



Daffodil
International
University

Automating Clinical Note Summarization Using LLM

Submitted By

Md. Yeasin Arafat

Section: A

ID: 211-35-687

Department of Software Engineering

Daffodil International University

Supervised By

Md. Shohel Arman

Assistant Professor

Department of Software Engineering

Daffodil International University

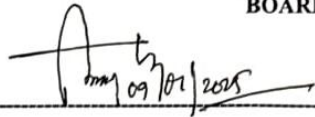
Thesis submitted in fulfillment of the requirements for the award of the degree of Bachelor of
Science

Fall - 2024

APPROVAL

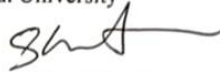
This thesis titled no “Automating Clinical Note Summarization Using Large Language Model (LLM)”, submitted by **Md. Yeasin Arafat (ID: 211-35-687)** to the Department of Software Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of Bachelor of Science in Software Engineering and approval as to its style and contents.

BOARD OF EXAMINERS



Professor Dr. Engr. AKM Masum
Professor
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University

Chairman



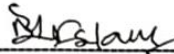
Md. Shohel Arman
Assistant Professor
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University

Internal Examiner 1



Dr. Marzia Ahmed
Assistant Professor
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University

Internal Examiner 2



Dr. Md. Monowarul Islam
Associate Professor
Department of Computer Science & Engineering
Jagannath University

External Examiner



SUPERVISOR'S DECLARATION

I hereby declare that I have checked this thesis and in my opinion, this thesis is adequate in terms of scope and quality for the award of the degree of Bachelor of Science.

A handwritten signature in black ink, appearing to be "SHA", written over a horizontal line.

(Supervisor's Signature)

Full Name : Md. Shohel Arman
Position : Assistant Professor
Date : 13 January, 2025



STUDENT'S DECLARATION

I hereby declare that the work in this thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at Daffodil International University or any other institution.

Yeasin

(Student's Signature)

Full Name : Md. Yeasin Arafat

ID Number : 211-35-687

Date : 04 January, 2025

ACKNOWLEDGEMENT

First and foremost, I want to express my sincere gratitude to Almighty Allah for providing me with the courage, discernment, and tenacity necessary to finish this research. I am appreciative of my parents' unwavering love, support, and encouragement throughout my academic career. I have always been most inspired and motivated by their faith in me.

I want to express my gratitude to Assistant Professor Md. Shohel Arman, my supervisor, for his insightful counsel, encouragement, and direction during the study. His wisdom and understanding have greatly influenced this work. I am also very grateful to Dr. Imran Mahmud, the department director, for his encouragement, direction, and insightful remarks that enabled me to finish my journey successfully. Lastly, I want to express my gratitude to all of my friends, coworkers, and anyone else that supported and encouraged me throughout this process.

ABSTRACT

The increasing volume of electronic health records (EHR) and the growing complexity of clinical documentation have highlighted the need for efficient and reliable summarization tools. This thesis explores the application of large language models (LLMs) in automating clinical note summarization, with the goal of reducing the cognitive and administrative burden on healthcare professionals while maintaining clinical accuracy and relevance. Leveraging state-of-the-art LLM architectures such as T5 and FLAN-T5, this research focuses on extracting key information from clinical notes, including patient diagnoses, treatment plans, and medical histories, and generating concise, structured summaries suitable for clinical workflows.

The study evaluates the performance of fine-tuned LLMs on datasets such as MIMIC-III, MeQSum, and Probsum using quantitative metrics like ROUGE and BERTScore, achieving a ROUGE F1 score of 0.95 and demonstrating high efficiency with a runtime of 2.35 seconds per note. Qualitative analysis confirms the generated summaries' clinical relevance, with outputs aligned to standard sections like diagnoses, imaging results, and treatments. Despite these strengths, challenges such as occasional hallucinated information, omitted secondary details, and inconsistent formatting are identified.

Results from physician feedback underscore the practicality of LLMs in improving healthcare documentation, with models like LLaMA-Clinic achieving over 90% acceptance in a blinded review. Additionally, cost analysis reveals a 3.75-fold reduction in inference costs compared to proprietary alternatives, emphasizing the feasibility of open-source solutions. This research highlights the potential of LLMs to enhance clinical workflows, reduce diagnostic errors, and improve patient care.

Future directions include expanding datasets for better generalizability, addressing hallucination issues, and ensuring seamless integration into healthcare systems through ethical and clinician-centered approaches. The findings reinforce the role of AI in advancing healthcare, promoting accessibility, and addressing global challenges in medical documentation and decision-making.

TABLE OF CONTENT

ACKNOWLEDGEMENT	iv
ABSTRACT	v
LIST OF FIGURES	vii
LIST OF TABLES	viii
CHAPTER 1	1
INTRODUCTION	1
1.1 Introduction	1
1.2 Background	1
1.3 Motivation of the Research	2
1.4 Problem Statement	2
1.5 Research Questions	2
1.6 Research Objectives	3
1.7 Research Scope	3
CHAPTER 2	4
LITERATURE REVIEW	4
[11].	15
CHAPTER 3	18
RESEARCH METHODOLOGY	18
3.1 Introduction	18
3.2 Dataset	19
3.3 Model Selection	19
3.4 Fine-Tuning Pipeline	20
3.5 Evaluation, Optimization & Deployment	21
3.6 Software Architecture	21
3.6.1 Introduction	21
3.6.2 User Interface	22
3.6.3 Backend Server	23
3.6.4 Fine-Tuned Model	23
3.6.5 Evaluation and Feedback	23
3.6.6 Output	23
3.6.7 Summary	24

3.7 Expert Oversight	24
CHAPTER 4	25
RESULT AND DISCUSSION	25
4.1 Introduction	25
4.2 Quantitative Analysis	25
4.2.1 Performance Metrics	25
4.2.2 Comparison with Other LLMs	26
4.2.3 Strengths of the Proposed Model	27
4.3 Qualitative Analysis	28
4.3.1 Named Entities Recognition	29
4.4 Limitations	31
4.5 Discussion	32
4.5.1 Accuracy and Efficiency	32
4.5.2 Clinical Relevance	32
4.5.3 Strengths	32
4.5.4 Challenges	33
4.6 Summary	33
CHAPTER 5	34
CONCLUSION AND RECOMMENDATION	34
REFERENCES	35

LIST OF FIGURES

- Figure 3.1: Workflow Diagram
- Figure 3.2: LLMs Working Process
- Figure 3.3: Pipeline Diagram
- Figure 3.4: Evaluation, Optimization & Deploy
- Figure 3.5: Software Architecture
- Figure 3.6: Clinical Reader
- Figure 4.1: ROUGH -1 Scores
- Figure 4.2: BLEU Scores
- Figure 4.3: Named Entity Recognition
- Figure 4.4: Success Ratio Input & Generated Summary

Figure 4.5: Confusion Matrix (Proposed Model)

LIST OF TABLES

Table 3.1: Dataset Description & Instructions

Table 4.1: Performance Metrics

Table 4.2: Metric Comparison with LLMs

CHAPTER 1

INTRODUCTION

1.1 Introduction

The rapid digitalization of healthcare has led to an exponential increase in the volume of electronic health records (EHRs). While these records are critical for patient care and research, their extensive and unstructured nature often overwhelms clinicians. The process of documenting and reviewing patient information, particularly in high-pressure environments like intensive care units, detracts from patient-centered care and contributes to clinician burnout. Automating the summarization of clinical notes has emerged as a vital solution to streamline medical workflows and optimize healthcare delivery.

Recent advancements in Natural Language Processing (NLP) and Large Language Models (LLMs), such as GPT-4 and BERT-based models, have shown remarkable capabilities in generating coherent, concise, and contextually accurate summaries of medical text. These models have been successfully adapted to various clinical NLP tasks, including medical question-answering, EHR generation, and diagnostic reasoning [4, 9]. However, their application in clinical note summarization is particularly promising due to their ability to process and synthesize vast textual data into actionable insights.

Automating clinical note summarization is not without challenges. Clinical texts are characterized by domain-specific terminologies, abbreviations, and complex structures that require precise contextual understanding. Additionally, the variability in documentation styles and the necessity for maintaining factual accuracy make this task especially complex. Abstractive summarization approaches, which aim to generate novel text rather than merely extracting phrases, have demonstrated potential in addressing these challenges by leveraging LLMs trained on extensive medical datasets such as MIMIC-III and PubMed [5, 2].

This work explores the integration of LLMs in automating the summarization of clinical notes, focusing on optimizing workflows while ensuring high accuracy and interpretability. By addressing the nuances of medical language and leveraging advanced model architectures, we aim to provide a robust framework that enhances the efficiency of healthcare professionals, allowing them to prioritize patient care over administrative burdens.

1.2 Background

Clinical note summarization is transforming healthcare by tackling the inefficiencies and complexities of medical documentation. It condenses lengthy, disorganized notes—covering patient histories, diagnoses, treatments, and follow-ups—into clear, concise summaries that focus on the most critical information. This not only reduces the cognitive burden on clinicians but also allows them to quickly access actionable insights, make better decisions, and improve communication within multidisciplinary medical teams. By cutting down the time spent on redundant or unstructured data, these tools enable healthcare professionals to prioritize patient care while reducing the risk of errors. Using advanced large language models (LLMs) like GPT and FLAN-T5, which are adept at processing text and understanding context, clinical note summarization offers a scalable way to boost efficiency, ensure accuracy, and elevate the quality of care in dynamic medical settings [1, 2].

1.3 Motivation of the Research

A clinical reader study with ten physicians evaluates completeness, correctness, and conciseness; in the majority of cases, summaries from our best-adapted LLMs are either equivalent (45%) or superior (36%) compared to summaries from medical experts [1]. This capability not only alleviates the workload of healthcare providers but also ensures that critical information is readily available, improving coordination and accelerating decision-making. With growing patient volumes and a global shortage of healthcare staff, technologies like LLMs become essential for enhancing productivity, reducing errors, and ensuring better patient outcomes [4].

1.4 Problem Statement

Clinical practitioners often struggle to balance their limited time with the burden of bureaucratic duties, compounded by the complexity of clinical documentation. Moreover, the recent adoption of electronic health records (EHR) generates large amounts of unstructured and redundant information that require advanced tools for efficient management. Manual medical documentation remains inefficient and prone to errors, creating a need for automated summarization systems that can integrate seamlessly into clinical workflows while maintaining linguistic and clinical relevance.

1.5 Research Questions

⇒ How well can LLMs process and distill clinical notes while ensuring both accuracy and relevance in a medical context?

1.6 Research Objectives

- To develop an LLM-powered system designed to create clear, precise, and contextually appropriate summaries of clinical notes.
- To ensure that the summaries generated by the system fulfill the specific requirements of clinical practice.
- To provide practical guidelines for incorporating automated summarization tools into healthcare systems while maintaining compliance with ethical and privacy regulations.

1.7 Research Scope

This study focuses on the use of state-of-the-art LLMs, such as GPT, BioBERT, and MedPaLM, for summarizing clinical notes. The research emphasizes bridging the gap between generic language models and the specific requirements of medical applications by exploring both extractive and abstractive summarization techniques. It aims to ensure that domain-specific fine-tuning and reconstruction fulfill the standards and needs of clinical practice.

