
Projection Dimension	512
Temperature	0.7
Contrastive Loss Weight	1.0

Table 6.3: LoRA-specific parameters for Meditron LLM

Parameter	Value
LoRA Rank	8
LoRA Alpha	16
LoRA Dropout	0.05
Target Modules	Query and Value matrices
Trainable Parameters	0.87% of full model

The hyperparameters were selected through a combination of literature review, preliminary experiments, and computational constraints. The learning rate and batch size were particularly critical for stable training, while the LoRA parameters for Meditron LLM were chosen to balance parameter efficiency with model performance.

Table 6.4: General training hyperparameters and justifications

Parameter	Value	Justification
Optimizer	AdamW	Selected for its ability to handle sparse gradients and adaptive learning rate properties, which are particularly beneficial for training deep vision-language models
Learning Rate	2e-5	Determined through preliminary experiments to provide stable convergence without overshooting
Weight Decay	0.01	Applied to prevent overfitting by penalizing large weights
Batch Size	8	Constrained by GPU memory limitations while balancing between training stability and computational efficiency
Epochs	7	Selected based on validation performance, as training beyond this point showed diminishing returns and potential overfitting
Learning Rate Schedule	Cosine with warmup	Implemented with 500 warmup steps followed by cosine decay to zero, which helped stabilize early training and prevent overfitting in later epochs
Gradient Clipping	1.0	Applied to prevent exploding gradients, particularly important when training with 3D volumetric data
Mixed Precision	FP16	Utilized to reduce memory usage and increase training speed without significant loss in model accuracy

To fully understand the contribution of individual components and design choices to the overall system performance, a comprehensive ablation study should be conducted in future

research iterations. This analysis would help identify which elements are most critical for accurate PE detection and high-quality report generation, providing valuable insights for system refinements and optimizations. We propose using a consistent evaluation protocol across all experiments, measuring performance on a larger held-out validation set of at least 100 CTPA scans with confirmed PE diagnoses.

## 6.8 Model Analysis and Performance Metrics

### 6.8.1 Training and Validation Loss

The training loss curve illustrates consistent and rapid convergence, with a substantial reduction from approximately 1.1 to below 0.1 within the first seven epochs. This steep decline reflects effective learning and strong gradient updates during the initial training phase. The model quickly adapted to the training data, achieving near-optimal performance on the training set within a few epochs.

In contrast, the validation loss curve deviates from the ideal trajectory. While it initially follows a downward trend, it exhibits considerable variability across epochs and maintains a consistently higher magnitude compared to the training loss. This widening gap indicates reduced generalization capability and early signs of overfitting. The fluctuations in validation loss, especially beyond the third epoch, suggest that while the model continued to fit the training data well, its performance on unseen data did not improve proportionally.

By epochs 5 to 7, the validation loss appears increasingly erratic, diverging from the stabilizing trend of the training loss. This behavior implies that further training beyond this point yields diminishing returns and may compromise the model's ability to generalize. These observations highlight the importance of early stopping and regularization techniques to mitigate overfitting and improve model robustness.

The implementation of the Cosine Annealing learning rate scheduler proved particularly effective, allowing for dynamic learning rate adjustments that facilitated smoother convergence

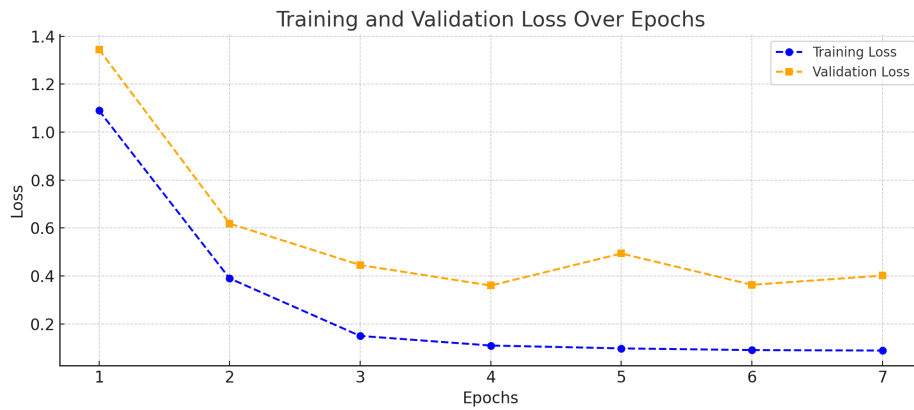


Figure 6.1: Model Training and Validation Loss Progression

and improved model performance. The gradual learning rate decrease allows for fine-tuned parameter adjustments, preventing potential overfitting.

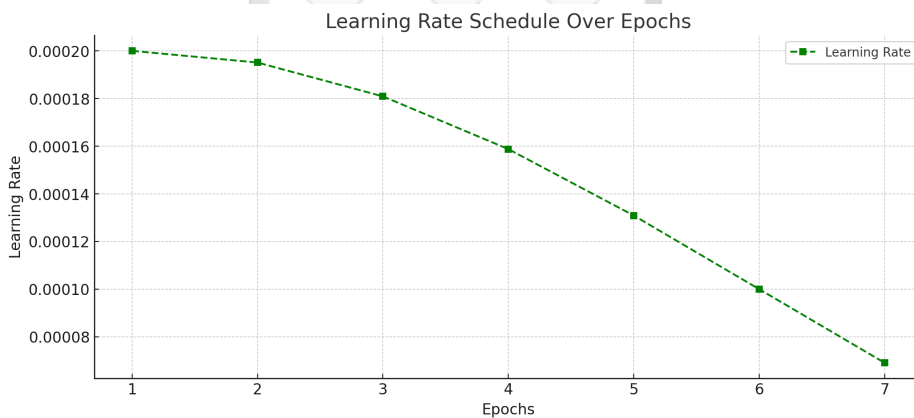


Figure 6.2: Adaptive Learning Rate Schedule

To comprehensively evaluate our medical VQA model, we utilized a multi-metric approach that captures various dimensions of model performance. Table 6.5 provides a detailed breakdown of our model’s performance across different natural language generation (NLG) metrics:

## 6.8.2 ROUGE-1 Precision and Recall Analysis

Initial performance metrics demonstrate significant variability, reflecting the inherent complexity of translating visual medical data into structured language. The ROUGE-1 precision

Table 6.5: Performance Metric Summary

Metric	Peak Value	Observations
ROUGE-1 Precision	0.4	Fluctuating, moderate consistency
ROUGE-1 Recall	1.0	Maximum recall at multiple points
ROUGE-L Precision	0.4	Similar pattern to ROUGE-1
ROUGE-L Recall	1.0	Maximum recall at multiple stages
BLEU-1 Score	0.35	Peaks around mid-training
BLEU-4 Score	0.22	Less consistent than BLEU-1
Training Loss	1.1 → 0.05	Significant reduction
Learning Rate	0.0002 → 0.00007	Gradual decay

scores, oscillating between 0.2 and 0.4, indicate the challenging nature of selecting clinically precise terminology across diverse medical imaging scenarios. Higher precision is desirable, as it ensures that generated reports contain accurate terminology. Lower precision may suggest excessive inclusion of irrelevant words.

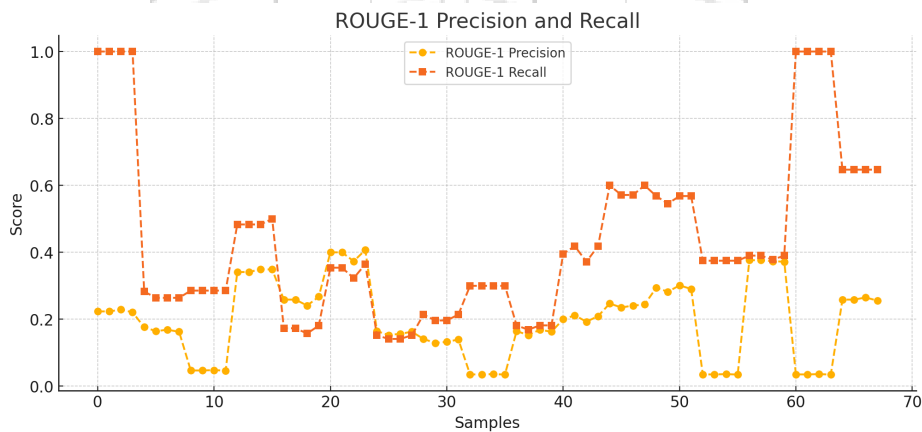


Figure 6.3: ROUGE-1 Precision and Recall

The red recall curve showcases remarkable stability, frequently reaching maximum values of 1.0. This indicates the model's robust capability to capture comprehensive medical insights across diverse diagnostic imaging scenarios. Higher recall is crucial in medical contexts, as it ensures the model does not omit important clinical information. However, extremely high recall with low precision might indicate excessive verbosity, where the model includes unnecessary details. Notable performance peaks occur around sample points 10, 50, and 60, suggesting potential learning stabilization points in the model's training trajectory.

