

## **Natural Language Processing**



## **Medical Transcriptions**

*Namratha Bolar* [nqb5534@psu.edu](mailto:nqb5534@psu.edu)

*Venkata Naga Sai Sri Vaishnavi Appala* [vka5190@psu.edu](mailto:vka5190@psu.edu)

*Sai Sruthi Avula* [spa5940@psu.edu](mailto:spa5940@psu.edu)

## **School of Graduate Professional Studies**

MPS/MS in Data Analytics

March, 2025

## Document Control

Work carried out by:

Name	Email Address	Other
Namratha Bolar	<a href="mailto:nqb5534@psu.edu">nqb5534@psu.edu</a>	
Venkata Naga Sai Sri Vaishnavi Appala	<a href="mailto:vka5190@psu.edu">vka5190@psu.edu</a>	
Sai Sruthi Avula	<a href="mailto:spa5940@psu.edu">spa5940@psu.edu</a>	

## Revision Sheet

Release No.	Date	Revision Description
1.0	03/23/2025	Defined the Project Proposal, Challenges, Solutions & Improvements
2.0	03/30/2025	Added Six Literature Reviews
3.0	04/06/2025	Updated the Literature Review References and added the dataset details
4.0	04/13/2025	Description of the Deep Learning Process and the Deep Neural Networks for the dataset
5.0	04/20/2025	Description of the Deep Learning Process, Enhanced and added description of Train/Test Strategy, Evaluation Metrics
6.0	04/27/2025	Updated the Deep Learning Process and added the codes and results of three Models
7.0	04/30/2025	Updated the Model Results and Evaluation Metrics and added references.

**Note:** For the project assignment deadline submission please refer to the syllabus. Each week's assignment is for 10 points and the project assignment is worth 40% of your final grade. The project presentation is worth 5% of your final grade. The course project should be completed as a team with no more than 3 members (in each team).

## TABLE OF CONTENTS

Project Objective and Proposal	5
Description and Challenges	6
Literature review	8
Dataset	18
Deep Neural Network Process	19
Evaluation Metrics	21
Results	25
Conclusion	30
References	38

## **General Guidelines**

1. Please use this template document to complete each deliverable assignment.
2. Each assignment must be submitted by the due date in the Course Schedule.
3. All figures should be followed by a brief description of the figure.
4. Figures can be hand-drawn and scanned in some circumstances, but the hand-drawn figure should be clear and legible to obtain full credit. Unclear hand-drawn figures will receive partial credit. For constructing figures and diagrams it is advised to use tools.
5. Figures and tables should have appropriate captions. For documenting and referencing styles please follow the APA or MLA writing style.
6. Please make sure that you provide a reference section.
7. Any material text or figure taken from books, journals or the Internet should be referenced. If you have a sentence or a figure that does not belong (authorship) to you, they must be clearly referenced. If you fail to do so your report will be considered as a case for plagiarism. It is your duty to make sure that your report is free from any activity related to plagiarism. Please see the section on Academic Integrity found below. The penalty for plagiarism will be a “0” awarded to your report. Thus, it is good to keep it simple, always have the principle to acknowledge people for their contributions.
8. Please note that the use of ChatGPT is prohibited. If the use of ChatGPT to complete the assignment has been determined, then you will receive NO credit for your project.

Throughout this course you will complete a team project where you will work on a problem from a *Neural Natural Processing perspective* with at least two of your classmates. You will be allowed to team up with other students in the class at your choice.

The project will consist of 3 parts (Specific details are in the section below):

- Project Proposal.
- Project Deliverables: Report, PowerPoint Slides, Code, etc (see details below)
- In-class/zoom Presentation

## Project Objectives

Defining and developing an NLP project is an important objective in this course. Teams are encouraged to define their own data science problems and demonstrate how to solve them with deep neural networks. This includes the collection and processing of datasets, building neural networks, training and evaluating their performance.

The first step is to collect or find datasets relevant to your project. You have many options. One option is to work on real-world data related to your job to get new insights into a problem of interest to your employer. Another option is to find dataset online. Fortunately, many open-data source platforms provide large datasets, covering a wide range of applications. You might be interested in checking websites such as:

- [Kaggle.](#)
- [Google Dataset Search.](#)
- [Awesome Public Datasets.](#)
- etc.

You are allowed to download and re-use existing datasets and extended code of existing projects. In all cases, you should explicitly refer to the source of your datasets and explain how you re-use code of existing projects. All projects submissions will be verified by [Turnitin](#) to check the originality of your work.

Defining the project is an iterative process, So, don't hesitate to contact me to discuss your ideas and frame your project objectives.

## Project Proposal

The project proposal will include:

- Project Title
- Team members
- The Problem from NLP Perspective
  - Description
  - Challenges

- Is there any available solution to this problem?
- Dataset: source, size, features description, etc.
- Deep Learning process
  - Choices of Deep Neural Networks and their justification
  - Training/Testing strategy
  - Evaluation metrics
- Expected results
- References to data source
- Justification of reusing existing code

The project proposal will be completed within the duration of week 1 to week 6

### **Project Assignment 1 (10 points)**

Provide the following details related to your project

- **Project Title : Medical Transcriptions**
- Team members: Venkata Naga Sai Sri Vaishnavi Appala, Sai Sruthi Avula and Namratha Bolar

#### **Description:**

Medical transcription contains critical patient information, including clinical diagnoses, treatment plans, and doctor-patient interactions. These transcriptions are often not reliable as they are unstructured and contain complex medical terminology, abbreviations, and contextual nuances which require a lot of manual processing in turn makes it time-consuming and prone to errors.

Here the purpose is to leverage Natural Language Processing (NLP) techniques to automatically process, analyze, and extract meaningful insights from these medical transcription data. The aim is to use Deep learning models and NLP methodologies, to improve the accuracy, efficiency, and accessibility of medical records.

#### **Challenges:**

Some of the important challenges we are facing are:

1. **Noisy & Unstructured text:** Medical transcriptions often contain errors, incomplete sentences, and inconsistencies due to human or automated transcription mistakes.
2. **Ambiguity in Medical Terms:** There are different meanings for each word based on the context making entity recognition challenging.
3. **De-identification of Patient Information:** The sensitivity of the patient details and privacy compliance need to be handled.

4. **Domain Adaptation:** The model might not perform well on medical transcriptions without fine-tuning domain-specific data.
5. **Lack of Standardized metrics:** Defining precise evaluation metrics for medical transcription analysis is challenging.

### **Available Solutions and Improvements:**

Yes, there are available solutions for this project. For example: Google's Medical AI (Google Health & Med-PaLM 2)

Google's AI models, like Med-PaLM, focus on understanding and summarizing medical text, including transcriptions.

These models have shown state-of-the-art performance in medical question-answering but are not specifically trained on physician transcriptions.