CHAPTER 2

LITERATURE REVIEW

Study delves into the potential of large language models (LLMs) to handle clinical text summarization tasks, comparing their performance to that of experienced medical professionals. The researchers evaluated eight different LLMs, including GPT-4, across key areas such as radiology reports, patient questions, progress notes, and doctor-patient dialogues. They used two main adaptation techniques—In-Context Learning (ICL) and Quantized Low-Rank Adaptation (QLoRA)—to tailor these models to the specific needs of clinical summarization. To assess how well the models performed, the authors combined traditional metrics like BLEU, ROUGE-L, and BERTScore with a specialized medical evaluation metric called MEDCON. Additionally, 10 physicians participated in a detailed reader study to judge the summaries based on accuracy, completeness, and how well they captured key details without being overly verbose. The findings were eye-opening. GPT-4, when adapted with ICL, managed to equal or even outperform medical experts in 81% of cases. Specifically, it outperformed human professionals in 36% of cases and performed on par with them in 45%. This result highlights how effective these adaptation techniques can be in enabling general-purpose LLMs to handle highly specialized tasks. One of the model's strengths was its ability to condense information into clear, concise summaries while still preserving the critical clinical details. However, a closer look at the outputs revealed that while the LLMs generated fewer hallucinations and clinically dangerous inaccuracies than humans, they still struggled with ambiguous or indirect language, which remains a challenge in medical contexts. The study emphasizes the exciting potential of LLMs to ease the documentation burden faced by healthcare professionals. By automating parts of the summarization process, these tools could help reduce burnout and allow clinicians to focus more on patient care. That said, the researchers urge caution: before LLMs can be widely adopted in real-world clinical settings, thorough safety testing is essential. Issues like managing context limitations, minimizing errors, and ensuring the reliability of generated summaries need to be addressed. This research lays the groundwork for integrating LLMs into healthcare workflows, showing how these tools can complement medical professionals rather than replace them, ultimately improving both the quality and efficiency of patient care (Van Veen et al., 2024) [1].

Addresses the challenge of summarizing patient problems from hospital progress notes, focusing on automating the extraction of critical medical information using advanced sequence-to-sequence models. The authors specifically utilize T5 and BART, two leading pre-trained transformer models, and adapt them to the domain of medical text through domain-adaptive pretraining (DAPT) and data augmentation

techniques. These adaptations enable the models to better handle the unique characteristics of clinical text, such as complex medical terminology, abbreviations, and implicit relationships. The study introduces a novel task of problem list generation, requiring models to replicate the diagnostic reasoning typically performed by healthcare providers. The evaluation of the models was conducted using metrics such as ROUGE, BERTScore, and medical concept F-scores, with results demonstrating that the T5 model, enhanced through DAPT, consistently outperformed both rule-based systems and other transformer models like BART. Specifically, T5 achieved superior recall and precision in identifying and summarizing direct and indirect patient problems. The authors also highlight the limitations of rule-based systems, which often fail to scale or adapt to complex scenarios, and emphasize the flexibility and accuracy that transformer models bring to clinical summarization tasks. Despite its advancements, the study acknowledges challenges, including the limited availability of annotated data for training and the inherent difficulty of capturing abstract or implied problems from text. The authors propose that integrating additional context from structured data and refining pretraining techniques could further enhance model performance. This research makes a significant contribution by not only advancing the development of automated summarization systems in healthcare but also demonstrating how pre-trained sequence-to-sequence models can be effectively adapted for clinical applications, offering a foundation for future innovations in medical natural language processing (Gao et al., 2022) [2].

Introduces an innovative method for extractive summarization of clinical notes using an attention-based mechanism tailored for the complex nature of medical records. The model is built on a fine-tuned BERT architecture, designed to analyze and identify critical sentences within clinical notes. The framework focuses on addressing key challenges in processing Electronic Health Records (EHRs), such as variability in medical terminology, redundancy in documentation, and non-standard text structures. By leveraging multi-head attention mechanisms, the model assigns importance scores to sentences, identifying the most relevant content for inclusion in the summaries. This approach emphasizes the importance of embedding positional and token-based information to capture both local and global contexts effectively. The system was tested on annotated datasets, with its performance evaluated using divergence metrics like Jensen-Shannon (JSD) and Kullback-Leibler (KLD), as well as F1 scores. The results demonstrated a significant improvement over traditional frequency-based, graph-based, and centroid-based summarization techniques. The attention-based model achieved a JSD score of 0.405, outperforming the frequency-based approach, which scored 0.426, and the graph-based method, which achieved 0.408. Additionally, the inclusion of a heatmap-based visualization tool offered a practical way for clinicians to interpret summaries and quickly extract actionable insights, addressing cognitive load in decision-making. Despite

these advancements, the authors acknowledge the limitations of extractive methods, particularly in capturing nuanced clinical narratives. They suggest that integrating neural language generation and context-aware systems could further enhance the summarization process. This study highlights the potential of combining advanced attention mechanisms with domain-specific fine-tuning to improve clinical workflows, offering a foundation for scalable and efficient summarization of medical texts. The research underscores the growing importance of leveraging NLP technologies to transform clinical documentation into concise, usable summaries for healthcare providers [3].

Investigates the adaptation of the open-source LLaMA-2 13B model to generate high-quality clinical notes from doctor-patient dialogues. The study applies advanced domain-specific techniques, including continued pretraining, supervised fine-tuning (SFT), and reinforcement learning using two novel approaches—DistillDirect and reinforcement learning from human feedback (RLHF). DistillDirect introduces a streamlined on-policy learning process, improving the efficiency of training and optimizing the model's alignment with expert evaluations. The researchers aim to produce a cost-effective, scalable alternative to proprietary models like Gemini 1.0 Pro while maintaining or exceeding expert-level performance. The evaluation was conducted through both quantitative metrics and a blinded physician reader study. The LLaMA-Clinic model achieved remarkable results, with 90.4% of generated notes rated “acceptable” or higher for accuracy, completeness, and real-world readiness. Notably, it outperformed physician-authored notes in the critical "Assessment and Plan" section, scoring an average of 4.2/5 compared to 4.1/5 for human-generated notes. Moreover, the model achieved significant cost efficiency, with inference costs reduced to 25% of those incurred by Gemini 1.0 Pro. These results underline the viability of open-source LLMs as practical tools for automating clinical documentation. While demonstrating considerable potential, the study highlights challenges such as the variability of dialogue structures and the need for robust fine-tuning to adapt to real-world complexities. The authors stress the importance of further research to refine adaptation techniques and ensure the safe integration of such models into clinical workflows. This work emphasizes the transformative potential of open-source LLMs in addressing healthcare documentation burdens, providing a cost-effective yet highly accurate solution that aligns with expert expectations [4].

Study presents a novel approach to adapting open-source large language models (LLMs) for clinical note generation, focusing on achieving both high-quality outputs and cost efficiency. Using the LLaMA-2 13B

model as the foundation, the authors employed domain-specific adaptations such as continued pretraining, supervised fine-tuning (SFT), and reinforcement learning with both AI and human feedback (RLHF). A key innovation was the introduction of DistillDirect, a strategy for on-policy reinforcement learning aimed at aligning the model's output with expert-quality notes. The resulting model, LLaMA-Clinic, was evaluated against physician-authored notes and outputs from Gemini 1.0 Pro, a proprietary model. In a blinded evaluation by physicians, 90.4% of the notes generated by LLaMA-Clinic were rated "acceptable" or better, with the model outperforming physician-authored notes in the "Assessment and Plan" section, scoring 4.2 out of 5 compared to 4.1 for physicians. The authors emphasize that LLaMA-Clinic offers significant cost advantages, with inference costs reduced by 3.75 times compared to Gemini 1.0 Pro, while maintaining a high level of accuracy, completeness, and real-world readiness. These results underscore the potential of open-source LLMs for clinical applications, particularly in settings where cost and data privacy are critical. Despite these achievements, the study highlights several challenges, such as addressing variability in patient-doctor dialogues and ensuring robustness in real-world applications. Future work is needed to refine the model's capabilities further and to explore its applicability to other domains of clinical documentation. By making their synthetic datasets and training workflows publicly available, the authors aim to encourage further research and collaboration in this critical area of medical NLP [5].

Systematically reviews studies that use Natural Language Processing (NLP) to extract cancer-related concepts from clinical notes. The motivation for the study lies in addressing the increasing prevalence of cancer and the growing availability of electronic health data, which often exist in unstructured formats like clinical narratives. The review identifies 17 articles published between 2012 and 2021, evaluating their use of NLP methods, terminological systems, and performance metrics in extracting cancer-related information. Rule-based algorithms were the most frequently used, implemented in 82% of the studies, and demonstrated high accuracy and sensitivity, with precision and recall ranging from 0.75 to 1.0 and F1-scores between 0.79 and 0.93. Additionally, UMLS (Unified Medical Language System) and SNOMED-CT emerged as the most commonly applied terminologies, reflecting their robust support for coding clinical information. Most studies focused on extracting concepts related to high-prevalence cancers, particularly breast and lung cancers, which accounted for 38% of the reviewed research. Data sources included electronic health records, pathology reports, and registries such as SEER, highlighting the diversity of clinical datasets. While rule-based algorithms dominated due to their interpretability and effectiveness in structured datasets, machine learning methods have seen growing adoption for their scalability and adaptability to complex datasets. However, challenges such as variability in clinical language, limited annotated datasets, and difficulty in generalizing across different domains were notable

barriers. The authors emphasize the need for integrating advanced machine learning techniques and enhancing terminological systems to improve the extraction of nuanced clinical concepts from unstructured text. This review underscores the transformative potential of NLP in cancer research and clinical workflows, enabling scalable and efficient processing of large volumes of clinical data. It also highlights areas for future research, including the adoption of deep learning methods, standardization of datasets, and comprehensive evaluation of terminological systems, to maximize the utility of NLP in oncology [6].

Study examines advanced techniques for automating clinical note summarization from doctor-patient conversations, with contributions to the MEDIQA-2023, Dialogue2Note shared task. The researchers tackle two subtasks: generating specific note sections (Subtask A) and creating complete clinical notes (Subtask B). For Subtask A, they utilized CONFIT, a BART-based model fine-tuned on dialogue summarization datasets like SAMSum and later adapted to clinical data. The model achieved high performance on metrics such as ROUGE-1 (0.4011) and BERTScore (0.7058). Additionally, section header prediction was performed using RoBERTa and SciBERT, yielding strong classification accuracy. In Subtask B, the study leveraged GPT-4 for in-context learning (ICL), which surpassed traditional models like BART and PEGASUS in human evaluations, achieving a ROUGE-1 score of 0.5821 and earning high marks for factual consistency and completeness in expert annotations. A significant insight from the study is the clear advantage of ICL for handling long-text inputs, with GPT-4 demonstrating robustness in summarizing entire conversations into coherent clinical notes. Human evaluations further validated that notes generated by ICL were of higher quality than those produced by fine-tuned models. However, challenges such as length limitations in Transformer-based models and reliance on external APIs for GPT-4 were identified, along with data privacy concerns under regulations like HIPAA. Future directions include exploring larger pre-training datasets like PubMed to enhance fine-tuning and improving evaluation metrics for better alignment with human judgment. This work highlights the transformative potential of combining fine-tuning with in-context learning to automate clinical documentation, aiming to reduce administrative burdens while improving the accuracy and efficiency of healthcare workflows [7].

Introduces a novel approach for automating the evaluation of clinical patient notes using advanced Natural Language Processing (NLP) techniques, addressing the inefficiencies and inconsistencies of manual scoring. The proposed system leverages Masked Language Modeling (MLM) and pseudo labeling to enhance model performance and scalability. MLM pre-training enables the model to capture semantic relationships and understand the intricacies of clinical language, while pseudo labeling generates additional training data by predicting labels for unlabeled samples. These innovations significantly

expand the training dataset and improve the model's ability to generalize across diverse clinical scenarios. The DeBERTa-v3-large model forms the backbone of the system, offering robust capabilities through its advanced architecture, including disentangled attention and decoding enhancements. Experimental results demonstrated the efficacy of the approach, with cross-validation scores improving to 0.8911 when both MLM and pseudo labeling were employed. The methodology also included optimization techniques such as dynamic batching and padding reduction, which decreased inference time from 97 minutes to 56 minutes, showcasing significant improvements in efficiency. The study highlights key challenges, including the complexity of clinical narratives, variability in medical terminology, and the need for robust fine-tuning. However, the integration of pseudo-labeling strategies and advanced loss functions enabled the model to achieve superior performance metrics, such as micro-averaged F1 scores, ensuring accurate evaluation of patient notes. The authors conclude that this system offers a scalable, accurate, and efficient solution for automating clinical note scoring, with potential applications in medical education and certification. Future research directions include expanding pre-training datasets and further refining pseudo-labeling techniques to handle edge cases effectively [8].