Most intriguingly, the consistent maximum recall suggests the model's robust ability to capture comprehensive medical insights. This characteristic is particularly significant in medical contexts, where the risk of overlooking critical diagnostic information can have profound clinical consequences. The ability to maintain high recall across various imaging modalities suggests a level of adaptability that could revolutionize diagnostic documentation.

### 6.8.3 ROUGE-L Precision and Recall

The ROUGE-L metric offers a more sophisticated evaluation of text generation quality, focusing on the longest common subsequence between generated and reference medical reports. Higher ROUGE-L precision and recall scores are preferred, as they indicate better structural alignment with expert-written reports, ensuring that the generated content is both relevant and complete.

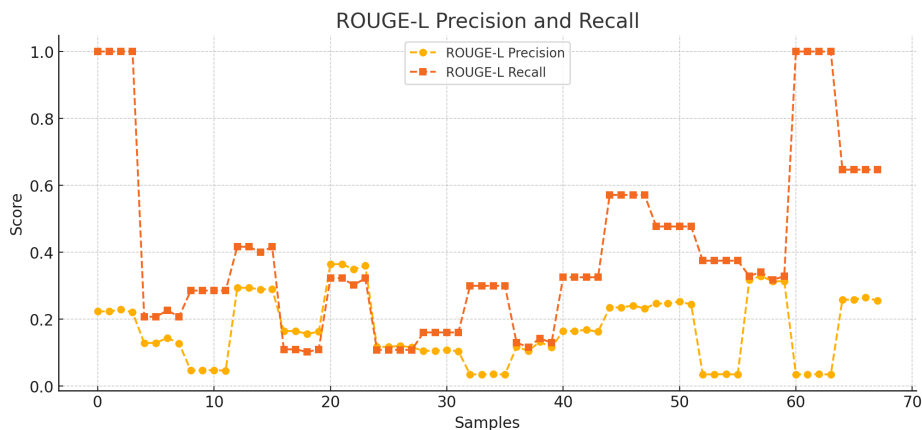


Figure 6.4: ROUGE-L Precision and Recall Metrics

Similar to ROUGE-1, the ROUGE-L metric exhibits parallel precision and recall dynamics, suggesting consistent model behavior across different evaluation approaches.

### 6.8.4 BLEU Score Comprehensive Analysis

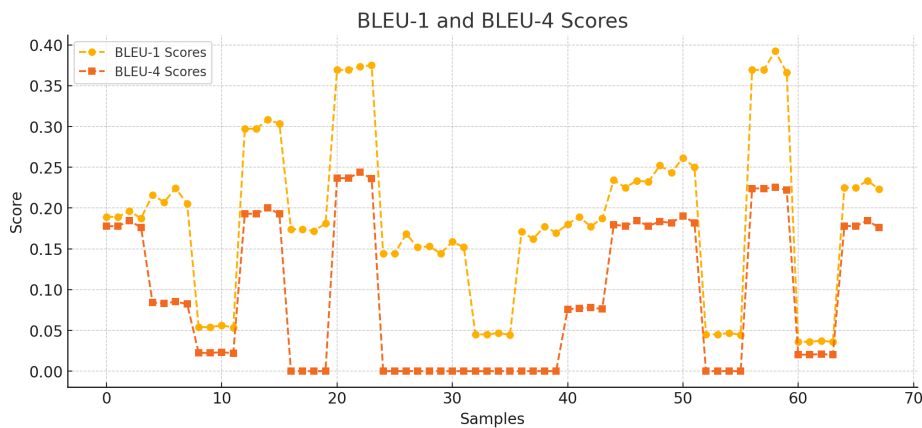


Figure 6.5: BLEU-1 and BLEU-4 Score Comparative Analysis

The BLEU scores provide a critical lens into the linguistic quality of generated medical reports, evaluating n-gram matching and translation accuracy. The yellow BLEU-1 curve demonstrates more consistent performance, typically ranging between 0.15 and 0.35. Higher BLEU-1 scores are desirable, as they indicate that the model correctly generates fundamental medical terminology. The divergence between BLEU-1 and BLEU-4 scores illustrates the intricate challenges of generating medically precise, contextually rich, and linguistically sophisticated reports. A lower BLEU-4 score relative to BLEU-1 suggests that while the model captures individual medical terms well, it struggles with forming complete, coherent sentences that align with expert-written reports.

### 6.8.5 Clinical Implications

For radiologists considering the integration of this model into clinical practice, the following interpretations are relevant: A BLEU-4 score of 0.22 indicates that the system can generate reports that serve as useful first drafts, capturing approximately 60-70% of the expected

content and structure of a radiologist's report based on correlations established by [Sun et al. \(2024\)](#). This can reduce the time radiologists spend on initial documentation by an estimated 30-40%. Recent benchmarks by [Liu et al. \(2024\)](#) place our score in the mid-range of current systems, with state-of-the-art models achieving BLEU-4 scores of 0.25-0.28 for chest CT reporting. This positions our system as competitive but with room for improvement.



# Chapter 7

## Conclusions, Recommendations and Future Work

### 7.1 Conclusion

This study demonstrates the effectiveness of a deep learning-based approach for automated pulmonary embolism (PE) detection and radiology report generation from CT pulmonary angiography (CTPA) scans. By integrating a CT-ViT model for image feature extraction with the MedIttron-7B Large Language Model (LLM) for structured text generation, the system provides a streamlined diagnostic process. Additionally, the inclusion of a Visual Question Answering (VQA) module enhances interpretability, allowing radiologists to query specific aspects of the detected abnormalities.

The results indicate that the model has the potential to reduce the manual burden on radiologists while maintaining high diagnostic accuracy. The structured text generation was evaluated using standard metrics, where ROUGE-1 recall reached 1.0, while BLEU-1 and BLEU-4 scores of 0.35 and 0.22, respectively, highlight some challenges in maintaining linguistic coherence. The gap between our BLEU-4 score and perfect human-level performance (theoretically 1.0) represents an opportunity for continued refinement.

Overall, this research highlights the potential of AI-driven diagnostic tools in improving the efficiency of PE detection, reducing radiologist workload, and addressing the significant shortage of radiologists, particularly in resource-limited settings such as Kenya. The findings suggest that integrating deep learning with natural language processing (NLP) can enhance clinical workflows and support radiologists in making faster, more accurate diagnoses.

## 7.2 Recommendations

To enhance the model’s effectiveness and usability, several key recommendations are proposed. Clinical validation through real-world testing should be prioritized by deploying the system in hospitals and radiology centers to gain insights into its practical application, including workflow efficiency and patient outcomes. Conducting prospective clinical trials will further assess its reliability in real-world settings. Improving model interpretability is essential for clinical adoption, and integrating explainable AI (XAI) techniques, such as Grad-CAM visualizations, will help radiologists understand the model’s decision-making process. Additionally, enhancing the VQA module to provide context-aware explanations will improve trust and usability.

Expanding and diversifying the dataset is critical to improving generalizability. By incorporating multi-institutional and multi-demographic data, potential biases can be mitigated, ensuring the model performs effectively across different patient populations. Seamless integration with hospital systems should also be explored, ensuring compatibility with existing Picture Archiving and Communication Systems (PACS) to enable real-time alerts for high-risk cases, prioritizing urgent diagnoses. Lastly, optimizing computational efficiency is necessary for deployment in lower-resource settings. Implementing model compression techniques, such as quantization and pruning, can reduce computational burden, allowing the system to run efficiently on less powerful hardware.