Amazon Transcribe Medical: AWS offers a HIPAA-compliant automatic speech recognition (ASR) system for medical dictation.

### **Project Assignment 2 (10 points)**

Provide the following details related to your project

- Conduct literature survey to identify works conducted in the past (if any) or any related work. I will recommend going through at least 5 journal/conference articles.
- Summarize those literature works indicating what work has been pursued, what tools and techniques have been employed in that study, reported performance measures and their finding.
- Conclude those summaries (for each surveyed article) of what you have learnt from their study.

Students will be provided with some examples of projects completed in the past.

### **Natural Language Processing to Identify the Creation and Impact of New Technologies in Patent text: Code, data, and new measures**

The objective of this paper is to focus on using NLP to identify the creation and impact of new technologies in the patent text. Uses computational techniques to analyze documents, extract insights, and determine the influence of emerging technologies.

### **Tools and techniques used:**

**NLP:** text mining, entity recognition, and topic modeling are used. **ML algorithm:** A couple of supervised and unsupervised models are applied to classify patent data and detect technologies. **Latent Dirichlet Allocation (LDA):** A topic modeling approach to identify themes in patent documents. **Word Embeddings (e.g., Word2Vec, BERT):** These techniques help in semantic analysis and trend detection.

### **Performance measures:**

Using accuracy scores like precision, recall, and F1 score along with Topic coherence scores. A couple of patent citation Metrics with network centrality measures.

### **Findings:**

**Emerging Technologies:** The study successfully identifies emerging technology trends within patents through NLP-driven approaches.

**Impact Assessment:** Citation-based analysis provides insights into the influence of patents and their contribution to technological advancement.

**Effectiveness of NLP:** The application of NLP enhances the ability to extract meaningful insights from complex patent text, making technology forecasting more efficient.

### **Reference:**

Shin, D., Moon, J., & Hwang, H. (2022). Identifying the creation and impact of new technologies in patent text using NLP and machine learning. *Discover Artificial Intelligence*, 2(1), 6.  
<https://doi.org/10.1007/s44163-022-00006-3>

## **Advancements in natural language processing: Implications, challenges, and future directions**

The Objective of this paper explore NLP applications, implications, and advancements. It discussed different methodologies in key areas like sentiment analysis, text classification, machine translation, and information retrieval.

### **Tools and technologies:**

deep learning (e.g., CNNs, RNNs, Transformers), traditional statistical models, and rule-based linguistic methods. Python-based NLP libraries (NLTK, spaCy, and Transformer-based models

### **Performance measures:**

Standard metrics such as accuracy, precision, recall, F1-score, and BLEU scores for translation tasks.

### **Findings:**

The literature review highlights that while deep learning models, particularly Transformers, have significantly advanced NLP performance, challenges remain in terms of contextual understanding.



## Reference:

Kaur, H., Tyagi, V., & Garg, R. (2021). *NLP Implementation: Current State, Challenges, and Perspectives*. 2021 6th International Conference on Signal Processing, Computing and Control (ISPCC), 383–388. IEEE. <https://doi.org/10.1109/ISPCC53510.2021.9609373>

## Using Natural Language Processing to Identify Stigmatizing Language in Labor and Birth Clinical Notes.

### Introduction:

In this paper disparities in maternal and perinatal outcomes particularly among racially and socially minoritized groups were well documented. While racism and discrimination are widely recognized as key contributors to these disparities, less attention has been given to the role of clinician bias as conveyed through language in electronic health records (EHRs). Negative language in clinical records can discreetly yet significantly impact patient treatment and healthcare provider perspectives. It frequently indicates underlying bias and strengthens detrimental societal hierarchies, particularly when the language pertains to race, immigration status, economic background, or compliance with conventional reproductive norms. Two main categories of stigmatizing language have been recognized: marginalizing language, which portrays patients negatively and privilege language, which emphasizes traits that conform to prevailing social standards. Previous studies on this topic have primarily concentrated on areas such as internal medicine and psychiatry, frequently employing qualitative methods or basic keyword searches, birth environments. Acknowledging this gap, the current research aimed to create and assess machine learning-driven natural language processing methods for the automatic identification of stigmatizing language in clinical notes related to labor and birth.

### Challenges:

This research is based on several important challenges. To begin with, there is no common, established technique for recognizing clinician bias in relation to documentation of birth-related healthcare. Conventional rule-based methods for NLP, like keyword spotting, have a narrow focus and are susceptible to false positive because they fail to consider context and tone. Moreover, existing literature on stigmatizing language in EHRs is sparse, with most studies focusing on other medical domains and relying on human review, which is time-consuming and not scalable. Another challenge lies in the lack of annotated datasets specific to labor and birth notes, which are essential for training supervised learning models. Furthermore, common text representation techniques like TF-IDF, while effective for certain classification tasks, may struggle to capture the nuances of word meaning, syntax, and sentence structure that are essential for understanding implicit bias and stigmatization in language. Lastly, the study is constrained by the limited generalizability of data from only two hospitals and a relatively small number of labeled notes.

### Data:

The dataset included notes from individuals who were more than 20 weeks pregnant and admitted for labor and delivery. A total of 1,117 clinical notes were reviewed. These notes were selected from seven document types like Admission Notes, Triage Notes, Nursing Notes, Postpartum Notes, deemed most likely to contain stigmatizing language. Reviewers labeled instances of both marginalizing and power/privilege language using a combination of theory-driven and emergent

coding strategies. Marginalizing language included subcategories such as questioning patient credibility, showing disapproval, stereotyping, labeling individuals as “difficult” and imposing unilateral clinical decisions. Power language referred to descriptors that elevated the patient’s social status, such as stable employment, nurturing marriage, or being appropriately groomed. Overall, 232 notes contained stigmatizing content: 205 with marginalizing language and 37 with power/privilege language. The high interrater agreement (Cohen’s  $K > 0.8$ ) indicated strong consistency among reviewers and validated the quality of the annotated dataset.

### **Methods Employed:**

The methodology combined qualitative human annotation with a machine learning-based NLP pipeline. In the first phase, a team of four expert reviewers used a qualitative descriptive approach to label stigmatizing language in the notes. Each note was reviewed by at least two experts, and team consensus was used to resolve discrepancies. Thematic content analysis was applied to identify linguistic patterns, and notes were classified based on the presence of marginalizing or power language. Once this annotated dataset was finalized, it served as the training and testing foundation for NLP modeling.

In the second phase, the text data underwent standard preprocessing steps such as converting all text to lowercase and removing punctuation and numeric characters. The text was then converted into a numerical format using Term Frequency-Inverse Document Frequency (TF-IDF), a method that assigns importance to words based on how frequently they appear in a document relative to the entire corpus. The TF-IDF vectors were split into training (70%) and testing (30%) sets using stratified sampling to preserve the proportion of stigmatizing categories. Three machine learning algorithms were tested to evaluate their ability to classify whether a note contained marginalizing or power/privilege language: Decision Trees (J48), Random Forests, and Support Vector Machines (SVM). These models are well-established for NLP tasks and offer a balance between performance and interpretability. Each model was trained and evaluated separately for the two types of stigmatizing language.