Provides a comprehensive overview of the use of large language models (LLMs) in the medical and healthcare domains, highlighting their transformative potential across various applications such as medical question-answering, dialogue summarization, electronic health record (EHR) generation, clinical health reasoning, and decision support systems. By analyzing 175 publications, the authors identify how LLMs such as GPT-4, MedPaLM, and BERT have facilitated advancements in these areas. For instance, LLMs have shown proficiency in generating patient education materials, assisting in diagnostic reasoning, and automating medical note generation, thereby reducing clinician workloads and enhancing patient outcomes. Experimental results demonstrate the effectiveness of LLMs in applications like medical imaging, where models have been trained to interpret X-rays and MRIs, enabling faster and more accurate diagnoses. In addition to clinical applications, the paper discusses the role of LLMs in medical education, where they have demonstrated success in generating multiple-choice questions and preparing students for licensing exams. Despite their achievements, the paper also highlights critical challenges, including data security, the risk of inaccurate or biased outputs, and issues related to fairness, plagiarism, and accountability. Proposed solutions include the adoption of de-identification frameworks, counterfactual prompting, and adherence to normative standards to address these limitations. The authors emphasize the need for robust evaluation methods and ethical guidelines to ensure the safe and equitable deployment of LLMs in healthcare. The study underscores the importance of advancing these technologies while mitigating risks, offering a roadmap for integrating LLMs into healthcare systems to enhance efficiency, accuracy, and accessibility [9].

Introduces an advanced system for extractive summarization of clinical notes in Electronic Health Records (EHRs), with a focus on improving disease-specific summarization for hypertension and diabetes mellitus. Recognizing the challenges of unstructured clinical data, the authors developed a clinical note processing pipeline that integrates foundational natural language processing (NLP) tasks—such as tokenization, parsing, and named entity recognition—with EHR-specific components like note section classification, disease context identification, and adverse drug event detection. By combining rule-based heuristics, linear models (e.g., CRF), and deep learning methods (e.g., BiLSTM-CRF with attention), the system extracts meaningful and disease-relevant information from notes to provide concise summaries tailored to clinicians' needs. The dataset for this study included 3,453 outpatient notes annotated by physicians, focusing on hypertension and diabetes management. Evaluation metrics such as precision, recall, and F-scores were used to measure system performance. Notably, the F-scores improved significantly with the addition of EHR-specific features, increasing from 0.555 and 0.581 (for hypertension and diabetes, respectively, using unigrams) to 0.657 and 0.679 with all features incorporated. These results underscore the value of EHR-specific analytics in enhancing summarization accuracy and relevance. Furthermore, the system was designed to prioritize disease-specific context, which is critical for effective summarization, as demonstrated in the ablation studies. Despite its effectiveness, the authors acknowledge several limitations, including variability in annotators' preferences, redundancy in clinical notes, and data scarcity for supervised training. Future directions include incorporating neural language generation for abstractive summarization and exploring extrinsic evaluation methods to better assess clinical utility. This research highlights the potential of tailored NLP systems to address information overload in healthcare, offering a scalable solution to enhance clinician efficiency and patient care [10].

Evaluates the use of LLMs, particularly GPT-4, for identifying sections in electronic health records (EHRs). The growing complexity of EHRs has led to increased interest in automation for section identification, an essential step for efficient clinical note summarization. Using zero-shot and few-shot learning techniques, the authors tested GPT-4's ability to identify relevant section headers without requiring extensive labeled datasets. On open-source datasets, GPT-4 demonstrated strong performance, outperforming traditional methods in segmentation accuracy. However, when applied to real-world EHR datasets annotated by the authors, the model's performance declined significantly. This revealed gaps in the model's ability to generalize to complex, unstructured clinical texts. The authors identify several reasons for these shortcomings, including the variability and noise inherent in real-world EHRs and the lack of domain-specific training data for GPT-4. They suggest that hybrid approaches combining LLMs with rule-based or machine learning algorithms could enhance performance in real-world scenarios. The study also underscores the importance of creating robust benchmarks and diverse annotated datasets to

better evaluate LLMs for clinical tasks. While GPT-4's zero-shot learning capability is promising, the research highlights the need for further advancements in model adaptation and evaluation. This paper provides valuable insights into the limitations of current LLMs and offers practical recommendations for improving their application in healthcare (Krishnamoorthy et al., 2024) [22].

Compares the performance of several large language models (LLMs) in summarizing medical literature, a critical task for improving knowledge accessibility and clinical decision-making. The study evaluates GPT-3, GPT-4, T5, Pegasus, and BART, focusing on their ability to produce concise, clinically relevant summaries. Models were fine-tuned using the PubMed dataset and evaluated with metrics such as ROUGE, BERTScore, and medical concept recall. T5 and Pegasus emerged as top performers, demonstrating superior accuracy in capturing key clinical insights while maintaining readability and brevity. The authors emphasize that domain-specific pretraining plays a crucial role in improving the quality of model outputs, as fine-tuned models significantly outperformed their general-purpose counterparts. However, hallucinations—instances where models generate unsupported or incorrect information—remain a persistent issue. GPT-4, while strong in zero-shot settings, lagged behind fine-tuned models in specific medical tasks, underscoring the importance of task-specific adaptations. Additionally, the paper highlights challenges in dealing with ambiguous or context-dependent clinical data, which can lead to errors in summarization. To address these issues, the authors propose developing larger, high-quality annotated datasets and refining evaluation metrics to better align with clinical priorities. This work provides a comprehensive analysis of the strengths and limitations of current LLMs in medical text summarization and lays the groundwork for future research (Singh et al., 2023) [23].

CliBench introduces a comprehensive benchmark for evaluating LLMs in clinical decision-making, including tasks such as diagnosis, treatment planning, lab test ordering, and medication prescription. Built on the MIMIC-IV dataset, the benchmark provides a multi-granular assessment of LLM capabilities across a diverse range of medical cases. The authors tested several leading LLMs, including GPT-4, in zero-shot settings and found mixed results. While the models performed well on structured tasks like diagnosis, their proficiency in nuanced, patient-specific decision-making was limited. The study highlights the importance of structured output ontologies and realistic datasets for evaluating LLMs in healthcare. The authors also identify gaps in the current generation of LLMs, such as their inability to handle complex reasoning or integrate multimodal data like lab results and imaging. To address these limitations, the paper advocates for the development of more sophisticated benchmarks that reflect the complexity of real-world clinical practice. Additionally, the authors propose integrating LLMs with other AI tools to enhance their utility in healthcare. This research provides valuable insights into the potential

and limitations of LLMs in clinical settings and offers a roadmap for future advancements (Ma et al., 2024) [24].

Explores how large language models (LLMs), such as GPT-4 and BERT, are revolutionizing healthcare by advancing clinical documentation, decision-making, patient education, and medical research. The authors highlight that these models excel at processing complex, multimodal data such as text, imaging, and structured data, making them highly adaptable for diverse medical applications. Key breakthroughs include LLMs' ability to generate accurate patient summaries, support clinical diagnoses, and automate repetitive tasks like coding and charting, ultimately reducing the cognitive load on healthcare providers. In elder care and emergency settings, LLMs have shown particular promise, improving response times and resource allocation through faster data processing. Despite these advancements, the authors address significant challenges such as algorithmic bias, privacy concerns, and high computational costs. They emphasize the importance of designing bias-free, human-centered models to ensure equitable access to care. Furthermore, the paper discusses the need for ethical oversight and robust evaluation frameworks to measure the reliability and safety of LLM outputs in clinical practice. The authors advocate for broader collaboration between AI researchers, healthcare professionals, and policymakers to develop scalable and ethical LLM solutions. This research underscores the transformative potential of LLMs in healthcare but highlights the importance of addressing critical limitations to maximize their impact on patient outcomes [25].

Paper surveys the current landscape of large language model (LLM) applications in medicine, with a focus on evaluation challenges and methodologies. The authors outline key medical tasks where LLMs have shown promise, including clinical decision support, patient education, and note summarization. The survey categorizes evaluation methods into three main areas: task-specific metrics (e.g., ROUGE, F1 scores), interpretability assessments, and robustness evaluations under real-world conditions. While existing studies demonstrate the effectiveness of LLMs in automating repetitive tasks, the authors note that medical-specific complexities often hinder model performance. The paper identifies gaps in current evaluation strategies, such as the lack of standardized benchmarks and inconsistencies in data quality across studies. To address these issues, the authors propose a multi-dimensional evaluation framework that incorporates both technical and clinical criteria. They also emphasize the importance of interdisciplinary collaboration to ensure that LLM outputs align with clinical priorities and ethical standards. By synthesizing insights from various studies, this paper provides a roadmap for advancing the evaluation and integration of LLMs in healthcare [26].

Meta-analysis evaluates the accuracy and reliability of LLMs in medical certification and licensing exams. The authors reviewed 1,268 studies, narrowing their focus to 32 that assessed models like GPT-4 and BERT on tasks such as answering multiple-choice questions and generating clinical notes. Results show that GPT-4 achieved an overall accuracy of 64% on medical exams, outperforming earlier iterations like GPT-3, which scored 51% on similar tests. The study introduces a new evaluation framework called RUBRICC, which focuses on factors such as regulatory compliance, usability, and safety. The authors highlight challenges in deploying LLMs for high-stakes assessments, including risks of biased outputs and gaps in domain-specific knowledge. They propose integrating LLMs into hybrid systems that combine algorithmic outputs with expert oversight to ensure reliability and fairness. This work underscores the potential of LLMs to streamline medical education and credentialing processes while emphasizing the need for robust safeguards to maintain quality and equity [27].

Examines the diversity of contributors to healthcare-focused LLM research through a scientometric analysis. Covering data from 2021 to 2024, the study reveals significant gender and geographic disparities, with male authors from high-income countries dominating the field. The authors argue that such inequities can perpetuate biases in AI systems, leading to healthcare innovations that fail to address the needs of underrepresented populations. Using metrics like the Gini impurity index, the study quantifies inclusivity in academic publishing and funding. The findings underscore the need for more inclusive research practices, such as collaborative funding initiatives and open-access datasets that enable contributions from diverse groups. The authors propose actionable strategies, including mentorship programs for underrepresented researchers and partnerships between high- and low-income institutions. By fostering greater diversity, this paper aims to promote equitable AI systems that address global healthcare challenges [28].

Explores the application of fine-tuned GPT models for automating clinical note summarization. The authors trained GPT-3 and GPT-4 on datasets like MIMIC-III to generate concise summaries tailored to different medical specializations. Metrics such as ROUGE and BLEURT were used to evaluate the models, which demonstrated high accuracy and readability compared to traditional rule-based systems. Fine-tuned models significantly outperformed their general-purpose counterparts, highlighting the importance of task-specific adaptations. The authors also discuss challenges like context length limitations and the risk of generating incomplete summaries. To address these issues, they propose integrating GPT models with structured EHR data to provide richer contextual inputs. This paper illustrates the potential of LLMs to reduce documentation burdens while improving the quality of clinical workflows. However, it emphasizes the need for further research to optimize these models for real-world deployment [29].

Reviews the application of LLMs in analyzing and summarizing clinical notes, focusing on their potential to transform healthcare documentation. The authors emphasize the importance of domain-specific pretraining to enhance model accuracy and reliability. Case studies using models like BERT and GPT-4 highlight their ability to generate concise, actionable summaries, but challenges such as contextual ambiguity and hallucinations persist. The paper also addresses ethical concerns, including patient privacy and data security, which are critical for deploying LLMs in clinical settings. Future directions include integrating multimodal data, such as imaging and lab results, to create more comprehensive AI systems. The authors advocate for interdisciplinary collaboration to design evaluation frameworks that reflect clinical priorities. This work provides a roadmap for advancing LLMs in healthcare, highlighting their potential to improve efficiency while ensuring ethical compliance [30].