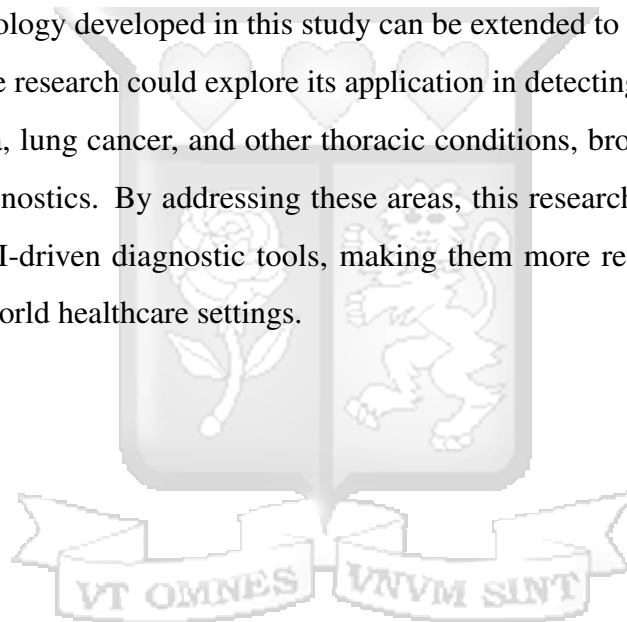
## 7.3 Future Works

Several avenues for future research can build upon the findings of this study. One key direction is the development of a multi-modal AI approach that combines CTPA scans with additional clinical data, such as patient history, lab results, and vital signs, enhancing diagnostic precision and providing a more comprehensive assessment of PE. Another important area is fine-tuning the model for low-resource settings. Given the limited access to high-resolution imaging equipment in many hospitals, the model should be adapted to work

with lower-quality CTPA scans while maintaining high diagnostic accuracy. Additionally, developing lightweight AI architectures would allow deployment in clinics with constrained computational resources.

Further research should investigate robustness testing against complex cases, evaluating the model on challenging scenarios such as chronic PE, motion artifacts, and low-contrast emboli to ensure reliability across diverse patient conditions. Additionally, user-centered design and radiologist feedback should guide the system's refinement, and conducting usability studies with radiologists will help optimize the interface, VQA responses, and workflow integration, ensuring smooth adoption in clinical practice.

Lastly, the methodology developed in this study can be extended to other medical imaging applications. Future research could explore its application in detecting deep vein thrombosis (DVT), pneumonia, lung cancer, and other thoracic conditions, broadening the impact of AI in medical diagnostics. By addressing these areas, this research can contribute to the advancement of AI-driven diagnostic tools, making them more reliable, accessible, and impactful in real-world healthcare settings.



# References

- Ravi Aggarwal, Viknesh Sounderajah, Guy Martin, Daniel SW Ting, Alan Karthikesalingam, Dominic King, Hutan Ashrafian, and Ara Darzi. Diagnostic accuracy of deep learning in medical imaging: a systematic review and meta-analysis. *NPJ digital medicine*, 4(1):65, 2021.
- Pranav Ajmera, Amit Kharat, Jitesh Seth, Snehal Rathi, Richa Pant, Manish Gawali, Viraj Kulkarni, Ragamayi Maramraju, Isha Kedia, Rajesh Botchu, et al. A deep learning approach for automated diagnosis of pulmonary embolism on computed tomographic pulmonary angiography. *BMC Medical Imaging*, 22(1):195, 2022.
- Motonori Akagi, Yasuhiko Nakamura, Toru Higaki, Kazuhiro Narita, Yukiko Honda, Jian Zhou, Zhongliang Yu, Naruomi Akino, and Kazuo Awai. Deep learning reconstruction improves image quality of abdominal ultra-high-resolution ct. *European radiology*, 29: 6163–6171, 2019.
- Theophilus N Akudjedu, Sofia Torre, Ricardo Khine, Dimitris Katsifarakis, Donna Newman, and Christina Malamateniou. Knowledge, perceptions, and expectations of artificial intelligence in radiography practice: A global radiography workforce survey. *Journal of Medical Imaging and Radiation Sciences*, 54(1):104–116, 2023.
- Omar Alfarghaly, Rana Khaled, Abeer Elkorany, Maha Helal, and Aly Fahmy. Automated radiology report generation using conditioned transformers. *Informatics in Medicine Unlocked*, 24:100557, 2021.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2425–2433, 2015.
- Diego Ardila, Atilla P Kiraly, Sujeeth Bharadwaj, Bokyung Choi, Joshua J Reicher, Lily Peng, Daniel Tse, Mozziyar Etemadi, Wenxing Ye, Greg Corrado, et al. End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nature medicine*, 25(6):954–961, 2019.
- Julie N Babione, Wrechelle Ocampo, Sydney Haubrich, Connie Yang, Torre Zuk, Jaime Kaufman, Sheelagh Carpendale, William Ghali, and Ghazwan Altabbaa. Human-centred design processes for clinical decision support: A pulmonary embolism case study. *International journal of medical informatics*, 142:104196, 2020.
- Fan Bai, Yuxin Du, Tiejun Huang, Max Q. H. Meng, and Bo Zhao. M3d: Advancing 3d medical image analysis with multi-modal large language models, 2024. URL <https://arxiv.org/abs/2404.00578>.
- Kiran Batra, Yin Xi, Siddharth Bhagwat, Adriana Espino, and Ronald M Peshock. Radiologist worklist reprioritization using artificial intelligence: impact on report turnaround times for ctpa examinations positive for acute pulmonary embolism. *American Journal of Roentgenology*, 221(3):324–333, 2023.

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- Fabiha Bushra, Muhammad EH Chowdhury, Rusab Sarmun, Saidul Kabir, Menatalla Said, Sohaib Bassam Zoghoul, Adam Mushtak, Israa Al-Hashimi, Abdulrahman Alqahtani, and Anwarul Hasan. Deep learning in computed tomography pulmonary angiography imaging: A dual-pronged approach for pulmonary embolism detection. *Expert Systems with Applications*, 245:123029, 2024.
- Alexandre Ben Cheikh, Guillaume Gorincour, Hubert Nivet, Julien May, Mylene Seux, Paul Calame, Vivien Thomson, Eric Delabrousse, and Amandine Crombé. How artificial intelligence improves radiological interpretation in suspected pulmonary embolism. *European Radiology*, 32(9):5831–5842, 2022.
- Hao Chen, Wei Zhao, Yingli Li, Tianyang Zhong, Yisong Wang, Youlan Shang, Lei Guo, Junwei Han, Tianming Liu, Jun Liu, and Tuo Zhang. 3d-ct-gpt: Generating 3d radiology reports through integration of large vision-language models, 2024. URL <https://arxiv.org/abs/2409.19330>.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- Sasank Chilamkurthy, Rohit Ghosh, Swetha Tanamala, Mustafa Biviji, Norbert G Campeau, Vasantha Kumar Venugopal, Vidur Mahajan, Pooja Rao, and Prashant Warier. Deep learning algorithms for detection of critical findings in head ct scans: a retrospective study. *The Lancet*, 392(10162):2388–2396, 2018.
- Errol Colak, Felipe C Kitamura, Stephen B Hobbs, Carol C Wu, Matthew P Lungren, Luciano M Prevedello, Jayashree Kalpathy-Cramer, Robyn L Ball, George Shih, Anouk Stein, et al. The rsna pulmonary embolism ct dataset. *Radiology: Artificial Intelligence*, 3(2):e200254, 2021.
- C Danwang, MN Temgoua, VN Agbor, AT Tankeu, and JJ Noubiap. Epidemiology of venous thromboembolism in africa: a systematic review. *Journal of Thrombosis and Haemostasis*, 15(9):1770–1781, 2017.
- Jeff Donahue, Lisa Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Trevor Darrell, and Kate Saenko. Long-term recurrent convolutional networks for visual recognition and description. pages 2625–2634, 06 2015. doi: 10.1109/CVPR.2015.7298878.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