Model performance was assessed using precision, recall, and F1-score. Additionally, the researchers used the Information Gain (InfoGain) method to identify which words were most influential in making classification decisions. This helped interpret the results and identify linguistic patterns that aligned with the expert-coded categories. For example, words such as “illicit,” “marijuana,” and “public housing” were strong predictors of marginalizing language, while terms like “appropriately groomed,” “employed,” and “cordial” signaled power/privilege language.

### **Results and Conclusion:**

The study found that machine learning algorithms could effectively classify stigmatizing language in clinical notes. For marginalizing language, Decision Trees performed best, achieving an F1 score of 0.73, indicating a good balance between precision and recall. For power/privilege language, Support Vector Machines outperformed the other models, reaching an F1 score of 0.91, demonstrating high accuracy in identifying positive social descriptors. The models were able to detect language patterns that aligned with expert-coded annotations and showed promising generalization on the test set.

The implications of these findings are both practical and academic. This is the first study to apply and evaluate machine learning-based NLP for identifying stigmatizing language specifically in labor and birth settings. The approach outperforms rule-based methods by capturing contextual

nuances and reducing false positives. The technology has potential applications in real-time clinical monitoring systems, where algorithms could flag stigmatizing language and alert clinicians to revise their notes. This could help reduce bias, improve patient trust—especially with the growing patient access to clinical notes under federal law—and ultimately enhance the quality of care. Furthermore, the feature analysis revealed that stigmatizing language was disproportionately associated with notes that referenced substance use, immigrant status, or lower socioeconomic conditions, echoing broader concerns about structural bias in healthcare.

The study also highlights several areas for future work. The authors suggest exploring more advanced NLP methods, such as transformer-based models (e.g., BERT or GPT), which can better capture semantic and syntactic context. Additionally, future research should examine whether stigmatizing language use varies by clinician characteristics or patient demographics, particularly across racial and ethnic lines. Expanding the dataset to include more hospitals and a broader range of clinical settings could also improve the generalizability of the findings. In summary, this study represents a foundational step in leveraging machine learning and NLP to detect and mitigate bias in clinical documentation, with potential to contribute significantly to addressing disparities in maternal and perinatal healthcare.

### **Reference:**

Barcelona, V., Scharp, D., Moen, H., Davoudi, A., Idnay, B. R., Cato, K., & Topaz, M. (2023). Using natural language processing to identify stigmatizing language in labor and birth clinical notes. *Maternal and Child Health Journal*, 28(3), 578–586. <https://doi.org/10.1007/s10995-023-03857-4>

## **Natural Language Processing of Clinical Notes on Chronic Diseases: Systematic Review**

### **Introduction:**

The paper begins by acknowledging the increasing burden of chronic diseases such as cancer, diabetes, hypertension, and heart conditions, which present a growing challenge to modern healthcare systems worldwide. While treatments and preventive measures have progressed, the rising incidence of these conditions demands innovative solutions that go beyond traditional evidence-based medicine. One promising approach is the secondary use of electronic health records (EHRs) to analyze patient data, uncover clinical patterns, and improve decision-making. Chronic diseases, due to their long-term nature, generate vast amounts of longitudinal data, including free-text clinical notes that capture nuanced, patient-specific information often missing in structured fields. However, the unstructured nature of clinical narratives poses challenges for direct computational use, which has sparked a growing interest in natural language processing (NLP) to convert these free-text records into structured, machine-readable data. The paper emphasizes that developing robust NLP methods to extract clinically meaningful insights from such narratives is essential for early detection, risk stratification, and personalized care in chronic disease management.

### **Challenges Identified:**

The review identifies multiple challenges that hinder the full potential of NLP in clinical research on chronic diseases. Firstly, clinical notes are highly variable, noisy, and sparse, making it difficult for algorithms to generalize across datasets. Many rule-based NLP systems rely on handcrafted dictionaries or ontologies, which are not only labor-intensive to build but often fail to capture the

linguistic nuances in clinical documentation. Furthermore, while machine learning methods are gaining traction, they are often underutilized due to data privacy issues, lack of large annotated corpora, and the clinical domain's preference for interpretability over predictive power. The scarcity of publicly available datasets further limits progress in this space, particularly for deep learning models that require large volumes of training data. Another major barrier is the limited ability of current systems to handle temporal information, entity relationships, and the transition from simple entity recognition to semantic understanding. Lastly, the field lacks a unified approach for integrating structured and unstructured data, which is crucial for modeling patient trajectories and comorbidities.

### **Data:**

This study is a systematic review that synthesizes findings from 106 selected research articles out of an initial pool of 2,652 identified through database searches following PRISMA guidelines. The included papers span work on 43 unique chronic diseases, which were grouped into 10 categories using the ICD-10 classification system. The most represented conditions were diseases of the circulatory system (n=38), neoplasms (n=34), and endocrine/metabolic disorders (n=14). Notably, the paper observes that although metabolic diseases are more prevalent in the population, they are underrepresented in the NLP literature due to the structured nature of their clinical data, unlike circulatory diseases, which are documented more in unstructured text. Most datasets analyzed in the reviewed studies were non-public institutional datasets, though some used standard open datasets like i2b2, MIMIC-II, THYME, and DeepPhe. These datasets typically include free-text clinical narratives, discharge summaries, radiology/pathology reports, and progress notes.

### **Methods:**

The review classifies the methods employed in the reviewed studies into rule-based, machine learning-based, and hybrid approaches. Among the rule-based systems, dictionary lookups, domain-specific ontologies, manually written rules, and regular expressions were common techniques for entity recognition and information extraction. These methods were favored for their high interpretability, which is critical in clinical applications but limited in scalability and generalization. On the machine learning side, support vector machines (SVMs) were the most frequently used algorithm, followed by Naïve Bayes, conditional random fields (CRFs), and random forests. These algorithms were used for tasks such as disease classification, phenotype recognition, and identifying risk factors or comorbidities. Only three studies employed deep learning methods, despite their growing success in general NLP, likely due to the data and interpretability constraints in the clinical domain. A significant number of studies adopted hybrid models that combined machine learning with rule-based elements, for instance, using SVMs with manually curated lexicons or combining CRFs with temporal rule sets for tracking disease progression. The main NLP tasks identified were text classification, entity recognition, negation detection, coreference resolution, and temporal information extraction.