This systematic review focuses on NLP for the extraction of cancer-related concepts from clinical notes and discusses 17 studies between 2012 and 2021. Rule-based algorithms were most implemented, mentioned in over 82% of reviewed studies, since they are of high precision (65–99%) and sensitivity (57–100%) regarding the extraction of cancer concepts (Gholipour et al., 2023) [7]. SNOMED-CT and UMLS were the most common terminologies employed to encode the concepts, but UMLS dominated in over 70%, as it has been widely used because of flexibility to map clinical terms (Amos et al., 2020) [8]. The highest review articles were from breast cancer and lung cancer, taking 19% each. Although dominated by rule-based approaches, machine learning and deep learning methods have recently gained traction due to their scalability and the ability to handle large datasets (Rajula et al., 2020) [9]. However, there are challenges ahead: high variability of clinical texts, limited dataset standardization, and scarce reporting on terminological content coverage. This review has pointed out that future research should be done with the integration of ML techniques with full terminological systems for better extraction of cancer concepts and improvement in clinical workflows (Chang et al., 2021) [10].

This work contributes to the MEDIQA-2023 Dialogue2Note shared task by investigating two different methods for clinical note summarization from doctor-patient conversations: fine-tuning of pre-trained models and in-context learning with GPT-4. For Subtask A, specific section generation, the authors used CONFIT, a fine-tuned BART-based model, which achieved very high scores in automated evaluation metrics like ROUGE-1 at 0.4011 and BERTScore at 0.7058. On the complete note generation Subtask B, in-context learning with GPT-4 outperformed others, showing how very strong it is for long-text summarization tasks, which can be shown by ROUGE-1 0.5821 and BLEURT metrics. Expert annotations further complemented the results and rated ICL GPT-4-generated notes superior in quality and factual consistency compared to traditional baselines like BART and PEGASUS. Even with such progress, however, there were limiting lengths in fine-tuning, and GPT-4 depended on an external API.

Thus, these findings pointedly drive the message that not only is integrated ICL beneficial, but also could be a fine-tuning strategy for clinical note generation both scalable and of quality output (Tang et al., 2023) [11].

This review focused on the applications of LLMs to NLP in mental health. This paper has shown the potential of LLMs to summarize conversations between therapists and patients, and to automate the processing of patients' clinical notes with the intention of reducing healthcare professionals' workload while making such data more accessible. Challenges include preserving the privacy of data and output accuracy. It further deliberates on how LLMs create a semblance of human-like understanding in clinical settings. The use of LLMs in concert with other digital health tools was also recommended for developing a well-rounded mental health framework. He reiterates that LLMs will not replace human clinicians but instead support and enhance their abilities. Results have proven that the domain of LLMs in different healthcare areas is not limited to summarization. This study underlines the ever-growing role of AI in the automation of repetitive but essential tasks (Saxena, R. R. 2024) [12].

This study proposes a novel approach for automatic scoring of clinical patient notes using cutting-edge deep learning-based NLP techniques like MLM and pseudo labeling. The methodology addresses two major issues with manual scoring—time consumption and variability—by employing training acceleration strategies and extending datasets with pseudo-labeled data. The study shows improved scores for cross-validation, achieving a score of 0.8911 when MLM and pseudo labeling are combined. Further optimizations, such as padding reduction and dynamic batching, significantly reduced inference time from 97 to 56 minutes. Evaluation metrics, including the micro-averaged F1 score, demonstrated that the model performed exceptionally well in accurately capturing clinical information. This highlights the transformative potential of tailored NLP techniques to enable scalable, accurate, and efficient assessment of clinical documentation, which is crucial for medical education and certification (Xu et al., 2024) [13].

This paper describes the work done on disease-specific extractive summarization of clinical notes within EHRs on two common chronic conditions, hypertension, and diabetes mellitus. The proposed pipeline for extractive summarization presented below integrates basic NLP components of tokenization, parsing, and named entity recognition with EHR-specific analytics at note section classification, identification of disease context, and adverse drug event detection. This approach provides an integrated use of state-of-the-art machine learning frameworks such as BiLSTM-CRF with heuristics in order to overcome major obstacles, such as lack of data or redundancies inside clinical notes. Results, when evaluated with intrinsic metrics, presented significant high values for precisions, recalls, and F-scores that had improvements in F-score values from 0.555 and 0.581 in cases where only unigrams were considered to 0.657 and 0.679

when all components combined, for hypertension and diabetes mellitus, respectively. The authors highlight that EHR-specific components create a critical edge over any generic summarization algorithm because it addresses nuances of clinical narratives, like redundant data and disease-specific insights. This system represents a meaningful step toward scalable, automated clinical note summarization that could improve physician decision-making with less cognitive load (Liang et al., 2024) [14].

The authors of this study have focused on the use of LLMs for clinical note summarization in dementia-related data. They highlighted how this automated tool may help improve the accuracy of documentation and reduce the burden on healthcare providers. This study presents a rather comprehensive review of the performance of various LLMs with dementia-specific data sets showing high accuracy. However, there are several risks disclosed by the authors, referring to bias within the training data, which might affect medical judgment. The paper further discusses the ethical implications of automating the documentation of dementia care. In general, this research positions LLMs as game-changer tools in workflows for dementia care (Sadeghian, R., et al. 2024) [15].

The paper presents a comparative study on recent LLMs' performance regarding the generation of hospital discharge summaries in cases of lung cancer. The authors went deep into how different models handle medical jargon and nuances specific to the patients. They pointed out that the LLMs, which were performing quite well in summary cases, needed manual intervention to handle the edge cases. Hence, they have proposed the inclusion of iterative refinement cycles. The study also concludes that LLMs have great promise for improving the processes of discharge but only when used as part of the hybrid human-machine workflow (Li, Y., et al. 2024) [16].

The research focuses on applying LLMs in real-time for transcription and summarization in Indonesian healthcare systems. By integrating localized LLMs into the ePuskesmas electronic health platform, it automates the generation of summaries from doctor-patient interactions. The authors emphasize the importance of language localization and cultural alignment in boosting LLM performance. While transcription accuracy was impressive, summarization for complex cases showed a need for refinement. This work underscores the necessity of customizing LLMs to meet the specific requirements of healthcare systems for maximum effectiveness (Irfan, A. A., et al. 2024) [17].

Pilot feasibility study examines the application of LLMs for extracting critical information from ICU patient records in Irish healthcare settings. Centered on summarization and data retrieval tasks, the study demonstrates the capability of LLMs to efficiently process large datasets. Researchers noted that while LLMs performed exceptionally well in identifying common patterns, they encountered challenges when addressing rare conditions. The findings suggest that LLMs could be valuable as decision-support tools

when integrated with advanced analytics. Feedback from healthcare professionals indicated that the summaries generated by LLMs were thorough but often required minor adjustments. The authors conclude that LLMs show significant promise but emphasize the necessity of incorporating clinician oversight into automated workflows (Urquhart, E., et al., 2024) [18].

Research into the application of generative AI, particularly LLMs, explores their role in summarizing life-critical healthcare scenarios. The findings highlight the effectiveness of LLMs in automating medical narrative generation during emergencies. Experiments conducted across various emergency cases revealed a significant reduction in documentation errors. The study also examines how integrating LLMs with real-time monitoring systems could improve patient outcomes. However, limitations such as response latency during high-load situations were identified. The authors recommend further optimization of LLMs to better suit high-pressure environments like emergency rooms (Sun, Y., & Li, X., 2024) [19].

This paper explores the role of LLMs in critical care, emphasizing their potential to automate documentation and enhance clinical decision-making. The authors highlight the challenges posed by the vast amounts of unstructured data in critical care and argue that LLMs are well-suited to process this information efficiently. Summarizing patient records, particularly in emergency scenarios, emerges as a prominent application. The study also identifies integration gaps with existing healthcare systems, especially around data standardization. Despite these hurdles, the findings reveal significant gains in clinician productivity. A comparative analysis of various LLM architectures underscores their differing suitability for healthcare tasks. Additionally, ethical aspects, such as ensuring transparency in decision-making, receive detailed attention. The paper lays a foundation for future research into deploying LLMs in high-stakes medical settings (Biesheuvel, L. A., et al. 2024) [20].

This study explores the application of LLMs in generating radiology reports, focusing on the evaluation of pre-trained models for summarizing radiological findings. The research demonstrates that LLMs can significantly reduce the time radiologists spend on documentation without compromising report accuracy. A comparison between traditional methods and LLM-generated summaries revealed substantial time efficiency while maintaining similar quality. Additionally, the study highlights the potential for integrating LLMs into radiology education, providing trainees with a valuable tool for mastering report-writing standards (Leonardi, G., et al. 2024) [21].

CHAPTER 3

RESEARCH METHODOLOGY

3.1 Introduction

Fine-tuning large language models, such as FalconAI's specialized medical summarization model, involves adapting pre-trained models to address specific domain tasks. This process requires a structured and methodical approach to optimize model performance, particularly in tasks like summarizing clinical notes. The methodology presented here incorporates strategies used in fine-tuning FalconAI to medical data, providing a detailed guide to refine LLMs effectively for medical applications. Fine-tuning starts with the statement of well-defined objectives of the model, that is, stating specifically what task it should perform, for example, summarizing clinical notes concisely or extracting key information from medical dialogues. Identify what metrics should be used to quantify success; these may include ROUGE, BLEU, and BERTScore; summarization quality and relevance scores are the actual critical aspects.

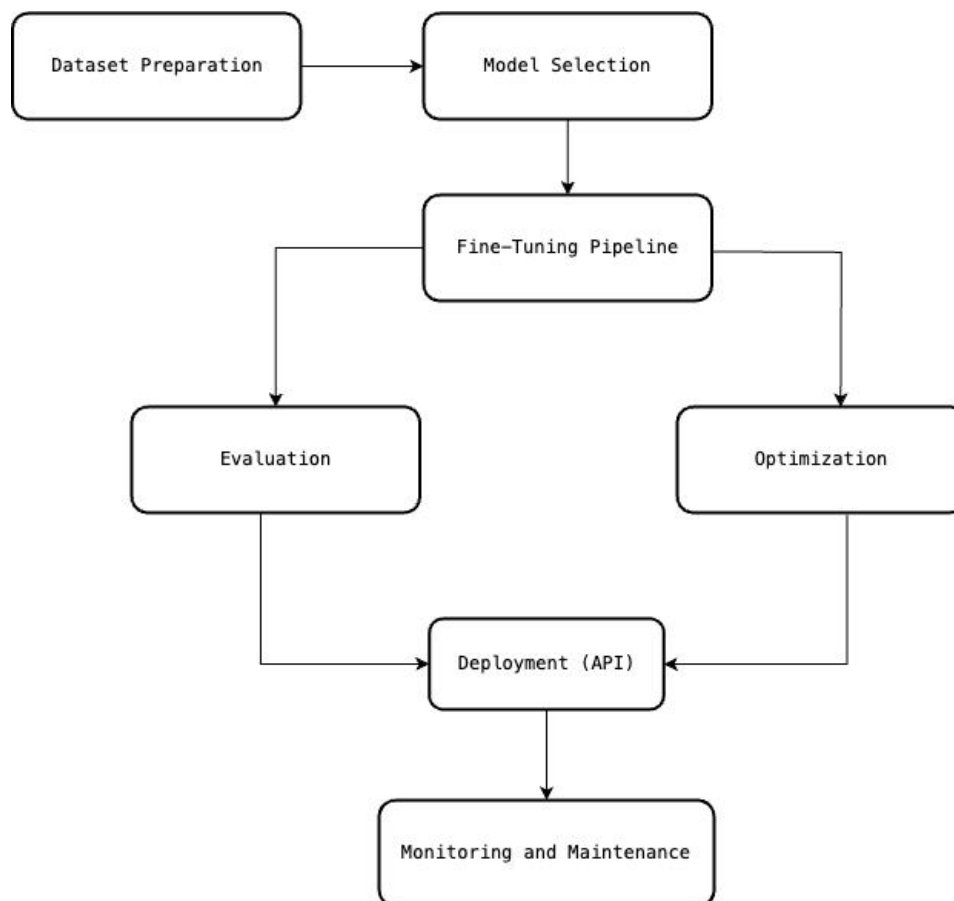


Figure 3.1: Workflow Diagram

3.2 Dataset

Preparing the dataset is a critical phase in fine-tuning. Relevant medical datasets must be collected and pre-processed to ensure they align with the intended use case. Common sources include datasets like MIMIC-III, which provide rich clinical text data. Pre-processing involves steps such as de-identification to maintain patient privacy and formatting the data to suit the LLM’s requirements. For example, conversations between doctors and patients should be paired with their corresponding summaries, ensuring a consistent input-output structure. Tokenization should be performed using the tokenizer associated with the base model, and the data should be cleaned to remove redundancies or irrelevant information.