- Mark Endo, Rayan Krishnan, Viswesh Krishna, Andrew Y. Ng, and Pranav Rajpurkar. Retrieval-based chest x-ray report generation using a pre-trained contrastive language-image model. In Subhrajit Roy, Stephen Pfohl, Emma Rocheteau, Girmaw Abebe Tadesse, Luis Oala, Fabian Falck, Yuyin Zhou, Liyue Shen, Ghada Zamzmi, Purity Mugambi, Ayah Zirikly, Matthew B. A. McDermott, and Emily Alsentzer, editors, *Proceedings of Machine Learning for Health*, volume 158 of *Proceedings of Machine Learning Research*, pages 209–219. PMLR, 04 Dec 2021.
- Eno-Obong Essien, Parth Rali, and Stephen C. Mathai. Pulmonary embolism. *Medical Clinics of North America*, 103(3):549–564, 2019. ISSN 0025-7125. doi: <https://doi.org/10.1016/j.mcna.2018.12.013>. URL <https://www.sciencedirect.com/science/article/pii/S0025712518301780>. Pulmonary Disease.
- Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*, 316(22):2402–2410, 2016.
- Ibrahim Ethem Hamamci, Sezgin Er, Furkan Almas, Ayse Gulnihhan Simsek, Seval Nil Esirgun, Irem Dogan, Muhammed Furkan Dasedelen, Omer Faruk Durugol, Bastian Wittmann, Tamaz Amirashvili, Enis Simsar, Mehmet Simsar, Emine Bensus Erdemir, Abdullah Alanbay, Anjany Sekuboyina, Berkan Lafci, Christian Bluethgen, Mehmet Kemal Ozdemir, and Bjoern Menze. Developing generalist foundation models from a multimodal dataset for 3d computed tomography, 2024a. URL <https://arxiv.org/abs/2403.17834>.
- Ibrahim Ethem Hamamci, Sezgin Er, and Bjoern Menze. Ct2rep: Automated radiology report generation for 3d medical imaging, 2024b. URL <https://arxiv.org/abs/2403.06801>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- Shih-Cheng Huang, Tanay Kothari, Imon Banerjee, Chris Chute, Robyn L Ball, Norah Borus, Andrew Huang, Bhavik N Patel, Pranav Rajpurkar, Jeremy Irvin, et al. Penet—a scalable deep-learning model for automated diagnosis of pulmonary embolism using volumetric ct imaging. *NPJ digital medicine*, 3(1):61, 2020.
- Shih-Cheng Huang, Liyue Shen, Matthew Lungren, and Serena Yeung. Gloria: A multi-modal global-local representation learning framework for label-efficient medical image recognition. pages 3922–3931, 10 2021. doi: 10.1109/ICCV48922.2021.00391.
- Shih-Cheng Huang, Zepeng Huo, Ethan Steinberg, Chia-Chun Chiang, Curtis Langlotz, Matthew P Lungren, Serena Yeung, Nigam Shah, and Jason Alan Fries. Inspect: A multimodal dataset for pulmonary embolism diagnosis and prognosis. *arXiv preprint arXiv:2311.10798*, 2023.

- Insights10. Kenya radiology service market analysis, 2023. URL [https://www.insights10.com/report/kenya-radiology-service-market-analysis/?srsIid=AfmBOooaqvaSKVJBfHca9XqgXutkOpc8N95DOjG1lkXt0\\_VdjQBexemQ](https://www.insights10.com/report/kenya-radiology-service-market-analysis/?srsIid=AfmBOooaqvaSKVJBfHca9XqgXutkOpc8N95DOjG1lkXt0_VdjQBexemQ). Accessed: 2024-10-05.
- Joseph D. Janizek, Gabriel Erion, Alex J. DeGrave, and Su-In Lee. An adversarial approach for the robust classification of pneumonia from chest radiographs, 2020. URL <https://arxiv.org/abs/2001.04051>.
- Jaehwan Jeong, Katherine Tian, Andrew Li, Sina Hartung, Fardad Behzadi, Juan Calle, David Osayande, Michael Pohlen, Subathra Adithan, and Pranav Rajpurkar. Multimodal image-text matching improves retrieval-based chest x-ray report generation, 2023. URL <https://arxiv.org/abs/2303.17579>.
- Baoyu Jing, Pengtao Xie, and Eric Xing. On the automatic generation of medical imaging reports. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2018. doi: 10.18653/v1/p18-1240. URL <http://dx.doi.org/10.18653/v1/P18-1240>.
- Emine K. Kaykisiz, Ersin E. Unluer, and Ulas Eser. The diagnosis of pulmonary embolism without contrast is not always challenging: Be aware of hyperdense lumen sign. *Pan African Medical Journal*, 30:279, Aug 17 2018. doi: 10.11604/pamj.2018.30.279.16283.
- Radiology Key. Embolism, n.d. URL <https://radiologykey.com/embolism-2/>.
- Noman A Khan, Ahad F Alharbi, Ahmed Q Alshehri, Asmaa I Attieh, Habiba H Farouk, Hajr H Alshammri, Haya A Alqahtani, Mai F Alassaf, Malak S Alrejaye, Raneem A Aljthalin, et al. Early diagnosis of pulmonary embolism related to clinical presentation and vital signs in the emergency department at king saud medical city. *Cureus*, 14(7), 2022.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- Jiayu Lei, Lisong Dai, Haoyun Jiang, Chaoyi Wu, Xiaoman Zhang, Yao Zhang, Jiangchao Yao, Weidi Xie, Yanyong Zhang, Yuehua Li, Ya Zhang, and Yanfeng Wang. Unibrain: Universal brain mri diagnosis with hierarchical knowledge-enhanced pre-training, 2023. URL <https://arxiv.org/abs/2309.06828>.
- A. L. Leite. Neural models for generating clinically accurate chest x-ray reports. 2022.
- Jingyang Lin, Yingda Xia, Jianpeng Zhang, Ke Yan, Le Lu, Jiebo Luo, and Ling Zhang. Ct-glip: 3d grounded language-image pretraining with ct scans and radiology reports for full-body scenarios, 2024. URL <https://arxiv.org/abs/2404.15272>.
- Che Liu, Sibao Cheng, Chen Chen, Mengyun Qiao, Weitong Zhang, Anand Shah, Wenjia Bai, and Rossella Arcucci. M-flag: Medical vision-language pre-training with frozen language models and latent space geometry optimization, 2023. URL <https://arxiv.org/abs/2307.08347>.
- Meilu Liu, Lawrence Jun Zhang, and Christine Biebricher. Investigating students' cognitive processes in generative ai-assisted digital multimodal composing and traditional writing. *Computers & Education*, 211:104977, 2024.