### **Results and Conclusion:**

The systematic review reveals a clear increase in the adoption of machine learning techniques for analyzing clinical notes in chronic disease research, though rule-based methods still dominate due to their simplicity and transparency. The most common application of NLP was to identify risk factors or disease phenotypes from clinical notes, while fewer studies focused on extracting

comorbidities or longitudinal patterns. Surprisingly, only a few studies used word embeddings or deep learning architectures, which the authors attribute to the lack of large-scale annotated corpora and the difficulty of interpreting deep models in clinical settings. Most studies still relied on shallow classifiers due to their lower data demands and easier integration into clinical workflows. The review also highlighted a gap between information extraction and understanding, with a need to shift from recognizing isolated entities to modeling their interactions, temporal context, and causal relationships.

In conclusion, the paper underscores the urgent need for more advanced and holistic NLP solutions that go beyond entity extraction to incorporate clinical reasoning, temporal modeling, and conceptual understanding. Recommendations include advancing toward semantic-level NLP, enabling relationship and temporal extraction, exploiting alternative knowledge sources (like textbooks or decision support tools), and promoting the creation of large, annotated, de-identified corpora through shared tasks or patient data donation initiatives. The authors expect the role of deep learning in clinical NLP to expand as more data becomes available and methods evolve to be more interpretable and trustworthy in healthcare environments.

### **References:**

Sheikhalishahi, S., Miotto, R., Dudley, J. T., Lavelli, A., Rinaldi, F., & Osmani, V. (2019). Natural Language Processing of Clinical Notes on Chronic Diseases: Systematic Review. *JMIR Medical Informatics*, 7(2), e12239. <https://doi.org/10.2196/12239>

## **Natural Language Processing in Medicine and Ophthalmology: A Review for the 21st-Century Clinician**

### **Introduction:**

The article offers a comprehensive overview of how Natural Language Processing (NLP) is transforming the medical community. It emphasizes the critical role of NLP in getting insights from unstructured clinical text such as Electronic Health Records (EHRs), which plays a vital role in the medical data and health sector. The authors highlight various NLP tasks like including Named Entity Recognition (NER), text summarization, question answering, and topic modeling and using their applications in clinical decision support, documentation automation, and patient communication. Currently special attention is getting added to the fields of the Large Language Models like GPT, and Clinical BERT showing good diagnostic accuracy and simplifying the complex medical content in every aspect. By focusing on real-world use cases, particularly in ophthalmology, the paper underscores the potential of NLP to improve efficiency, reducing the manual work, time and support data-driven decision-making in healthcare. It also acknowledges current challenges such as data quality, model bias, enhancing the integration of NLP technologies into clinical practice.

### **What work had been pursued:**

Natural Language Processing (NLP) has been explored and applied across multiple factors of modern medicine, focusing particularly on ophthalmology. It bridges the gap between cutting-edge AI technologies and healthcare applications, with a good impact on real-world implementation



potential. The paper defines the foundational components of NLP, such as Natural Language Understanding (NLU) and Natural Language Generation (NLG), and discusses how these components support aspects like entity recognition, text classification, summarization, semantic analysis and mapping them to the clinical functions. It then maps these tasks to their clinical functions, for example, using NER to extract diagnoses, treatments, and lab test results from unstructured clinical notes. A major focus of the work is on the role of Large Language Models (LLMs), particularly models like BERT, ClinicalBERT, GPT, and BioCPT, dealing with complex language tasks in the medical domain. The paper explores how these models are pre-trained on massive datasets and fine-tuned to perform question answering automated diagnosis, and information extraction with high accuracy. Followed up with Specialized applications like Automated triaging of ophthalmology referrals using CNNs and Simplifying pathology reports using GPT-4 and Bard followed up with Sentiment and emotion analysis in online patient discussions using IBM Watson NLP and integration of chatbots.

### **Tools and Techniques:**

A wide range of Natural Language Processing (NLP) tools and techniques have been developed and employed for clinical applications, in interpreting and extraction of information. Classical NLP Preprocessing Techniques consist of Tokenization, Stemming, and Lemmatization-grams and Bag-of-Words Models, TF-IDF, Word Embeddings and Contextual Representations define Word2Vec and GloVe, Contextual Embeddings via Transformers and Named Entity Recognition (NER) Techniques are Rule-based NER, Machine Learning-based NER models which contains categories like Support Vector Machines (SVM), Conditional Random Fields (CRF), Long Short-Term Memory (LSTM), Bidirectional LSTM (BiLSTM), Transformer-based NER, SpaCy Along with all these IBM Watson NLP Toolkit is applied to perform sentiment and emotion analysis on patient discussions and Domain-adapted models trained on biomedical literature to enhance the EHR's accuracy in clinical NLP Tasks.

### **Reported Performance Measures:**

For the Named Entity Recognition (NER), for domain-specific knowledge where the results show that F1-Score achieved more than 85% across all medical entity types like symptoms, treatments, and lab tests. Incorporating medical knowledge into the embedding layer significantly improved NER performance. CNN-Based Triage System, Used NLP to categorize ophthalmology referrals into urgent and non-urgent cases achieved an accuracy of 81%, and got the area under curve (AUC) with a value 0.83. This results shows the best accuracy values of GPT-4 with 97.4% and Bard with 87.6%. This tells that Chatbots improved readability and patient understanding. Sentiment and Emotion Analysis, based on the IBM Watson NLP on Patient Forums applied to online discussions about ophthalmology procedures shows the effectiveness in capturing emotional responses related to the medical conditions and procedures. NLP can reduce documentation work, improving communication and fast detection with high accuracy and generalization.

### **Conclusion:**

NLP is Essential for Extracting Value from Unstructured Medical Text: The data in free text format within Electronic Health Records had been underutilized and provides the tools to structure, summarize, and analyze the data making it useful for the tasks, diagnosis, and research.

Enhancement of NLP Capabilities with Large Language Models: Advance transformer-based models like GPT, and BERT have improved the performance of NLP in medicine by handling the complex language, and understanding the clinical context. Model Performance can be high but depends on Data Quality: Clean, representative domain-specific datasets are necessary for the model performance and successful deployments of NLP systems in healthcare. Adoption of NLP Applications is Still Limited: While NLP is being used in areas like diagnostics, sentiment analysis, triage, clinical trials, and even public health monitoring, real-world clinical adoption remains low (~5%) due to challenges like model bias, real-world clinical adoption remains low (~5%) due to challenges like model bias. Totally underlying the immense potential of NLP in revolutionizing healthcare, especially when integrated thoughtfully and responsibly. It emphasizes that while technology is maturing rapidly with accuracy and diverse applications, real-world integration still demands attention to model reliability, data ethics, interpretability, and clinical validation.