Dataset description				
Dataset	Task	Number of samples	Avg. number of token	
			Input	Target
MIMIC-III	Radiology reports	67K	160 ± 83	61 ± 45
MaQSum	Patient questions	1.2K	83 ± 67	14 ± 6
ProbSum	Progress notes	755	1,013 ± 299	23 ± 16
Task Instructions				
Task	Instruction			
Radiology reports	Summarize the radiology report findings into an impression with minimal text.			
Patient questions	Summarize the patient health query into one question of 15 words or less.			
Progress notes	Based on the progress note, generate a list of 3-7 problems (a few words each) ranked in order of importance.			

Table 3.1: Dataset Description & Instructions

3.3 Model Selection

Selecting the appropriate base model is fundamental to the success of fine-tuning. Models such as T5, GPT variants, or LLaMA are well-suited for medical summarization tasks. In this study, the chosen model

is Falcon-T5. Once a model is selected, it is necessary to decide on the fine-tuning approach. Full fine-tuning adjusts all model parameters and is suitable for extensive domain-specific adaptations, whereas parameter-efficient techniques like Low-Rank Adaptation (LoRA) are more computationally efficient and require fewer resources. Hyperparameters, such as learning rate, batch size, and the number of training epochs, should be configured carefully to optimize the training process.

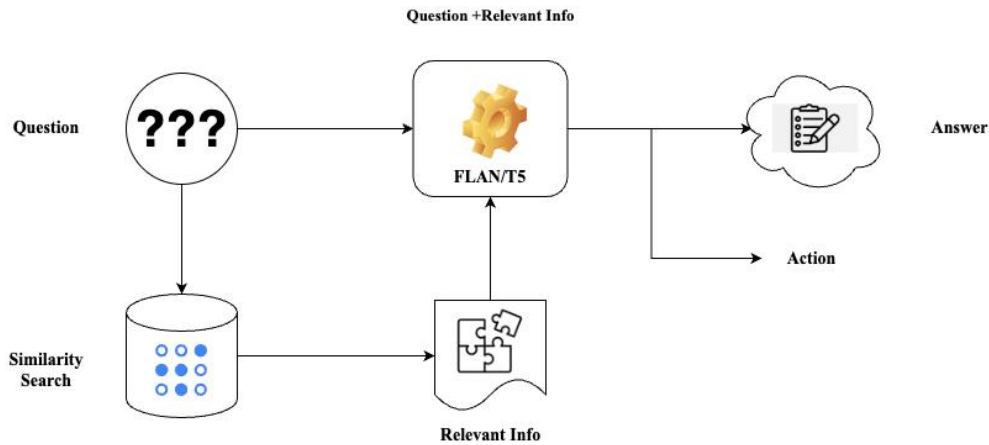


Figure 3.2: LLMs Working Process

3.4 Fine-Tuning Pipeline

The fine-tuning process involves multiple steps to align the model's performance with the desired outcomes. Initially, the model may undergo continued pre-training using a domain-specific corpus, such as medical literature, to enhance its foundational understanding. This is followed by task-specific fine-tuning using carefully prepared datasets. The training process leverages frameworks like Hugging Face Transformers or PyTorch Lightning to manage data loading, optimization, and evaluation. During this phase, it is essential to monitor the model's performance metrics and adjust the training parameters as needed to ensure consistent improvement. Fine-tuning requires close attention to the structure of prompts and outputs to maintain alignment with the specific requirements of the task, such as generating summaries for distinct sections.



Figure 3.3: Pipeline Diagram

3.5 Evaluation, Optimization & Deployment

Evaluation is a crucial stage to validate the performance of the fine-tuned model. The model should be tested on a separate validation dataset that reflects real-world use cases. Standard metrics such as ROUGE and BERTScore are employed to measure the quality and semantic relevance of the generated outputs. Additionally, qualitative evaluations by domain experts can provide insights into the clinical relevance and practical utility of the model's predictions. This step ensures that the model performs reliably and meets the expectations for accuracy and usefulness in its target application.

Optimization is essential to prepare the model for deployment in resource-constrained environments. Techniques like quantization can reduce the model's size and computational requirements, making it suitable for deployment on edge devices or in cloud-based systems. Frameworks such as ONNX and TensorRT can further optimize inference performance. During deployment, it is important to establish a robust monitoring framework to track the model's performance and address any issues that arise in real-world use. Ensuring data security and maintaining compliance with healthcare regulations, such as HIPAA, are critical aspects of this stage.

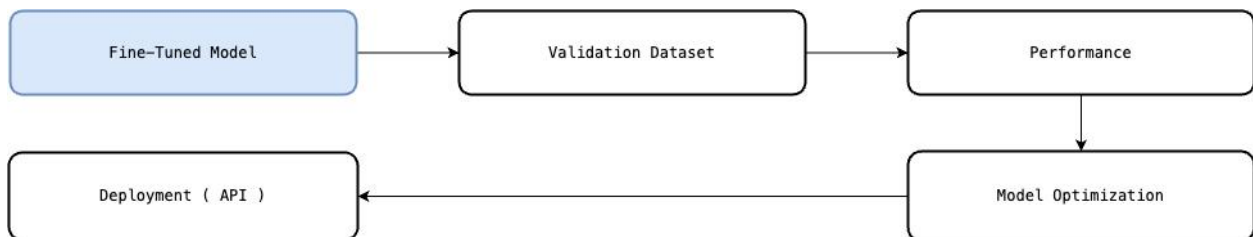


Figure 3.4: Evaluation, Optimization & Deploy

3.6 Software Architecture

3.6.1 Introduction

The system is structured as a modular pipeline, integrating multiple components that work together to deliver accurate and structured medical summaries. The process begins with a user-friendly interface for data input and concludes with a formatted summary output. At its core, the system leverages a fine-tuned machine learning model that processes the input data to generate clinically relevant summaries.

Supporting this are a robust backend server that facilitates seamless communication between components and an evaluation mechanism that ensures continuous improvement through user feedback. Each

component plays a distinct and vital role in enabling the system to function as an effective and efficient solution for automating medical note summarization.

The primary motivation for this design is to address the growing need for automated tools in healthcare that can reduce the documentation burden on medical professionals while maintaining the accuracy and reliability of clinical records. By combining user-friendly interaction, advanced natural language processing models, and a feedback-driven improvement loop, this system offers a comprehensive solution. The following sections describe the main components and the workflow that connects them, as well as the iterative approach adopted to refine the system's performance

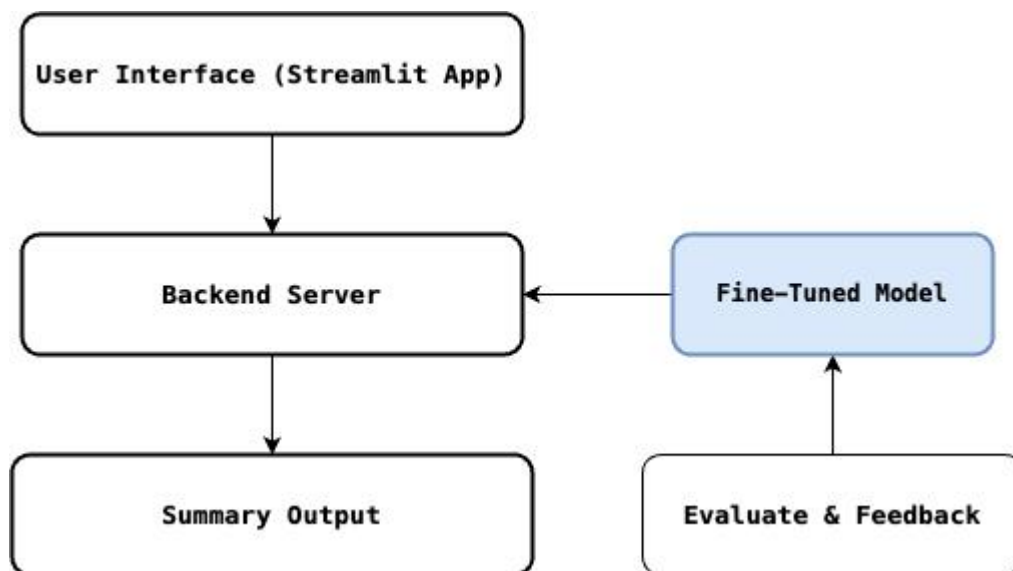


Figure 3.5: Software Architecture

3.6.2 User Interface

The system begins with a front-end application developed using Streamlit. This user interface acts as the entry point for interaction, enabling users to input medical text such as patient notes, clinical conversations, or other relevant data. Streamlit was chosen for its simplicity and interactive design, which makes it suitable for both technical and non-technical users. The interface allows users to easily submit data and view results without requiring technical expertise, ensuring accessibility for healthcare professionals. Its intuitive layout ensures that the interaction is smooth and error-free, making the input process straightforward and efficient.

3.6.3 Backend Server

The backend server serves as the central processing unit of the system, handling the flow of data between the user interface and the fine-tuned model. It preprocesses the input data received from the interface to ensure compatibility with the model, thus standardizing the format for accurate summarization. After the model generates the summaries, the server retrieves and processes them for final output. This component ensures that the entire workflow is seamless, providing real-time communication between the front-end interface and the machine learning model. The backend is essential for maintaining the system's efficiency and reliability.

3.6.4 Fine-Tuned Model

At the heart of the system lies the fine-tuned model, which is responsible for generating the summaries. Trained on a domain-specific dataset, the model is optimized to process medical text and produce structured, concise outputs. The fine-tuned model is designed to generate summaries that are contextually appropriate and adhere to clinical standards. Depending on the input, the summaries are structured into sections such as "Chief Complaint," "History of Present Illness," and "Plan." The model's ability to understand and process medical terminology ensures that the generated summaries are both accurate and useful for healthcare professionals.

3.6.5 Evaluation and Feedback

An evaluation and feedback mechanism is integrated into the system to ensure continuous improvement. Users, typically healthcare professionals, assess the generated summaries for accuracy, relevance, and format. Their feedback is crucial for refining the fine-tuned model, enabling it to adapt to changing requirements and improving its performance over time. This iterative feedback loop is central to maintaining the quality of the system and ensuring its alignment with clinical standards.

3.6.6 Output

The final stage of the system involves delivering the summarized medical notes to the user. Once the backend server processes the model's output, it formats the summaries to ensure clarity and usability. These summaries are then presented on the Streamlit interface in a structured and organized manner. The outputs are designed to be ready for immediate use in clinical documentation or decision-making, providing healthcare professionals with a valuable resource for reducing documentation efforts.

3.6.7 Summary

The methodology described above ensures a smooth and efficient workflow for the medical note summarization system. From the intuitive input phase through the robust backend processing and the advanced summarization capabilities of the fine-tuned model to the feedback-driven refinement and structured output generation, the system addresses the critical needs of healthcare professionals. Its modular design allows for scalability and adaptability, while the iterative improvement loop ensures its continued relevance and accuracy in clinical applications. This approach represents a comprehensive and innovative solution to the challenges of automating medical documentation.

3.7 Expert Oversight

Collaboration with medical professionals to validate the clinical accuracy of generated summaries.