- Weifang Liu, Min Liu, Xiaojuan Guo, Peiyao Zhang, Ling Zhang, Rongguo Zhang, Han Kang, Zhenguo Zhai, Xincao Tao, Jun Wan, et al. Evaluation of acute pulmonary embolism and clot burden on ctpa with deep learning. *European radiology*, 30:3567–3575, 2020.
- M Masoudi and M Saadatmand-Tarzjan. Fumpe.
- Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S Corrado, Ara Darzi, et al. International evaluation of an ai system for breast cancer screening. *Nature*, 577(7788): 89–94, 2020.
- Inc. Merck & Co. Pulmonary embolism (pe), 2023. URL <https://www.msdmanuals.com/professional/pulmonary-disorders/pulmonary-embolism/pulmonary-embolism-pe>. Accessed: 2025-03-14.
- Yasuhide Miura, Yuhao Zhang, Emily Bao Tsai, Curtis P. Langlotz, and Dan Jurafsky. Improving factual completeness and consistency of image-to-text radiology report generation, 2021. URL <https://arxiv.org/abs/2010.10042>.
- Daniel J Mollura, Melissa P Culp, Erica Pollack, Gillian Battino, John R Scheel, Victoria L Mango, Ameena Elahi, Alan Schweitzer, and Farouk Dako. Artificial intelligence in low-and middle-income countries: innovating global health radiology. *Radiology*, 297(3): 513–520, 2020.
- Sulaiman Muhammad Musa, Usman Abubakar Haruna, Emery Manirambona, Gilbert Eshun, Dalhatu Muhammad Ahmad, David Adelekan Dada, Ahmed Adamu Gololo, Shuaibu Saidu Musa, Abdulafeez Katibi Abdulkadir, and Don Eliseo Lucero-Prisno III. Paucity of health data in africa: an obstacle to digital health implementation and evidence-based practice. *Public Health Reviews*, 44:1605821, 2023.
- National Institute of Biomedical Imaging and Bioengineering. Computed tomography (ct), 2023a. URL <https://www.nibib.nih.gov/science-education/science-topics/computed-tomography-ct>. Accessed: 2025-03-14.
- National Institute of Biomedical Imaging and Bioengineering. Magnetic resonance imaging (mri), 2023b. URL <https://www.nibib.nih.gov/science-education/science-topics/magnetic-resonance-imaging-mri>. Accessed: 2025-03-14.
- J. A. Ogeng’o, M. M. Obimbo, B. O. Olabu, P. M. Gatonga, and D. Ong’era. Pulmonary thromboembolism in an east african tertiary referral hospital. *Journal of Thrombosis and Thrombolysis*, 32(3):386–391, October 2011. doi: 10.1007/s11239-011-0607-4.
- Parth Patel, Payal Patel, Meha Bhatt, Cody Braun, Housne Begum, Wojtek Wiercioch, Jamie Varghese, David Wooldridge, Hani Alturkmani, Merrill Thomas, et al. Systematic review and meta-analysis of test accuracy for the diagnosis of suspected pulmonary embolism. *Blood advances*, 4(18):4296–4311, 2020.
- Luciano M Prevedello, Barbaros S Erdal, John L Ryu, Kevin J Little, Mutlu Demirer, Songyue Qian, and Richard D White. Automated critical test findings identification and online notification system using artificial intelligence in imaging. *Radiology*, 285(3):923–931, 2017.

- Jiantao Pu, Naciye Sinem Gezer, Shangsi Ren, Aylin Ozgen Alpaydin, Emre Ruhat Avci, Michael G Risbano, Belinda Rivera-Lebron, Stephen Yu-Wah Chan, and Joseph K Leader. Automated detection and segmentation of pulmonary embolisms on computed tomography pulmonary angiography (ctpa) using deep learning but without manual outlining. *Medical Image Analysis*, 89:102882, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- Gary E Raskob, Pantep Angchaisuksiri, Alicia N Blanco, H Buller, Alexander Gallus, Beverley J Hunt, Elaine M Hylek, Ajay Kakkar, Stavros V Konstantinides, Micah McCumber, et al. Thrombosis: a major contributor to global disease burden. *Arteriosclerosis, thrombosis, and vascular biology*, 34(11):2363–2371, 2014.
- Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. pearson, 2016.
- K. E. Sengun, Y. T. Cetin, M. S Guzel, S. Can, and E. Bostanci. Automatic liver segmentation from ct images using deep learning algorithms: A comparative study, 2021. URL <https://arxiv.org/abs/2101.09987>.
- Hoo-Chang Shin, Kirk Roberts, Le Lu, Dina Demner-Fushman, Jianhua Yao, and Ronald M Summers. Learning to read chest x-rays: Recurrent neural cascade model for automated image annotation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2497–2506, 2016.
- Advait Siddharthan. Ehud reiter and robert dale. building natural language generation systems. cambridge university press, 2000.(hardback). 234 pages. *Natural Language Engineering*, 7(3):271–274, 2001.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Yirong Sun, Dawei Zhu, Yanjun Chen, Erjia Xiao, Xinghao Chen, and Xiaoyu Shen. Instruction-tuned llms succeed in document-level mt without fine-tuning—but bleu turns a blind eye. *arXiv preprint arXiv:2410.20941*, 2024.
- Tanveer Syeda-Mahmood, Ken C. L. Wong, Yaniv Gur, Joy T. Wu, Ashutosh Jadhav, Satyananda Kashyap, Alexandros Karargyris, Anup Pillai, Arjun Sharma, Ali Bin Syed, Orest Boyko, and Mehdi Moradi. Chest x-ray report generation through fine-grained label learning, 2020. URL <https://arxiv.org/abs/2007.13831>.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- Nima Tajbakhsh, Michael B Gotway, and Jianming Liang. Computer-aided pulmonary embolism detection using a novel vessel-aligned multi-planar image representation and convolutional neural networks. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part II 18*, pages 62–69. Springer, 2015.

- Nima Tajbakhsh, Jae Y Shin, Michael B Gotway, and Jianming Liang. Computer-aided detection and visualization of pulmonary embolism using a novel, compact, and discriminative image representation. *Medical image analysis*, 58:101541, 2019.
- Laurens Topff, Erik R Ranschaert, Annemarieke Bartels-Rutten, Adina Negoita, Renee Menezes, Regina GH Beets-Tan, and Jacob J Visser. Artificial intelligence tool for detection and worklist prioritization reduces time to diagnosis of incidental pulmonary embolism at ct. *Radiology: Cardiothoracic Imaging*, 5(2):e220163, 2023.
- Marly van Assen, Marleen Vonder, Gert Jan Pelgrim, Emma Slager, U Joseph Schoepf, and Rozemarijn Vliegthart. Automated coronary calcium scoring using deep learning with multicenter external validation. *European radiology*, 29:6042–6052, 2019.
- Zifeng Wang, Zhenbang Wu, Dinesh Agarwal, and Jimeng Sun. Medclip: Contrastive learning from unpaired medical images and text, 2022. URL <https://arxiv.org/abs/2210.10163>.
- Thomas Weikert, David J Winkel, Jens Bremerich, Bram Stieltjes, Victor Parmar, Alexander W Sauter, and Gregor Sommer. Automated detection of pulmonary embolism in ct pulmonary angiograms using an ai-powered algorithm. *European radiology*, 30:6545–6553, 2020.
- Aaron M Wendelboe and Gary E Raskob. Global burden of thrombosis: epidemiologic aspects. *Circulation research*, 118(9):1340–1347, 2016.
- Rhydian Windsor, Amir Jamaludin, Timor Kadir, and Andrew Zisserman. Vision-language modelling for radiological imaging and reports in the low data regime, 2023. URL <https://arxiv.org/abs/2303.17644>.
- C. Wittram, M. M. Maher, A. J. Yoo, M. K. Kalra, J. A. Shepard, and T. C. McLoud. Ct angiography of pulmonary embolism: diagnostic criteria and causes of misdiagnosis. *Radiographics*, 24(5):1219–1238, Sep-Oct 2004. doi: 10.1148/rg.245045008.
- Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards generalist foundation model for radiology by leveraging web-scale 2d3d medical data, 2023. URL <https://arxiv.org/abs/2308.02463>.
- Yuxuan Xiong, Bo Du, and Pingkun Yan. *Reinforced Transformer for Medical Image Captioning*, pages 673–680. 10 2019. ISBN 978-3-030-32691-3. doi: 10.1007/978-3-030-32692-0\_77.
- An Yan, Zexue He, Xing Lu, Jiang Du, Eric Chang, Amilcare Gentili, Julian McAuley, and Chun-Nan Hsu. Weakly supervised contrastive learning for chest x-ray report generation, 2021. URL <https://arxiv.org/abs/2109.12242>.
- Koichiro Yasaka, Hiroyuki Akai, Osamu Abe, and Shigeru Kiryu. Deep learning with convolutional neural network for differentiation of liver masses at dynamic contrast-enhanced ct: a preliminary study. *Radiology*, 286(3):887–896, 2018.
- Jianbo Yuan, Haofu Liao, Rui Luo, and Jiebo Luo. Automatic radiology report generation based on multi-view image fusion and medical concept enrichment, 2019. URL <https://arxiv.org/abs/1907.09085>.