### **References:**

Rojas-Carabali, W., Agrawal, R., Gutierrez-Sinisterra, L., Baxter, S. L., Cifuentes-González, C., Wei, Y. C., Abisheganaden, J., Kannapiran, P., Wong, S., Lee, B., de-la-Torre, A., & Agrawal, R. (2024). Natural Language Processing in Medicine and ophthalmology: A review for the 21st-century clinician. *Asia-Pacific Journal of Ophthalmology*, 13(4), 100084. <https://doi.org/10.1016/j.apjo.2024.100084>

## **Extracting Medical Information From Free-Text and Unstructured Patient-Generated Health Data Using Natural Language Processing Methods: Feasibility Study With Real-world Data**

### **Introduction:**

In recent years, Patient-Generated Health Data (PGHD) health-related data collected directly from patients or caregivers through devices such as smartphones, wearables, and home-based health monitoring systems. These data are often recorded in free-text form. Due to this the unstructured data remains underutilized in clinical workflows due to the lack of structured formats. This paper deals about the critical gap by proposing the Natural Language Processing pipeline to extract meaningful clinical information like symptoms, and medication etc., This study explores a zero-shot learning approach which eliminates the need of labeled training data for low-resource settings and applications where proper datasets are not available.

Another research that was followed in this article is analysis of data collected from 24 parents of children with special health care needs (CSHCN) who used a voice-enabled app over two weeks to record their daily experiences, symptoms, and medication routines. The voice data were transcribed using Amazon Web Services and processed through an NLP pipeline that incorporated SciSpaCy for biomedical NER, RXNorm and SNOMED CT for entity linking, and dependency parsing for relation extraction. The pipeline demonstrated the promising performance in

identifying the key health components with high precision and recall for symptoms mentions etc., The paper contributes a feasibility-focused exploration of how NLP can bridge the gap between patient-side health data and clinical decision-making, through automation and extraction techniques.

### **What work had been pursued:**

Unlike traditional NLP research, this study shifts its focus to free-text notes created by patients and caregivers, particularly parents of children with special health care needs (CSHCN). These notes, recorded via voice or text on a mobile app, define the insights about the symptoms, medication and the relevant data. A hybrid NLP pipeline with zero-shot capabilities model was implemented, defining it does not rely on domain-specific annotated training data. The process follows collecting real-world PGHD through a voice-enabled app from 24 caregivers over a two-week period then transcribing the voice data using AWS Transcribe into text notes. Then NLP pipeline was developed that combines Named Entity Recognition (NER), medical ontology mapping (using RXNorm for medications and SNOMED CT for symptoms), and dependency parsing for relational extraction. SciSpaCy (a domain-specific NLP model suite) and the date parser Python library to enhance syntactic and temporal information extraction from patient notes. Evaluating the pipeline's performance in terms of precision, recall, and F1 score across different components like medication names, units, quantities, dates, and symptoms. There were a few challenges like the use of Informal language use, unstructured variable sentences, the presence of noise and errors, and the different sizes of the dataset.

The automated extraction of medication and symptom information from patient-authored content is not only feasible but effective, even without domain-specific training data. Future systems that integrate PGHD into formal healthcare systems, enable better remote monitoring, personalized care, and real-time clinical decision support.

### **Tools and Techniques:**

Named Entity Recognition (NER), the NLP pipeline relied on NER to identify mentions of medications and symptoms in patient notes. This task was performed using SciSpaCy, which is optimized for processing scientific and clinical text. Dependency Parsing, Sentence-level dependency parsing to extract relationships between entities. Useful for identifying dosages, units, and temporal attributes. Zero-Shot Learning Approach tells that The pipeline did not require any training data, instead operating in a zero-shot setting. Ideal for low-resource environments where large, labeled datasets are unavailable. Patient entries made via voice were automatically transcribed using Amazon Web Services (AWS) Transcribe. The transcriptions were then processed through the NLP pipeline, supporting both audio and text inputs. The performance of the pipeline was evaluated by comparing output against ground-truth interpretations. And the metrics that are used are Precision, Recall, and F1 Score.

### **Reported Performance Measures and Findings:**

Performance of NLP pipeline in extracting medication and symptom-related information from real-world, patient-generated health data (PGHD). Precision, recall, and F1 score standard metrics in NLP tasks that reflect accuracy, and overall effectiveness. The pipeline focus is on Medication instances, Symptoms. Medication Instance: Precision – 0.83, Recall – 0.77, F1 Score- 0.80.



Symptoms: Precision –0.65, Recall- 0.82, F1 Score-0.72. The performance results indicate that the NLP pipeline is robust and reliable for medication name and symptom extraction, even when applied to noisy, real-world PGHD in a zero-shot learning context.

### **Conclusion:**

Right combination of rule-based and pre trained tools like SciSpaCy, RXNorm, and SNOMED CT. NLP pipelines can handle noisy, unstructured, and variable data authored by non-professionals. This expands the applicability of NLP beyond traditional EHRs. The zero-shot pipeline is deployable in scenarios where data labeling is not appropriate such as small populations etc. providing cost-effective and scalable solutions. PGHD when processed correctly, can complement traditional clinical data, providing a view of patient health, for chronic conditions or at-home care scenarios. Even with minimal resources and no prior data labeling, NLP can be effectively used to process messy, real-world health data. As healthcare becomes increasingly digital and decentralized, such systems will play a vital role in personalized, proactive, and patient-centric care.

### **Reference:**

Sezgin, E., Hussain, S.-A., Rust, S., & Huang, Y. (2023). Extracting medical information from free-text and unstructured patient-generated health data using natural language processing methods: Feasibility study with real-world data. *JMIR Formative Research*, 7, e43014. <https://doi.org/10.2196/43014>

## **Project Assignment 3 (10 points)**

Provide the following details related to your project

- Dataset: source, size, features description, etc.

Students will be provided with some examples of projects completed in the past.

The dataset used in this project is extracted from Kaggle, specifically from a collection of medical transcription records. These records represent real-world clinical documentation transcribed from interactions between healthcare professionals and patients. Each record captures essential details of a medical encounter, including patient complaints, medical history, surgical procedures, physical examinations, diagnostic assessments, and treatment plans. The transcriptions span

multiple medical specialties, offering a diverse and comprehensive view of clinical language and structure. It consists of 4,998 rows and 6 columns with each row corresponding to an individual medical transcription record, representing a unique patient encounter documented by a healthcare professional.

The dataset is structured in a tabular format and includes several key fields.

- **Description:** The description column provides a brief overview of the patient's chief complaint or the purpose of the visit.
- **Medical\_specialty:** The medical\_specialty field identifies the area of medicine relevant to the transcription, such as Bariatrics, Allergy / Immunology, or Cardiovascular / Pulmonary.
- **Sample\_name:** The sample\_name gives a short, descriptive title for each record, often indicating the type of consultation or procedure.
- **Transcription:** The transcription column is the most substantial feature, containing the full body of the transcribed medical note. This text may include structured sections like subjective/objective notes, physical exams, assessments, plans, medication history, and social or family histories.
- **keywords:** The keywords column lists important medical terms extracted from the transcription, which may include diagnoses, treatments, anatomical terms, and medications.