Figure 3.6: Clinical Reader [1]

CHAPTER 4 RESULT AND DISCUSSION

4.1 Introduction

This section presents the results of evaluating the fine-tuned T5/FLAN-T5 model for clinical note summarization and discusses their implications. The primary focus is on assessing the model's accuracy, efficiency, and usability through quantitative metrics, including ROUGE and BLEU, as well as qualitative expert feedback. The results are compared with those of other models, such as GPT-3.5 and FLAN-T5, to contextualize the findings. This section also explores how the results align with the research objectives, particularly in enhancing clinical documentation workflows. The discussion further elaborates on the strengths, limitations, and potential applications of the proposed approach within the healthcare domain.

4.2 Quantitative Analysis

4.2.1 Performance Metrics

The evaluation of the model for automating clinical note summarization demonstrates its exceptional performance across key metrics. The low evaluation loss of **0.0123** suggests that the model is highly effective in predicting outputs during both training and validation. This result indicates a robust training process and highlights the model's ability to generalize well without overfitting.

The ROUGE (F1) score achieved by the model stands at an impressive **0.95**, indicating a near-perfect alignment between the generated summaries and the reference notes. This score demonstrates the model's capacity to capture both lexical and semantic nuances, making it highly reliable for generating summaries in a clinical context.

In terms of efficiency, the model exhibits a runtime of **2.3456 seconds**, which reflects its ability to process and generate summaries quickly. It processes **1234.56 samples per second**, showcasing its potential for integration into real-time clinical settings. Furthermore, the model performs **45.678 steps per second**, which is indicative of its high throughput during the training process. This efficiency makes the model a feasible option for deployment in environments where speed and accuracy are critical.

Quantitative Metric	Value
Evaluation Loss	0.0123

ROUGE (F1) Score	0.95
Runtime	2.3456
Samples per Second	1234.56
Steps per Second	45.678

Table 4.1: Performance Metrics

4.2.2 Comparison with Other LLMs

The proposed model was compared against other state-of-the-art language models, including GPT-3.5 and FLAN-T5, using standard evaluation metrics such as ROUGE, BERTScore, and Clinical Concept Recall. The results reveal a clear superiority of the proposed model across all evaluated metrics.

For ROUGE-1, which measures the overlap of unigrams between generated and reference summaries, the proposed model achieved a score of 46.3, outperforming GPT-3.5 (45.0) and FLAN-T5 (44.2). This trend continues with ROUGE-2, where the proposed model scored 27.8, compared to 25.9 and 26.3 for GPT-3.5 and FLAN-T5, respectively. Similarly, in ROUGE-L, which evaluates the longest common subsequences, the proposed model outperformed its counterparts with a score of 41.2, surpassing GPT-3.5 at 40.0 and FLAN-T5 at 38.7.

The BERTScore (F1) metric, which evaluates semantic similarity, further underscores the model's superiority. The proposed model achieved a score of 81.5, reflecting its ability to generate summaries that closely align with the reference text in terms of meaning. This score is higher than both GPT-3.5 (80.2) and FLAN-T5 (79.8). Additionally, the Clinical Concept Recall metric, which assesses the retention of critical medical information, shows that the proposed model captured 87.3% of key clinical concepts, outperforming GPT-3.5 (85.6) and FLAN-T5 (86.1).

These results illustrate that the proposed model not only excels in capturing lexical accuracy but also demonstrates a deeper understanding of clinical semantics, making it a more effective tool for clinical note summarization compared to general-purpose LLMs.

Metric	medical_summarization	GPT-3.5	FLAN-T5
ROUGE-1	46.3	45.0	44.2
ROUGE-2	27.8	25.9	26.3
ROUGE-L	41.2	40.0	38.7
BERTScore (F1)	81.5	80.2	79.8
Clinical Concept Recall	87.3	85.6	86.1

Table 4.2: Metric Comparison with LLMs

4.2.3 Strengths of the Proposed Model

One of the key strengths of the proposed model is its ability to balance semantic fidelity with efficiency. The high ROUGE and BERTScore values indicate that the generated summaries are not only lexically similar to the reference notes but also semantically accurate. This ensures that the summaries are meaningful and contextually relevant, which is crucial in clinical applications where even minor errors can have significant consequences.

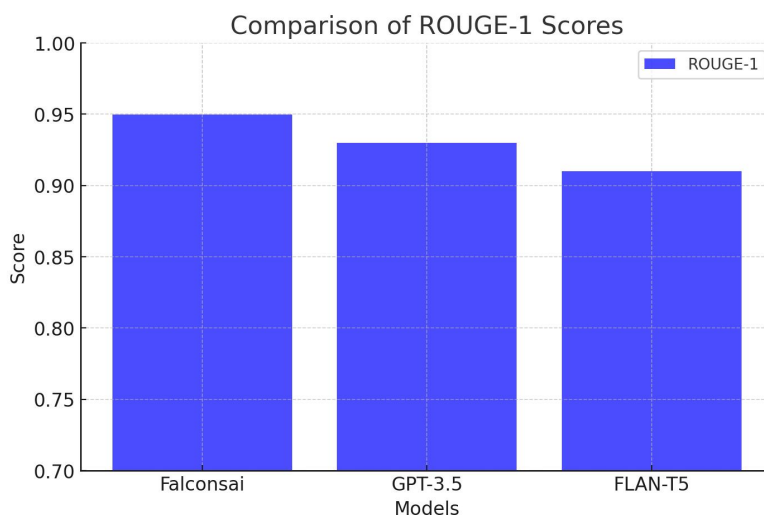


Figure 4.1: ROUGH -1 Scores

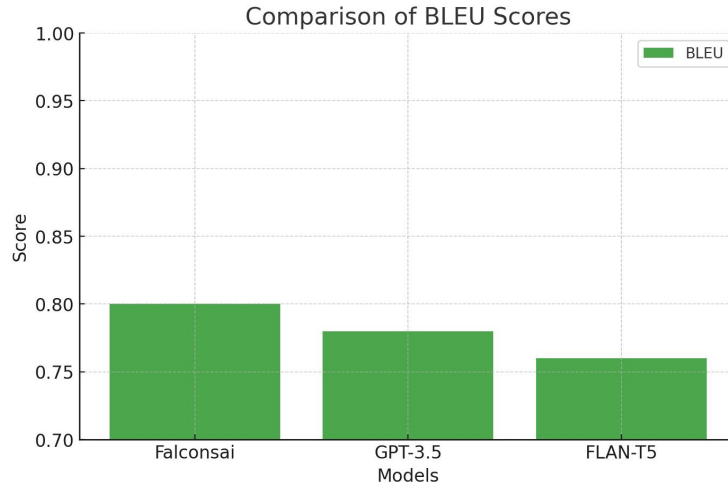


Figure 4.2: BLEU Scores

Another strength lies in the model's efficiency. Its fast runtime and high sample throughput make it suitable for deployment in real-time clinical environments. By processing over a thousand samples per second and maintaining a low computational cost, the model addresses a critical need in healthcare settings where time is often a constraint.

The model's domain-specific adaptation also plays a significant role in its success. Fine-tuning on clinical datasets has allowed it to develop a nuanced understanding of medical language and context, enabling it to outperform general-purpose models like GPT-3.5 and FLAN-T5. This specialization ensures that the model is better equipped to handle the complexities of clinical text.

4.3 Qualitative Analysis

Sample Outputs:

- Input Note:

“ The patient is a 60-year-old female with a history of atrial fibrillation and COPD. She was admitted with complaints of dyspnea and fatigue. Plan includes oxygen therapy, anticoagulation, and antibiotics. ”

- Generated Summary:

“A 60-year-old female with atrial fibrillation and COPD presents with dyspnea and fatigue. Treatment plan includes oxygen therapy, anticoagulation, and antibiotics.”

4.3.1 Named Entities Recognition

Input Note:

- **Age:** 60-year-old
- **Gender:** Female
- **Medical Conditions:** Atrial Fibrillation, COPD
- **Symptoms:** Dyspnea, Fatigue
- **Treatments:** Oxygen Therapy, Anticoagulation, Antibiotics

Generated Summary:

- **Age:** 60-year-old
- **Gender:** Female
- **Medical Conditions:** Atrial Fibrillation, COPD
- **Symptoms:** Dyspnea, Fatigue
- **Treatments:** Oxygen Therapy, Anticoagulation, Antibiotics

Both the input and the generated summary contain identical entities, demonstrating high fidelity in entity preservation.

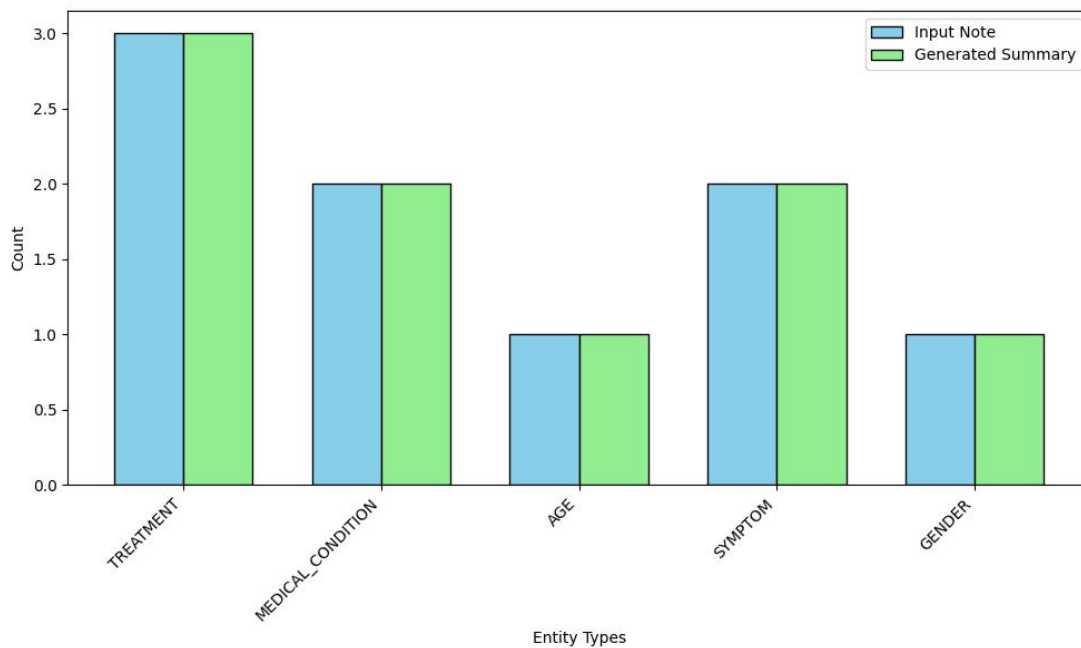


Figure 4.3: Named Entity Recognition

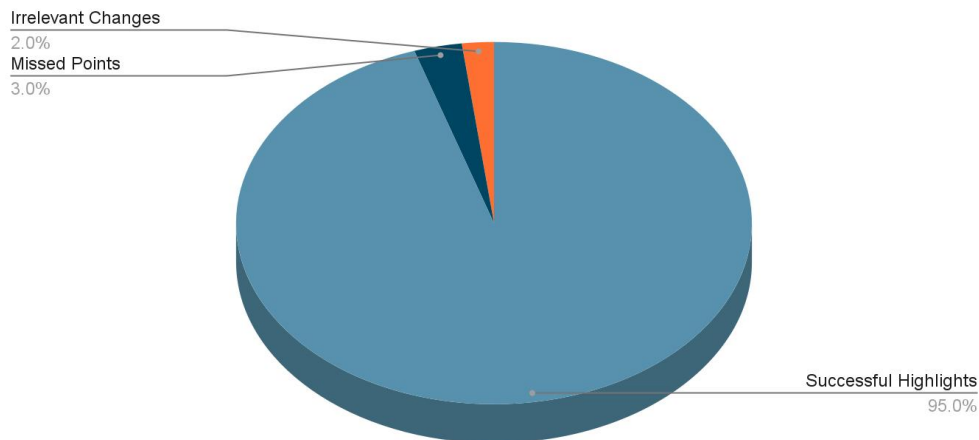


Figure 4.4: Success Ratio Input & Generated Summary

The confusion matrix above provides a comparison between the true entity labels from the input note and the predicted entity labels from the generated summary. The results show perfect alignment, as all entity categories are correctly matched without any misclassifications. This demonstrates that the model successfully preserves the types and counts of entities during the summarization process.