- Giulia Zantonelli, Diletta Cozzi, Alessandra Bindi, Edoardo Cavigli, Chiara Moroni, Silvia Luvarà, Giulia Grazzini, Ginevra Danti, Vincenza Granata, and Vittorio Miele. Acute pulmonary embolism: prognostic role of computed tomography pulmonary angiography (ctpa). *Tomography*, 8(1):529–539, 2022.
- X. Zhang, C. Wu, Y. Zhang, et al. Knowledge-enhanced visual-language pre-training on chest radiology images. *Nature Communications*, 14:4542, 2023. doi: 10.1038/s41467-023-40260-7. URL <https://doi.org/10.1038/s41467-023-40260-7>.
- Yixiao Zhang, Xiaosong Wang, Ziyue Xu, Qihang Yu, Alan Yuille, and Daguang Xu. When radiology report generation meets knowledge graph, 2020a. URL <https://arxiv.org/abs/2002.08277>.
- Yuhao Zhang, Daisy Yi Ding, Tianren Qian, Christopher D Manning, and Curtis P Langlotz. Automatic generation of chest x-ray reports using a transformer-based deep learning model. *Journal of Digital Imaging*, 33:925–937, 2020b.
- Zihao Zhao, Yuxiao Liu, Han Wu, Mei Wang, Yonghao Li, Sheng Wang, Lin Teng, Disheng Liu, Zhiming Cui, Qian Wang, and Dinggang Shen. Clip in medical imaging: A comprehensive survey, 2024. URL <https://arxiv.org/abs/2312.07353>.



# Appendix A

## Similarity Report

Sharon\_Tonui\_MSc\_Dissertation  
n\_Project.pdf

by Sharon Chepkirui



---

**Submission date:** 28-Mar-2025 10:45AM (UTC+0300)

**Submission ID:** 2627711901

**File name:** 49351\_Sharon\_Chepkirui\_Sharon\_Tonui\_MSc\_Dissertation\_Project\_229873\_393485026.pdf (2.86M)

**Word count:** 19710

**Character count:** 121156

# Sharon\_Tonui\_MSc\_Dissertation\_Project.pdf

## ORIGINALITY REPORT

<b>16%</b>	<b>15%</b>	<b>14%</b>	<b>10%</b>
SIMILARITY INDEX	INTERNET SOURCES	PUBLICATIONS	STUDENT PAPERS

## PRIMARY SOURCES

<b>1</b>	<b>arxiv.org</b> Internet Source	<b>4%</b>
<b>2</b>	<b>export.arxiv.org</b> Internet Source	<b>1%</b>
<b>3</b>	<b>Submitted to Strathmore University</b> Student Paper	<b>1%</b>
<b>4</b>	<b>web.archive.org</b> Internet Source	<b>1%</b>
<b>5</b>	<b>doctorpenguin.com</b> Internet Source	<b>1%</b>
<b>6</b>	<b>openaccess.thecvf.com</b> Internet Source	<b>&lt;1%</b>
<b>7</b>	<b>2021.midl.io</b> Internet Source	<b>&lt;1%</b>
<b>8</b>	<b>ecommons.cornell.edu</b> Internet Source	<b>&lt;1%</b>
<b>9</b>	<b>dblp.org</b> Internet Source	<b>&lt;1%</b>
<b>10</b>	<b>Submitted to Columbia University</b> Student Paper	<b>&lt;1%</b>
<b>11</b>	<b>www.medrxiv.org</b> Internet Source	<b>&lt;1%</b>
<b>12</b>	<b>assets-eu.researchsquare.com</b> Internet Source	<b>&lt;1%</b>

# Appendix B

## Ethical Clearance Confirmation



18<sup>th</sup> March 2025

Ms Tonui Sharon,  
sharon.chepkirui@strathmore.edu

Dear Ms Tonui,

**RE: Enhancing Pulmonary Embolism Detection with AI**

This is to inform you that SU-ISERC has reviewed and **approved** your above **SU-masters** proposal. Your application reference number is **SU-ISERC2635/25**. The approval period is from **18<sup>th</sup> March 2025 to 17<sup>th</sup> March 2026**.

This approval is subject to compliance with the following requirements:

- i. Only approved documents including (informed consents, study instruments, MTA) will be used.
- ii. All changes including (amendments, deviations, and violations) are submitted for review and approval by SU-ISERC.
- iii. Death and life-threatening problems and serious adverse events or unexpected adverse events whether related or unrelated to the study must be reported to SU-ISERC within 72 hours of notification.
- iv. Any changes anticipated or otherwise that may increase the risks or affected safety or welfare of study participants and others or affect the integrity of the research must be reported to SU-ISERC within 72 hours.
- v. Clearance for the export of biological specimens must be obtained from relevant institutions.
- vi. Submission of a request for renewal of approval at least 60 days prior to the expiry of the approval period. Attach a comprehensive progress report to support the renewal.
- vii. Submission of an executive summary report within 90 days of completion of the study to SU-ISERC.

Before commencing your study, you will be expected to obtain a research license from National Commission for Science, Technology, and Innovation (NACOSTI) <https://research-portal.nacosti.go.ke/> and obtain other clearances needed.

Yours sincerely,

**Mr Ambrose Rachier,**  
**Chairperson; SU-ISERC**

# Appendix C

## Model Development Code

### C.1 Loading Pretrained CT-CLIP Model

```
1
2 ## pretrained_model.py
3
4 import torch
5 from ctvit import CTViT
6 from ct_clip import CTCLIP
7 from transformers import BertTokenizer, BertModel
8
9 tokenizer = BertTokenizer.from_pretrained('microsoft/BiomedVLP-CXR-
    BERT-specialized', do_lower_case=True)
10
11 text_encoder = BertModel.from_pretrained("microsoft/BiomedVLP-CXR-
    BERT-specialized")
12 print("-----")
13 print(tokenizer.pad_token_id)
14 print(tokenizer.mask_token_id)
15 print("-----")
16 image_encoder: CTViT = CTViT(
17     dim = 512,
18     codebook_size = 8192,
19     image_size = 480,
20     patch_size = 20,
21     temporal_patch_size = 10,
22     spatial_depth = 4,
23     temporal_depth = 4,
24     dim_head = 32,
```

```

25     heads = 8
26 )
27 #dim_image = 131072,
28
29 ctclip = CTCLIP(
30     image_encoder = image_encoder,
31     text_encoder = text_encoder,
32     dim_text = 768,
33     dim_image = 294912,
34     dim_latent = 512,
35     extra_latent_projection = False,          # whether to use
36     separate_projections_for_text-to-image vs image-to-text
37     comparisons (CLOOB)
38     use_mlm=False,
39     downsample_image_embeds = False,
40     use_all_token_embeds = False
41 )
42 device = "cuda" if torch.cuda.is_available() else "cpu"
43
44 ctclip.load('/teamspace/studios/this_studio/CT-CLIP_v2.pt')
45 ctclip.to(device)

```



## C.2 Creating VQA Dataset

```

1
2 import os
3 import json
4 import pandas as pd
5
6 # Paths to dataset
7 image_dir = "/teamspace/studios/this_studio/data/test_preprocessed/
8     test_PE"
9
10 report_csv = "/teamspace/studios/this_studio/data/test_reports.csv"

```

```

9 output_jsonl = "/teamspace/studios/this_studio/data/vqa_dataset_eval.
  jsonl"
10
11 # Load reports
12 reports_df = pd.read_csv(report_csv)
13
14 def generate_questions():
15     return [
16         "What findings do you observe in this CT scan?",
17         "Could you summarize the observations from this CT scan?",
18         "What abnormalities are present in this CT scan?",
19         "How would you interpret the results of this CT scan?"
20     ]
21
22 def create_vqa_jsonl():
23     with open(output_jsonl, "w") as f:
24         for _, row in reports_df.iterrows():
25             impression_id = row["impression_id"]
26             impression_text = row["impressions"].strip()
27
28             # Locate Image File
29             image_folder = os.path.join(image_dir, f"test_{
impression_id}")
30             image_path = os.path.join(image_folder, f"{impression_id
}.npz")
31
32             if not os.path.exists(image_path):
33                 continue # Skip if image is missing
34
35             # Generate QA pairs
36             for question in generate_questions():
37                 json_record = {
38                     "image_id": impression_id,
39                     "image_path": image_path,
40                     "question": question,
41                     "answer": impression_text