With nearly 5,000 records, the dataset is sufficiently large to support meaningful statistical analysis and machine learning model training along with Natural Language Processing (NLP) applications within the healthcare domain. It can be used to train classification models to categorize notes by medical specialty, perform named entity recognition (NER) to extract drugs, symptoms, or procedures, generate summaries of lengthy medical notes, or even build information retrieval systems for clinical documentation. Additionally, the dataset can serve as a training base for specialized language models in healthcare, helping to fine-tune models for tasks such as clinical predictions.

Overall, the dataset provides a comprehensive snapshot of clinical documentation practices across various specialties, making it an asset for researchers, data scientists, and healthcare technologists interested in developing intelligent systems to support medical documentation, decision support, and patient care analytics.

### Project Assignment 4 (10 points)

Provide the following details related to your project

- Deep Learning process
  - Choices of Deep Neural Networks and their justification
- Students will be provided with some examples of projects completed in the past.

The considered Dataset modality is **free-text natural language**, which is unstructured and domain-specific.

- **Transcription** gives the core input for most tasks and is composed of long-form medical narratives.
- **Medical specialty**, categorical label suitable for multi-class classification.
- **Description and Keywords** are both textual targets, useful for text summarization and sequence generation tasks respectively.

The Deep Learning Process covers NLP Tasks like Multi-class text classification, Sequence-to-sequence generation (summarization), Keyword extraction, or tagging. These models can cover automatic learning hierarchical representations of text, understanding contextual dependencies through mechanisms like attention or self-attention, Generalizing better from large unstructured text datasets without requiring heavy feature engineering.

Given the length, complexity, and domain-specific vocabulary in medical transcripts, traditional bag-of-words or TF-IDF-based models are often insufficient. Deep learning models, like recurrent, convolutional, or transformer-based architectures, are more suitable because they can capture long-range dependencies.

- **Recurrent Neural Networks (RNNs)**

**LSTM (Long Short-Term Memory):** It learns sequential dependencies in text data and remembers long-term information using gated memory cells. It can classify medical\_specialty based on the full transcription text can manage long input sequences and handle the sequential nature of language.

**BiLSTM (Bidirectional Long Short-Term Memory):** Medical transcriptions are sequential in nature, and follow a pattern like symptoms, diagnosis, and treatment. It can read the sequence both forward and backward, capturing dependencies from past and future tokens. Suitable for multi-class classification (medical-specialty) and NER tasks (e.g., extracting symptoms or procedures).

- **TextCNN (Convolutional Neural Network for Text)**

It can capture local patterns and n-gram features effectively (for example: Heart Failure). It is less computationally intensive than RNNs and very effective for classification tasks and can work well even with limited training data and is a strong baseline for classifying medical\_specialty.

- **Transformer-based Models**

**BioBERT / ClinicalBERT(Bidirectional Encoder Representations from Transformers)**

These models are trained on biomedical corpora (PubMed, MIMIC-III), making them adept at understanding medical language, abbreviations, and clinical terminology. It is ideal for classifying medical specialties, summarizing medical notes, and extracting keywords or conditions.

**T5 (Text-to-Text Transfer Transformer)**

A flexible, encoder-decoder architecture that frames every NLP task as a “text-to-text” problem. Suitable for summarization and keyword generation. Handles varying sequence lengths and complex relationships within the text.

- **Longformer / BigBird**

It overcomes BERT’s 512-token limit and is built for very long documents, which many medical transcriptions are and maintains the global and local aspects, allowing the model to learn structure across lengthy reports. Useful for tasks involving full transcription input.

We can use some more models like

- **Sequence Tagging with CRF + BERT/BiLSTM**

Combining contextual word embeddings (from BERT or BiLSTM) with a Conditional Random Field (CRF) output layer for sequence labelling. It is especially effective for keyword extraction, NER, or segmenting medical sections and tag each token in transcription with BIO Labels.

- **Multi-Task Learning (MTL) with Shared BERT or LSTM Backbone**

Medical data often involves interrelated tasks, and sharing a backbone model allows better generalization and resource efficiency. Regularizes training, improving robustness in small or imbalanced classes.

### **Project Assignment 5 (10 points)**

Provide the following details related to your project

- Deep Learning process
  - Training/Testing strategy

- Evaluation metrics
- Students will be provided with some examples of projects completed in the past.

We want to automatically classify medical transcriptions into corresponding medical specialties using advanced Natural Language Processing (NLP) techniques. Due to the unstructured nature of clinical descriptions and the varied length of transcriptions, there are some specifications again of choosing the deep learning models.

Finally, we are proceeding with 3 models and want to define our accuracy, as some models are not compatible with our data and objectives. Suitable for domain-specific classification because of biomedical pretraining.

### **DistilBERT (Hugging Face Transformers)**

DistilBERT is a smaller, faster version of BERT that retains 97% of its language understanding performance while being more efficient in computation.

- **Model:** distilbert-base-uncased
- **Tokenizer:** DistilBertTokenizerFast with truncation and padding
- **Input:** Tokenized text (max length = 512 tokens)

#### **Architecture:**

- DistilBERT transformer encoder layers (6 layers vs. BERT's 12)
- Fully-connected classification head with softmax activation

### **Framework: Hugging Face Transformers + TensorFlow**

- **Batch Size:** 16
- **Loss:** SparseCategoricalCrossentropy (from logits = True)
- **Optimizer:** Adam with learning rate 2e-5 (using Hugging Face create\_optimizer with weight decay)
- **Training Strategy:** Token-level embeddings averaged through attention for sentence classification

### **Bio\_ClinicalBERT (Hugging Face Transformers)**

Bio\_ClinicalBERT is a domain-specific BERT model pretrained on **clinical notes** from the MIMIC-III dataset, in addition to biomedical literature. It's designed for **clinical text** understanding and outperforms general BERT on medical tasks.

- **Model:** emilyalsentzer/Bio\_ClinicalBERT
- **Tokenizer:** AutoTokenizer (cased model) with truncation and padding

- **Input:** Tokenized clinical text (max length = 512 tokens)

#### **Architecture:**

- Bio\_ClinicalBERT transformer encoder layers (12 layers, similar to BERT-Base)
- Fully-connected classification head with softmax activation for multi-class output

#### **Framework: Hugging Face Transformers + TensorFlow**

- **Batch Size:** 8 or 16 (depending on GPU memory)
- **Loss:** SparseCategoricalCrossentropy (from\_logits=True)
- **Optimizer:** Adam (learning rate = 2e-5), with optional weight decay
- **Training Strategy:** Fine-tuned on clinical narratives, capturing medical terminology and relationships specific to healthcare.

#### **BiLSTM (Bidirectional Long Short-Term Memory)**

BiLSTM is a sequential deep learning model designed to capture both past (left) and future (right) context in text. It is effective for tasks where word order and dependencies matter, such as medical text classification.