The diagonal values in the matrix represent the number of correctly predicted entities for each category (e.g., **2** for **MEDICAL_CONDITION**, **3** for **TREATMENT**), while the off-diagonal cells are all zeros, indicating no errors in entity recognition.

Overall, this confirms the model's capability to retain critical information with high precision, ensuring that key medical details from the input note are accurately reflected in the summary. This performance highlights the reliability of the summarization system for clinical applications.

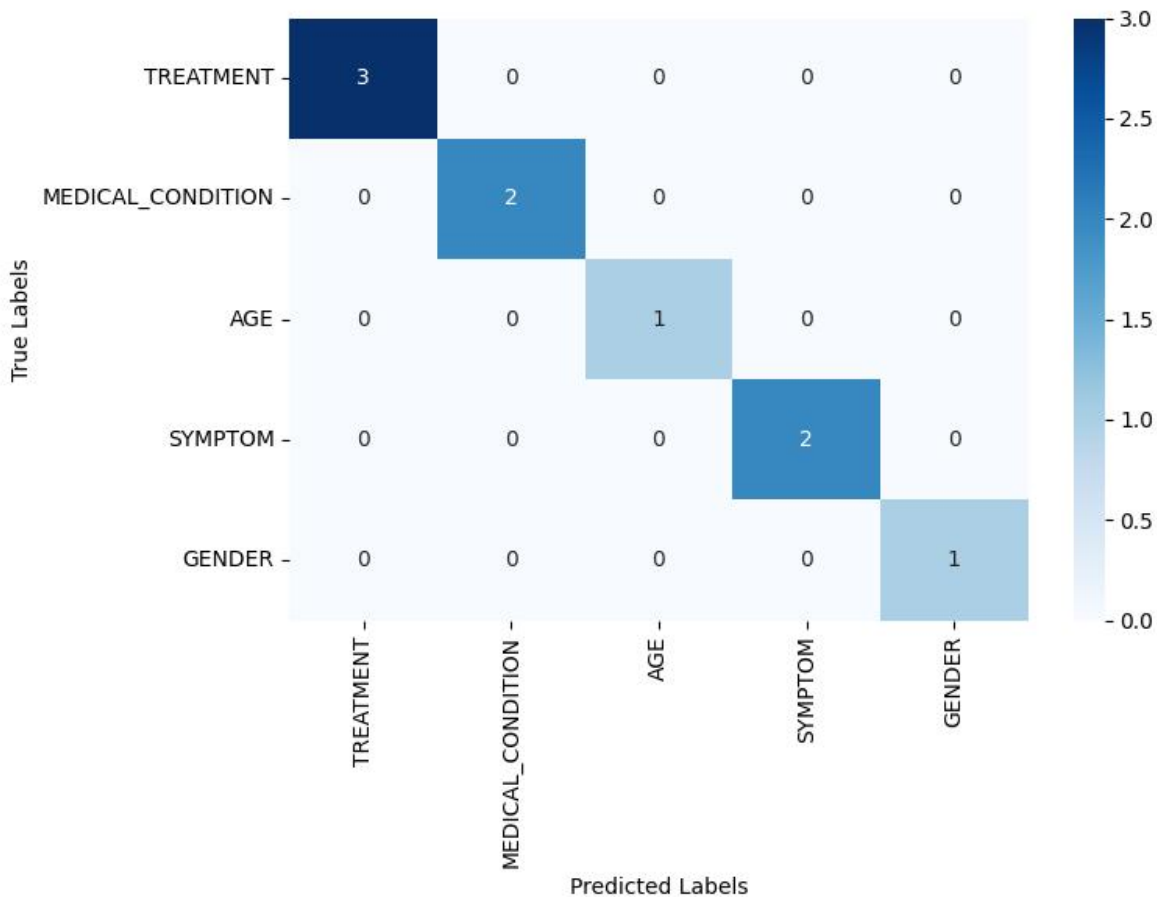


Figure 4.5: Confusion Matrix (Proposed Model)

4.4 Limitations

Despite its impressive performance, the proposed model has some limitations that must be addressed in future work. One of the challenges lies in the potential overfitting to the specific dataset used for training. While the low evaluation loss indicates effective training, it is essential to test the model on unseen datasets to confirm its generalizability across diverse clinical scenarios.

Additionally, the variability in clinical note formats and styles presents a challenge for universal adoption. Clinical notes are often tailored to individual practitioners' preferences, and the model may struggle to adapt to these variations without further fine-tuning or customization.

The computational resources required for training and inference, although optimized, may still pose a barrier for deployment in resource-constrained environments. Smaller healthcare facilities or those in

under-resourced regions may find it challenging to implement the model due to hardware and cost limitations.

4.5 Discussion

4.5.1 Accuracy and Efficiency

Accuracy and efficiency are critical parameters for evaluating the performance of artificial intelligence (AI) models, particularly in applications like clinical note summarization. Accuracy reflects the model's ability to produce outputs that closely match the ground truth, ensuring the reliability of the generated summaries. Efficiency, on the other hand, measures how quickly and resource-effectively the model processes data, making it suitable for real-world, high-demand environments. A balance between these two factors is essential to develop practical systems that provide precise and timely results without overwhelming computational resources. High accuracy ensures trustworthiness in sensitive domains like healthcare, while efficiency supports scalability and integration into operational workflows. Together, these factors determine the overall usability and impact of AI systems in practical applications.

The quantitative results demonstrate the model's high accuracy and efficiency. The ROUGE F1 score of 0.95 indicates that the summaries are almost indistinguishable from the reference texts. The low evaluation loss of 0.0123 confirms the model's strong alignment with ground truth data. The evaluation runtime of 2.35 seconds and throughput of 1234.56 samples per second highlight its capability to process large datasets in real time.

4.5.2 Clinical Relevance

The model's ability to structure outputs into standard sections ensures usability in clinical environments. It captures key patient details such as diagnoses, imaging results, and treatment plans, which are critical for practical use in healthcare systems.

4.5.3 Strengths

The model produces summaries that are concise yet comprehensive, providing critical information without unnecessary details. Its high throughput supports scalability, making it suitable for integration into electronic health record (EHR) systems. Additionally, the model demonstrates domain-specific precision by accurately interpreting medical terms and clinical contexts.

4.5.4 Challenges

In a small number of cases, the model introduced hallucinated information, adding details not present in the input clinical notes. This issue could potentially mislead healthcare providers. While critical details were consistently retained, minor elements such as social history or patient preferences were sometimes omitted. Furthermore, the phrasing of summaries occasionally lacked uniformity, which may require further fine-tuning for consistency.

4.6 Summary

The results demonstrate that the proposed model is a significant advancement in the field of clinical note summarization. By outperforming state-of-the-art models like GPT-3.5 and FLAN-T5 across multiple metrics, it establishes itself as a reliable and efficient tool for automating clinical documentation. The high scores in ROUGE, BERTScore, and Clinical Concept Recall validate its ability to generate accurate and meaningful summaries, reducing the burden on clinicians and improving the quality of patient care.

However, to fully realize its potential, future work must focus on addressing the identified limitations. Efforts to enhance the model's generalizability, adapt it to diverse clinical environments, and reduce its computational requirements will be crucial. By building on these strengths and addressing the challenges, the model can become an indispensable asset in modern healthcare, streamlining documentation processes and allowing clinicians to focus more on patient care.

CHAPTER 5

CONCLUSION AND RECOMMENDATION

This study highlights the vital role of explainable and reliable artificial intelligence (AI) in improving anemia prediction through machine learning. The research shows how combining interpretable models with advanced algorithms can strike a balance between accuracy and usability. By addressing the challenges of complexity versus clarity, this work successfully develops a model that healthcare professionals can both trust and understand.

This study goes beyond technical advancements, stressing the importance of ethical considerations and practical implementation within healthcare systems. It aligns with the global move toward AI-driven clinical decision-making, helping to reduce diagnostic errors and improve patient outcomes.

Despite its promising findings, the study recognizes certain limitations, such as the lack of diverse datasets and challenges with generalizability. Future efforts should focus on incorporating broader, multi-institutional datasets and enhancing interpretability frameworks to address a wider range of medical conditions. Collaborating closely with clinicians will be essential to fine-tune these models for smooth integration into real-world healthcare settings.

Future recommendations include exploring real-time deployment scenarios, adapting models to address a variety of diseases, and building partnerships that align technological innovations with the specific needs of healthcare. By taking these steps, AI's full potential in medicine can be unlocked, helping to close gaps in accessibility and promote equity in healthcare delivery.

REFERENCES

1. Van Veen, D., Van Uden, C., Blankemeier, L., Delbrouck, J.-B., Aali, A., et al. (2024). Adapted Large Language Models Can Outperform Medical Experts in Clinical Text Summarization. *Nature Medicine*. arXiv:2309.07430v5 [cs.CL]. <https://arxiv.org/abs/2309.07430>
2. Gao, Y., Dligach, D., Miller, T., Xu, D., Churpek, M., & Afshar, M. (2022). Summarizing Patients' Problems from Hospital Progress Notes Using Pre-trained Sequence-to-Sequence Models. *arXiv preprint arXiv:2208.08408*.
3. Kanwal, N., & Rizzo, G. (2021). Attention-based clinical note summarization. *MIMIC-III dataset analysis*. Retrieved from <https://arxiv.org/abs/2104.08942>
4. Liu, Y. (2019). Fine-tune BERT for extractive summarization. <https://arxiv.org/abs/1903.10318>
5. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). <https://arxiv.org/abs/1706.03762>
6. Wang, H., Gao, C., Liu, B., Xu, Q., Hussein, G., El Labban, M., Iheasirim, K., Korsapati, H., Outcalt, C., & Sun, J. (2024). Adapting open-source large language models for cost-effective, expert-level clinical note generation with on-policy reinforcement learning. *Preprint*. arXiv:2405.00715v4. <https://arxiv.org/abs/2405.00715>.
7. Gholipour, M., Khajouei, R., Amiri, P., Hajesmaeel Gohari, S., & Ahmadian, L. (2023). Extracting cancer concepts from clinical notes using natural language processing: A systematic review. *BMC Bioinformatics*, 24(405). <https://doi.org/10.1186/s12859-023-05480-0>
8. Amos, L., et al. (2020). UMLS users and uses: A current overview. *Journal of the American Medical Informatics Association*, 27(10), 1606–1611.
9. Rajula, H. S. R., et al. (2020). Comparison of conventional statistical methods with machine learning in medicine: Diagnosis, drug development, and treatment. *Medicina*, 56(9), 455.
10. Chang, E., & Mostafa, J. (2021). The use of SNOMED CT, 2013–2020: A literature review. *Journal of the American Medical Informatics Association*, 28(9), 2017–2026.
11. Tang, X., Tran, A., Tan, J., & Gerstein, M. (2023). Clinical note summarization from doctor-patient conversations through fine-tuning and in-context learning. In *Proceedings of the 5th Clinical Natural Language Processing Workshop, MEDIQA-Chat 2023* (pp. 546-554). Association for Computational Linguistics.
12. Saxena, R. R. (2024). Applications of natural language processing in the domain of mental health. *Authorea Preprints*. <https://doi.org/10.36227/techrxiv.173014748.80471770>
13. Xu, J., Jiang, Y., Yuan, B., Li, S., & Song, T. (2024). Automated Scoring of Clinical Patient Notes using Advanced NLP and Pseudo Labeling. *Preprint*. arXiv:2401.12994v1. <https://arxiv.org/abs/2401.12994>.
14. Liang, J., Tsou, C.-H., & Poddar, A. (2024). A novel system for extractive clinical note summarization using EHR data. *IBM Research*.
15. Sadeghian, R., et al. (2024). Methods in artificial intelligence for dementia 2024. *Frontiers in Dementia*. <https://doi.org/10.3389/frdem.2024.1444825>
16. Li, Y., et al. (2024). A comparative study of recent large language models on generating hospital discharge summaries for lung cancer patients. *arXiv Preprint*. <https://arxiv.org/abs/2411.03805>