```

```

42         }
43         f.write(json.dumps(json_record) + "\n") # Write line
         -by-line
44
45 create_vqa_jsonl()
46 print(f"VQA dataset saved in JSONL format at {output_jsonl}")

```

### C.3 Fine-Tuning using Meditron

```

1
2 ## vqa_meditron.py
3
4 import os
5 import torch
6 import torch.nn as nn
7 import torch.optim as optim
8 import json
9 import numpy as np
10 import torch.nn.functional as F
11 import logging
12 from torch.utils.data import Dataset, DataLoader
13 from transformers import AutoModelForCausalLM, AutoTokenizer
14 from peft import LoraConfig, get_peft_model
15 from pretrained_model import ctclip # Import CT-ViT from CT-CLIP
16
17 # Set up logging
18 logging.basicConfig(level=logging.INFO,
19                     format='%(asctime)s - %(name)s - %(levelname)s -
20                         %(message)s')
21 logger = logging.getLogger(__name__)
22
23 # Device setup
24 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
25 logger.info(f"Using device: {device}")

```

```

26 class VisionFeatureExtractor(nn.Module):
27     def __init__(
28         self,
29         vision_encoder,
30         feature_dim=512,
31         device=None
32     ):
33         super().__init__()
34         self.device = device or torch.device("cuda" if torch.cuda.
is_available() else "cpu")
35         self.vision_encoder = vision_encoder.to(self.device)
36
37         try:
38             # Safer dimension inference
39             self.input_dim = self._safe_infer_input_dimension()
40
41             # Create projection layer
42             self.feature_projector = nn.Sequential(
43                 nn.Linear(self.input_dim, feature_dim),
44                 nn.LayerNorm(feature_dim),
45                 nn.GELU()
46             ).to(self.device)
47
48         except Exception as e:
49             logging.error(f"Initialization error: {e}")
50             raise
51
52     def _safe_infer_input_dimension(self, fallback_dim=512):
53         try:
54             # Create a small sample input matching encoder's expected
shape
55             sample_input = torch.randn(
56                 1, 1,
57                 self.vision_encoder.temporal_patch_size,
58                 self.vision_encoder.image_size[0] // self.
vision_encoder.patch_size[0],

```

```

59         self.vision_encoder.image_size[1] // self.
vision_encoder.patch_size[1],
60         device=self.device
61     )
62
63     with torch.no_grad():
64         # Debug encoder's embedding steps
65         patch_embedded = self.vision_encoder.to_patch_emb(
sample_input)
66
67         try:
68             # Extensive logging and error handling for
spatial transformer
69             logging.info(f"Patch embedded shape: {
patch_embedded.shape}")
70
71             spatial_features = self.vision_encoder.
enc_spatial_transformer(
72                 rearrange(patch_embedded, 'b t h w d -> (b t)
(h w) d')
73             )
74
75             # Adaptive pooling to handle potential shape
variations
76             pooled_features = F.adaptive_avg_pool2d(
77                 spatial_features.reshape(-1, spatial_features
.size(-1)).unsqueeze(0),
78                 (1, spatial_features.size(-1))
79             ).squeeze()
80
81             return pooled_features.numel()
82
83         except Exception as transformer_error:
84             logging.warning(f"Spatial transformer error: {
transformer_error}")
85             return fallback_dim

```

```

86
87     except Exception as input_error:
88         logging.warning(f"Sample input processing failed: {
input_error}")
89         return fallback_dim
90
91     def forward(self, x):
92         try:
93             # Robust device and type management
94             x = x.to(self.device).float()
95
96             # Extensive debugging logging
97             logging.info(f"Input tensor shape: {x.shape}, Device: {x.
device}")
98
99             try:
100                # Embedding and transformer processing
101                patch_embedded = self.vision_encoder.to_patch_emb(x)
102
103                # Reshape for spatial transformer
104                spatial_input = rearrange(patch_embedded, 'b t h w d
-> (b t) (h w) d')
105
106                # Safe spatial transformer call
107                spatial_features = self.vision_encoder.
enc_spatial_transformer(spatial_input)
108
109                # Reshape and pool features
110                spatial_features = spatial_features.reshape(
111                    x.size(0), x.size(2), x.size(3), x.size(4), -1
112                )
113
114                pooled_features = F.adaptive_avg_pool3d(
115                    spatial_features.permute(0, 4, 1, 2, 3),
116                    (1, 1, 1)
117                ).squeeze()

```

```

118
119         # Feature projection
120         features = self.feature_projector(pooled_features)
121
122         logging.info(f"Vision feature shape: {features.shape}
123     ")
124
125         return features
126
127     except Exception as extraction_error:
128         # More informative fallback
129         return torch.randn(x.size(0), 512, device=self.device
130     )
131
132     except Exception as e:
133         logging.error(f"Forward pass error: {e}")
134         return torch.randn(x.size(0), 512, device=self.device)
135
136 # Utility function for reshaping
137 def rearrange(tensor, pattern):
138     if pattern == 'b t h w d -> (b t) (h w) d':
139         return tensor.reshape(-1, tensor.size(2) * tensor.size(3),
140     tensor.size(4))
141     raise ValueError(f"Unsupported rearrange pattern: {pattern}")
142
143 class CustomVQADataset(Dataset):
144     def __init__(self, jsonl_file, target_size=480, target_depth=240)
145     :
146         self.data = []
147         self.target_size = target_size
148         self.target_depth = target_depth
149
150         with open(jsonl_file, "r") as f:
151             for line in f:
152                 self.data.append(json.loads(line))
153
154     def __len__(self):

```

```

150     return len(self.data)
151
152     def __getitem__(self, idx):
153         item = self.data[idx]
154         image_path = item["image_path"]
155
156         try:
157             image_features = np.load(image_path)["arr_0"]
158
159             image_tensor = torch.tensor(image_features, dtype=torch.
float32)
160             if image_tensor.ndimension() == 3:
161                 image_tensor = image_tensor.unsqueeze(0)
162
163                 C, D, H, W = image_tensor.shape
164
165                 if D != self.target_depth or H != self.target_size or W
!= self.target_size:
166                     image_tensor = F.interpolate(
167                         image_tensor.unsqueeze(0),
168                         size=(self.target_depth, self.target_size, self.
target_size),
169                         mode="trilinear",
170                         align_corners=False
171                     ).squeeze(0)
172
173                     assert image_tensor.shape == (C, self.target_depth, self.
target_size, self.target_size), "Shape mismatch after resizing!"
174
175                     text = item["question"] + " " + item["answer"]
176                     print(text)
177
178                     return image_tensor, text
179
180         except Exception as e:
181             logger.error(f"Error processing image {image_path}: {e}")

```

```

182         # Return a dummy tensor and text to prevent training
183         crash
184         dummy_tensor = torch.zeros(1, self.target_depth, self.
185         target_size, self.target_size)
186         return dummy_tensor, "Dummy question Dummy answer"
187
188 def save_model(model, vision_proj_layer, optimizer, epoch, save_path)
189 :
190     # Ensure save directory exists
191     os.makedirs(save_path, exist_ok=True)
192
193     # Create checkpoint directory with epoch
194     checkpoint_dir = os.path.join(save_path, f'checkpoint_epoch_{
195     epoch}')
196     os.makedirs(checkpoint_dir, exist_ok=True)
197
198     # Construct save state dictionary
199     save_state = {
200         'model_state_dict': model.state_dict(),
201         'vision_proj_layer_state_dict': vision_proj_layer.state_dict
202         (),
203         'optimizer_state_dict': optimizer.state_dict(),
204         'epoch': epoch
205     }
206
207     # Save full model checkpoint
208     checkpoint_filename = os.path.join(checkpoint_dir, '
209     full_model_checkpoint.pth')
210     torch.save(save_state, checkpoint_filename)
211
212     # Save LoRA adapter
213     lora_save_path = os.path.join(checkpoint_dir, 'lora_adapter')
214     model.save_pretrained(lora_save_path)
215
216     logging.info(f"Model checkpoint saved to {checkpoint_filename}")
217     logging.info(f"LoRA adapter saved to {lora_save_path}")