#### **Model Type: RNN-based BiLSTM (Bidirectional LSTM Network)**

- **Tokenizer:** Keras Tokenizer (word-level) with padding/truncation.
- **Input:** Padded sequences of tokenized text (max length = 512 tokens).

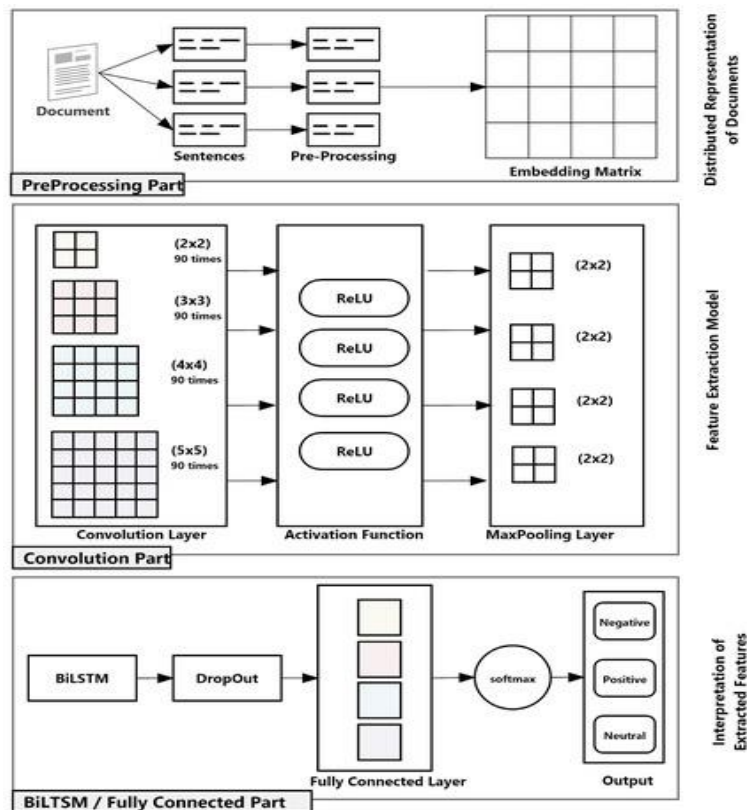
#### **Architecture:**

1. **Embedding Layer:**
  - Input Dimension: Vocabulary size (e.g., 20,000).
  - Output Dimension: 128 or pre-trained embeddings (e.g., GloVe).
  - Input Length: 512 tokens.
2. **BiLSTM Layer:**
  - Units: 64 (forward + backward)
  - Return Sequences: False
  - Captures context from both directions.
3. **Dropout Layer:** 0.3 (to prevent overfitting)
4. **Dense Layer:**
  - Units: 64
  - Activation: ReLU
5. **Output Layer:**
  - Units: Number of classes
  - Activation: Softmax (for multi-class classification)

#### **Framework: TensorFlow + Keras**

- **Batch Size:** 32

- Loss: CategoricalCrossentropy (one-hot encoded labels) or SparseCategoricalCrossentropy
- Optimizer: Adam (learning rate = 0.001)
- Training Epochs: Typically 5-10
- Class Weights: Used to handle class imbalance.



### Training/Testing Strategy:

We are considering an 80/20 stratified train-test split to maintain class balance for all models. Tokenized using either Hugging Face tokenizer or Keras tokenizer, depending on model type. The Handling Class Imbalance is done using `class_weight` in Keras models and ensures that rare classes are not ignored during training. Model training will be using validation accuracy. Following the Hyper Parameters: Epochs, Max Sequence, Batch Size, Optimizer, Learning Rate.

### Evaluation Metrics:

To evaluate the performance of the deep learning models applied to medical transcription classification, several standard metrics were used. We are using the primary baseline metric “Accuracy,” but accuracy cannot stand alone to reflect the model's effectiveness across all categories. Additional metrics such as **precision**, **recall**, and **F1-score** were employed for a more balanced outcome.

**Accuracy:** Overall correctness of predictions.

**Precision:** How many predicted specialties were correct?

**Recall:** How many actual specialties were successfully retrieved?

**F1 Score:** Harmonic mean of precision and recall.

**Macro F1:** F1 score averaged equally across classes.

**Weighted F1:** F1 score averaged by class frequency.

**Confusion Matrix:** Visualizes prediction vs. ground truth across all specialties.

### References:

Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2019). BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4), 1234–1240. <https://arxiv.org/abs/1901.08746>

Nguyen, V., Nguyễn, A., & Yang, H.-J. (2019). Real-time event detection using recurrent neural network in social sensors. *International Journal of Distributed Sensor Networks*, 15(6), 155014771985649. [https://www.researchgate.net/publication/333752473\\_Real-time\\_event\\_detection\\_using\\_recurrent\\_neural\\_network\\_in\\_social\\_sensors](https://www.researchgate.net/publication/333752473_Real-time_event_detection_using_recurrent_neural_network_in_social_sensors)

Rhanoui, M., Mikram, M., Yousfi, S., & Barzali, S. (2023). A CNN-BiLSTM model for document-level sentiment analysis. *Information*, 14(1), 50. <https://www.mdpi.com/2504-4990/1/3/48>

## Project Assignment 6 (10 points)

Provide the following details related to your project

- Obtained results
- References/Appendix
- Conclusion of your study
- Appendix should include the codebase and instructions to run your application.
- Prepare a demo version of your application which will be presented in week 7

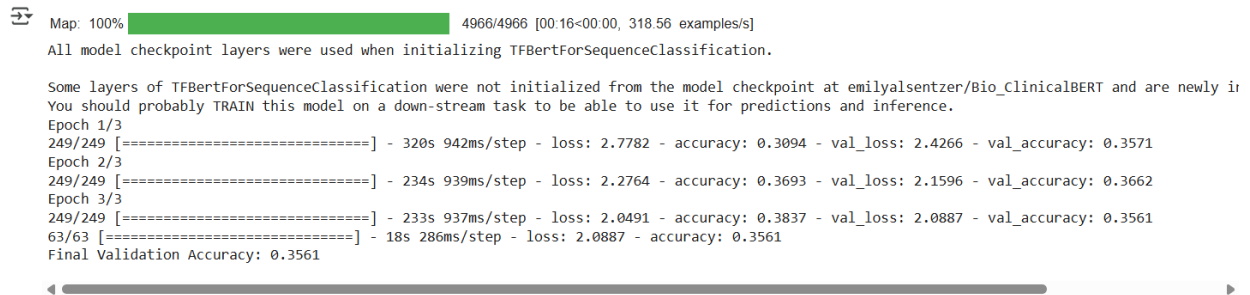
## Results

### Bio-ClinicalBERT

The Bio-ClinicalBERT model, specifically fine-tuned for medical text classification, achieved a final validation accuracy of 35.61% after 3 training epochs. The training process demonstrated consistent improvement across epochs, starting from 30.94% accuracy in the first epoch to 38.37% by the end of the third. The model's loss value decreased from 2.77 to 2.04 during this process, indicating effective learning. On unseen medical text data, the model successfully classified clinical scenarios into the appropriate specialties.