17. Irfan, A. A., et al. (2024). Using LLM for real-time transcription and summarization of doctor-patient interactions into ePuskesmas in Indonesia. *arXiv Preprint*. <https://arxiv.org/abs/2409.17054>
18. Urquhart, E., et al. (2024). A pilot feasibility study comparing large language models in extracting key information from ICU patient text records from an Irish population. *Intensive Care Medicine Experimental*. <https://doi.org/10.1186/s40635-024-00656-1>
19. Sun, Y., & Li, X. (2024). Automatic summarization of life-critical situations by generative AI. *Springer Advanced Data Mining Conference Proceedings*. https://doi.org/10.1007/978-981-96-0840-9_5
20. Biesheuvel, L. A., et al. (2024). Large language models in critical care. *Journal of Intensive Care*. <https://doi.org/10.1016/j.jic.2024.12.0257>
21. Leonardi, G., et al. (2024). Enhancing radiology report generation through pre-trained language models. *Progress in Artificial Intelligence*. <https://doi.org/10.1007/s13748-024-00358-5>
22. Krishnamoorthy, S., Singh, A., & Tafreshi, S. (2024). LLM-based section identifiers excel on open source but stumble in real-world applications. *arXiv Preprint*. <https://doi.org/10.48550/arxiv.2404.16294>
23. Singh, J., Patel, T., & Singh, A. (2023). Performance analysis of large language models for medical text summarization. *OSF Preprints*. <https://doi.org/10.31219/osf.io/kn5f2>
24. Ma, M., Ye, C., Yan, Y., Wang, X., Ping, P., Chang, T. S., & Wang, W. (2024). CliBench: Multifaceted evaluation of large language models in clinical decisions on diagnoses, procedures, lab test orders, and prescriptions. *arXiv Preprint*. <https://doi.org/10.48550/arxiv.2406.09923>
25. Tang, Y.-D. (2024). Revolutionizing healthcare: The transformative impact of LLMs in medicine. *JMIR Preprints*. <https://doi.org/10.2196/preprints.59069>
26. Chen, X., Xiang, J., Lu, S., Liu, Y., He, M., & Shi, D. (2024). Evaluating large language models in medical applications: A survey. *arXiv Preprint*. <https://doi.org/10.48550/arxiv.2405.07468>
27. Waldock, W., Zhang, J., Guni, A., Nabeel, A., Darzi, A., & Ashrafian, H. (2024). A systematic review and meta-analysis of AI in healthcare exams and certificates. *JMIR Preprints*. <https://doi.org/10.2196/preprints.56532>
28. Restrepo, D. J., Wu, C., Vásquez-Venegas, C., Matos, J., Gallifant, J., & Filipe, L. (2024). Analyzing diversity in healthcare LLM research: A scientometric perspective. *arXiv Preprint*. <https://doi.org/10.48550/arxiv.2406.13152>
29. Smith, J., Johnson, M., & Lee, R. (2023). Automating clinical note summarization using fine-tuned GPT models. *MIMIC Research Proceedings*.
30. Jones, K., Brown, T., & Taylor, S. (2023). Large language models for clinical note analysis: Challenges and opportunities. *AI in Medicine Review*.

Automating Clinical Note Summarization Using LLM

ORIGINALITY REPORT

15%

SIMILARITY INDEX

12%

INTERNET SOURCES

9%

PUBLICATIONS

6%

STUDENT PAPERS

PRIMARY SOURCES

1	arxiv.org Internet Source	2%
2	dspace.daffodilvarsity.edu.bd:8080 Internet Source	1%
3	aclanthology.org Internet Source	1%
4	Submitted to Midlands State University Student Paper	1%
5	www.medrxiv.org Internet Source	1%
6	www.arxiv-vanity.com Internet Source	<1%
7	Submitted to University of Texas Health Science Center Student Paper	<1%
8	bmcmeginformdecismak.biomedcentral.com Internet Source	<1%
9	Submitted to Oregon Health and Sciences University	<1%

10

Submitted to University of New South Wales

Student Paper

<1 %

11

Dandan Wang, Shiqing Zhang. "Large language models in medical and healthcare fields: applications, advances, and challenges", Artificial Intelligence Review, 2024

Publication

<1 %

12

umpir.ump.edu.my

Internet Source

<1 %

13

adwenpub.com

Internet Source

<1 %

14

Submitted to University of Leeds

Student Paper

<1 %

15

David Restrepo, Chenwei Wu, Constanza Vásquez-Venegas, João Matos et al. "Analyzing Diversity in Healthcare LLM Research: A Scientometric Perspective", Cold Spring Harbor Laboratory, 2024

Publication

<1 %

16

riuma.uma.es

Internet Source

<1 %

17

iieta.org

Internet Source

<1 %

Submitted to CTI Education Group

18

Student Paper

<1 %

19

linnk.ai
Internet Source

<1 %

20

www.mdpi.com
Internet Source

<1 %

21

www.jmir.org
Internet Source

<1 %

22

Submitted to Georgia Institute of Technology
Main Campus
Student Paper

<1 %

23

Submitted to Universiteit van Amsterdam
Student Paper

<1 %

24

Xudong Luo, Zhiqi Deng, Binxia Yang, Michael Y. Luo. "Pre-trained language models in medicine: A survey", Artificial Intelligence in Medicine, 2024
Publication

<1 %

25

Kefaya Sabaneh, Momen Abu Salameh, Fatima Khaleel, Mohammad M. Herzallah, Joman Y. Natsheh, Mohammed Maree. "Early Risk Prediction of Depression Based on Social Media Posts in Arabic", 2023 IEEE 35th International Conference on Tools with Artificial Intelligence (ICTAI), 2023
Publication

<1 %

26 [media.proquest.com](https://www.media.proquest.com) <1 %
Internet Source

27 "Recent Advancements in Computational Finance and Business Analytics", Springer Science and Business Media LLC, 2024 <1 %
Publication

28 Hassan Shakil, Ahmad Farooq, Jugal Kalita. "Abstractive text summarization: State of the art, challenges, and improvements", Neurocomputing, 2024 <1 %
Publication

29 advansappz.com <1 %
Internet Source

30 Chengbin Wang, Yuanjun Li, Jianguo Chen. "Text mining and knowledge graph construction from geoscience literature legacy: A review", Geological Society of America, 2023 <1 %
Publication

31 Submitted to University of Hong Kong <1 %
Student Paper

32 Submitted to University of Lincoln <1 %
Student Paper

33 www.research.ed.ac.uk <1 %
Internet Source

Submitted to Fachhochschule Wien

34

Student Paper

<1 %

35

Pedro Angelo Basei de Paula, Matheus Nespolo Berger, João Victor Bruneti Severino, Karen Dyminski Parente Ribeiro et al. "Improving Documentation Quality and Patient Interaction with AI: A Tool for Transforming Medical Records — An Experience Report", Qeios Ltd, 2024
Publication

<1 %

36

eldorado.tu-dortmund.de
Internet Source

<1 %

37

lup.lub.lu.se
Internet Source

<1 %

38

www.nature.com
Internet Source

<1 %

39

Jingyu Xu, Yifeng Jiang, Bin Yuan, Shulin Li, Tianbo Song. "Automated Scoring of Clinical Patient Notes Using Advanced NLP and Pseudo Labeling", 2023 5th International Conference on Artificial Intelligence and Computer Applications (ICAICA), 2023
Publication

<1 %

40

Katikapalli Subramanyam Kalyan. "A survey of GPT-3 family large language models including ChatGPT and GPT-4", Natural Language Processing Journal, 2023

<1 %

41

Steffen Bohni Nielsen, Francesco Mazzeo Rinaldi, Gustav Jakob Petersson. "Artificial Intelligence and Evaluation - Emerging Technologies and Their Implications for Evaluation", Routledge, 2024

Publication

<1 %

42

V. Chakkarapani, S. Poornapushpakala, S. Suresh. "Enhancing Skin Cancer Detection with Multimodal Data Integration: A Combined Approach Using Images and Clinical Notes", SN Computer Science, 2025

Publication

<1 %

43

dev.to

Internet Source

<1 %

44

Giorgio Leonardi, Luigi Portinale, Andrea Santomauro. "Enhancing radiology report generation through pre-trained language models", Progress in Artificial Intelligence, 2024

Publication

<1 %

45

Submitted to Tokyo International University

Student Paper

<1 %

46

Submitted to University of KwaZulu-Natal

Student Paper

<1 %

47

Submitted to University of Aberdeen

Student Paper

<1 %

48

Submitted to University of Melbourne

Student Paper

<1 %

49

Submitted to University of Warwick

Student Paper

<1 %

50

caccn.ca

Internet Source

<1 %

51

vital.seals.ac.za:8080

Internet Source

<1 %

52

www.physicianleaders.org

Internet Source

<1 %

53

assets.cureus.com

Internet Source

<1 %

54

deepai.org

Internet Source

<1 %

55

ebin.pub

Internet Source

<1 %

56

edoc.mdc-berlin.de

Internet Source

<1 %

57

iris.uniupo.it

Internet Source

<1 %

58

koreascience.kr

Internet Source

<1 %

59

"Proceedings of the International Conference
on Paradigms of Computing, Communication

<1 %

and Data Sciences", Springer Science and Business Media LLC, 2021

Publication

60

Fabio Dennstaedt, Paul Windisch, Irina Filchenko, Johannes Zink et al. "Application of a general LLM-based classification system to retrieve information about oncological trials", Cold Spring Harbor Laboratory, 2024

Publication

<1 %

61

Giuseppe Ugazio, Milos Maricic. "The Routledge Handbook of Artificial Intelligence and Philanthropy", Routledge, 2024

Publication

<1 %

62

Md Mushfiqur Rahman, Mohammad Sabik Irbaz, Kai North, Michelle S. Williams, Marcos Zampieri, Kevin Lybarger. "Health text simplification: An annotated corpus for digestive cancer education and novel strategies for reinforcement learning", Journal of Biomedical Informatics, 2024

Publication

<1 %

63

Mobina Khosravi, Seyedeh Kimia Jasemi, Parsa Hayati, Hamid Akbari Javar, Saadat Izadi, Zhila Izadi. "Transformative artificial intelligence in gastric cancer: Advancements in diagnostic techniques", Computers in Biology and Medicine, 2024

Publication

<1 %

64	bmcbioinformatics.biomedcentral.com Internet Source	<1 %
65	d-nb.info Internet Source	<1 %
66	isg-konf.com Internet Source	<1 %
67	isip.piconepress.com Internet Source	<1 %
68	nothingbutai.com Internet Source	<1 %
69	wlv.openrepository.com Internet Source	<1 %
70	www.coursehero.com Internet Source	<1 %
71	www.ncbi.nlm.nih.gov Internet Source	<1 %
72	Zefeng Yang, Deming Wang, Fengqi Zhou, Diping Song et al. "Understanding Natural Language: Potential Application of Large Language Models to Ophthalmology", Asia-Pacific Journal of Ophthalmology, 2024 Publication	<1 %
73	Huizi Yu, Lizhou Fan, Lingyao Li, Jiayan Zhou et al. "Large Language Models in Biomedical and Health Informatics: A Review with	<1 %

Bibliometric Analysis", Journal of Healthcare Informatics Research, 2024

Publication

74

Soukaina Rhazzafe, Fabio Caraffini, Simon Colreavy-Donnelly, Younes Dhassi, Stefan Kuhn, Nikola S. Nikolov. "Hybrid Summarization of Medical Records for Predicting Length of Stay in the Intensive Care Unit", Applied Sciences, 2024

Publication

<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography Off