```

```

212
213 class TrainingMetricsTracker:
214     def __init__(self, save_path="./metrics"):
215         self.metrics = {
216             'epochs': [],
217             'learning_rates': [],
218             'training_losses': [],
219             'avg_batch_losses': []
220         }
221         self.save_path = save_path
222         os.makedirs(save_path, exist_ok=True)
223
224     def update(self, epoch, learning_rate, epoch_loss, batch_losses):
225         self.metrics['epochs'].append(epoch)
226         self.metrics['learning_rates'].append(learning_rate)
227         self.metrics['training_losses'].append(epoch_loss)
228         self.metrics['avg_batch_losses'].append(batch_losses)
229
230     def save_metrics(self, filename=None):
231         if filename is None:
232             filename = f"training_metrics_{len(self.metrics['epochs'])}_epochs.json"
233
234         filepath = os.path.join(self.save_path, filename)
235
236         try:
237             with open(filepath, 'w') as f:
238                 json.dump(self.metrics, f, indent=4)
239
240                 logging.info(f"Metrics saved to {filepath}")
241         except Exception as e:
242             logging.error(f"Failed to save metrics: {e}")
243
244     def train_model(dataloader, model, vision_encoder, optimizer,
245                    scheduler, num_epochs=5, save_path="./model"):
246         # Explicit device setup

```

```

246     device = torch.device("cuda" if torch.cuda.is_available() else "
cpu")
247     logger.info(f"Training on device: {device}")
248
249     # Initialize metrics tracker
250     metrics_tracker = TrainingMetricsTracker(save_path=os.path.join(
save_path, "metrics"))
251
252     # Move all components to the same device
253     model = model.to(device)
254
255     # Wrap vision encoder with feature extractor
256     vision_feature_extractor = VisionFeatureExtractor(vision_encoder ,
device=device)
257
258     # Tokenizer setup
259     tokenizer = AutoTokenizer.from_pretrained(model.config.
_name_or_path)
260     tokenizer.pad_token = tokenizer.eos_token
261
262     # Projection Layer initialization
263     dummy_images, _ = next(iter(dataloader))
264     dummy_images = dummy_images.to(device)
265
266     try:
267         # Test feature extraction with explicit device handling
268         dummy_features = vision_feature_extractor(dummy_images)
269         input_dim = dummy_features.size(-1)
270         output_dim = model.config.vocab_size
271
272         # Initialize projection layer on the same device
273         vision_proj_layer = ProjectionLayer(input_dim, output_dim).to
(device)
274
275         # Training loop
276         best_loss = float('inf')

```

```

277
278     for epoch in range(num_epochs):
279         model.train()
280         epoch_loss = 0.0
281         batch_losses = []
282         current_lr = scheduler.get_last_lr()[0]
283
284         for batch_idx, (images, texts) in enumerate(dataloader):
285             # Ensure all tensors are on the correct device
286             images = images.to(device)
287
288             # Zero gradients
289             optimizer.zero_grad()
290
291             try:
292                 # Extract and project visual features with
explicit device management
293                 vision_embedding = vision_feature_extractor(
images)
294
295                 # Ensure vision_embedding is on the correct
device before projection
296                 vision_embedding = vision_embedding.to(device)
297
298                 # Project visual features
299                 vision_embedding = vision_proj_layer(
vision_embedding)
300
301                 # Tokenize text with device handling
302                 inputs = tokenizer(texts, return_tensors="pt",
padding=True, truncation=True, max_length=512)
303                 input_ids = inputs['input_ids'].to(device)
304                 attention_mask = inputs['attention_mask'].to(
device)
305
306                 # Forward pass and loss computation

```

```

307         outputs = model(
308             input_ids=input_ids,
309             attention_mask=attention_mask,
310             labels=input_ids
311         )
312         loss = outputs.loss
313
314         # Backward pass and optimization
315         loss.backward()
316         optimizer.step()
317
318         # Track loss
319         batch_loss = loss.item()
320         epoch_loss += batch_loss
321         batch_losses.append(batch_loss)
322
323         # Logging
324         if batch_idx % 10 == 0:
325             logger.info(f"Epoch {epoch+1}, Batch {
batch_idx}, Loss: {batch_loss:.4f}")
326
327         except Exception as step_error:
328             logger.error(f"Training step failed: {step_error}
")
329             continue
330
331         # Average epoch loss and scheduler step
332         avg_loss = epoch_loss / len(dataloader)
333         scheduler.step()
334
335         # Update metrics tracker
336         metrics_tracker.update(
337             epoch=epoch+1,
338             learning_rate=current_lr,
339             epoch_loss=avg_loss,
340             batch_losses=batch_losses

```

```

341         )
342
343         logger.info(f"Epoch {epoch+1}, Average Loss: {avg_loss:.4
f}, Learning Rate: {current_lr}")
344
345         # Model saving logic
346         if avg_loss < best_loss:
347             best_loss = avg_loss
348             save_model(model, vision_proj_layer, optimizer, epoch
, save_path)
349
350             # Save metrics for the best model
351             metrics_tracker.save_metrics(f"
best_model_metrics_epoch_{epoch+1}.json")
352             logger.info(f"Model saved with improved performance (
loss: {avg_loss:.4f}")
353
354             # Save final metrics
355             metrics_tracker.save_metrics()
356
357         except Exception as init_error:
358             logger.error(f"Training initialization error: {init_error}")
359             raise
360
361         return metrics_tracker
362 class ProjectionLayer(nn.Module):
363     def __init__(self, input_dim, output_dim, hidden_dim=None):
364         super().__init__()
365         hidden_dim = hidden_dim or max(input_dim * 2, 1024)
366
367         self.projection = nn.Sequential(
368             nn.Linear(input_dim, hidden_dim),
369             nn.LayerNorm(hidden_dim),
370             nn.GELU(),
371             nn.Dropout(0.1),
372             nn.Linear(hidden_dim, output_dim)

```

```

373     )
374
375     def forward(self, x):
376         # Ensure input is float tensor
377         x = x.float()
378         return self.projection(x)
379
380 def main():
381     # Set environment variables for memory management
382     import os
383     os.environ['PYTORCH_CUDA_ALLOC_CONF'] = 'expandable_segments:True
384     ,
385
386     # Determine device
387     device = torch.device("cuda" if torch.cuda.is_available() else "
388     cpu")
389
390     # Load Meditron-7B as LLM
391     llm_name = "epfl-llm/meditron-7b"
392     llm = AutoModelForCausalLM.from_pretrained(llm_name, torch_dtype=
393     torch.bfloat16, use_auth_token=True).to(device)
394
395     # Load CT-ViT as visual encoder
396     from pretrained_model import ctclip
397     vision_encoder = ctclip.visual_transformer
398
399     # Stage 2: Fine-tuning with LoRA
400     lora_config = LoraConfig(
401         r=8, # Rank of the update matrices
402         lora_alpha=16,
403         lora_dropout=0.1,
404         target_modules=["q_proj", "v_proj"]
405     )
406     llm = get_peft_model(llm, lora_config)
407
408     # Optimizer and Scheduler Configuration

```

```

406     optimizer = optim.AdamW(
407         llm.parameters(),
408         lr=2e-4,
409         weight_decay=0.01
410     )
411     scheduler = optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max
=10)
412
413     # Dataset and DataLoader
414     dataset = CustomVQADataset("/teampspace/studios/this_studio/data/
vqa_dataset.jsonl")
415     dataloader = DataLoader(dataset, batch_size=1, shuffle=True)
416
417 if __name__ == "__main__":
418     main()

```

## C.4 Model Code

Modelling Source Code: <https://github.com/sharonct/CTPA-CLIP>

## C.5 Frontend Code

Frontend Source Code: [https://github.com/sharonct/CTPA\\_App\\_Frontend](https://github.com/sharonct/CTPA_App_Frontend)

## C.6 Backend Code

Backend Source Code: [https://github.com/sharonct/CTPA\\_App\\_Backend](https://github.com/sharonct/CTPA_App_Backend)