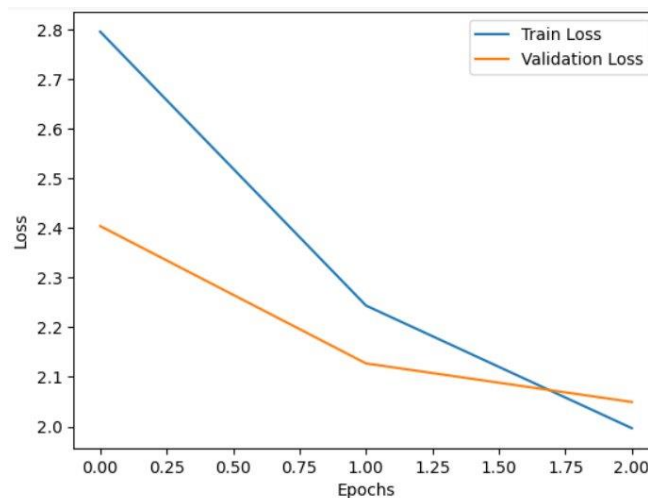
For instance, a text describing chest pain and myocardial infarction was correctly predicted as Cardiovascular / Pulmonary, and a case involving knee replacement was accurately labeled as Surgery. However, some misclassifications were observed, such as predicting Cardiovascular / Pulmonary for a brain tumor case instead of Neurosurgery. These results reflect Bio-ClinicalBERT's strong performance in capturing domain-specific medical contexts, with room for further improvement through extended fine-tuning or data augmentation.



1/1 [=====] - 0s 88ms/step  
 Text: The patient was admitted for chest pain and diagnosed with myocardial infarction.  
 Predicted Specialty: Cardiovascular / Pulmonary

Text: The patient underwent a total knee replacement and showed good post-operative recovery.  
 Predicted Specialty: Surgery

Text: MRI results indicate a brain tumor in the left frontal lobe, scheduled for neurosurgery.  
 Predicted Specialty: Cardiovascular / Pulmonary



## DistilBert

The DistilBERT model, fine-tuned on a dataset of medical transcriptions, demonstrated promising results for the task of multi-class classification across various medical specialties. After training for 5 epochs, the model achieved a final validation accuracy of 36.82%, with notable class-wise variations. The macro average F1-score stood at 0.08, indicating challenges with class imbalance and difficulty in accurately classifying less frequent medical specialties. The weighted average F1-score was 0.27, reflecting the model's stronger performance on more prevalent classes like "Surgery" (precision: 0.44, recall: 0.58)

In unseen text evaluation, the model successfully predicted relevant specialties, such as classifying chest pain and myocardial infarction under Cardiovascular / Pulmonary, and correctly assigning knee replacement recovery to Surgery.

Some weights of the PyTorch model were not used when initializing the TF 2.0 model TFDistilBertForSequenceClassification. This IS expected if you are initializing TFDistilBertForSequenceClassification from a PyTorch model trained on another task or with another data. This IS NOT expected if you are initializing TFDistilBertForSequenceClassification from a PyTorch model that you expect to be exactly identical. Some weights or buffers of the TF 2.0 model TFDistilBertForSequenceClassification were not initialized from the PyTorch model and are newly initialized. You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Epoch 1/5  
249/249 [=====] - 263s 980ms/step - loss: 2.9518 - accuracy: 0.3079 - val\_loss: 2.5408 - val\_accuracy: 0.3602

Epoch 2/5  
249/249 [=====] - 243s 974ms/step - loss: 2.5417 - accuracy: 0.3676 - val\_loss: 2.3465 - val\_accuracy: 0.3853

Epoch 3/5  
249/249 [=====] - 242s 972ms/step - loss: 2.3432 - accuracy: 0.3900 - val\_loss: 2.1994 - val\_accuracy: 0.3783

Epoch 4/5  
249/249 [=====] - 242s 974ms/step - loss: 2.1565 - accuracy: 0.4058 - val\_loss: 2.0815 - val\_accuracy: 0.3742

Epoch 5/5  
249/249 [=====] - 242s 972ms/step - loss: 2.0444 - accuracy: 0.4235 - val\_loss: 2.0461 - val\_accuracy: 0.3682

63/63 [=====] - 20s 299ms/step

	precision	recall	f1-score	support
Allergy / Immunology	0.00	0.00	0.00	1
Autopsy	0.00	0.00	0.00	2
Bariatrics	0.00	0.00	0.00	4
Cardiovascular / Pulmonary	0.32	0.22	0.26	74
Chiropractic	0.00	0.00	0.00	3
Consult - History and Phys.	0.31	0.79	0.45	103
Cosmetic / Plastic Surgery	0.00	0.00	0.00	5
Dentistry	0.00	0.00	0.00	5
Dermatology	0.00	0.00	0.00	6
Diets and Nutrition	0.00	0.00	0.00	2
Discharge Summary	0.66	0.86	0.75	22
ENT - Otolaryngology	0.00	0.00	0.00	19
Emergency Room Reports	0.00	0.00	0.00	15
Endocrinology	0.00	0.00	0.00	4
Gastroenterology	0.00	0.00	0.00	45
General Medicine	0.08	0.02	0.03	52
Hematology - Oncology	0.00	0.00	0.00	18
Hospice - Palliative Care	0.00	0.00	0.00	1
IME-QME-Work Comp etc.	0.00	0.00	0.00	3
Lab Medicine - Pathology	0.00	0.00	0.00	1
Letters	0.25	0.20	0.22	5
Nephrology	0.00	0.00	0.00	16
Neurology	0.22	0.13	0.17	45
Neurosurgery	0.00	0.00	0.00	19
Obstetrics / Gynecology	0.00	0.00	0.00	31
Office Notes	0.00	0.00	0.00	10
Ophthalmology	0.00	0.00	0.00	17
Orthopedic	0.23	0.14	0.17	71
Pain Management	0.00	0.00	0.00	12
Pediatrics - Neonatal	0.00	0.00	0.00	14
Physical Medicine - Rehab	0.00	0.00	0.00	4
Podiatry	0.00	0.00	0.00	9
Psychiatry / Psychology	0.00	0.00	0.00	11
Radiology	0.31	0.53	0.39	55
Rheumatology	0.00	0.00	0.00	2
SOAP / Chart / Progress Notes	0.34	0.39	0.37	33
Sleep Medicine	0.00	0.00	0.00	4
Speech - Language	0.00	0.00	0.00	2
Surgery	0.44	0.87	0.58	218
Urology	0.00	0.00	0.00	31
accuracy			0.37	994
macro avg	0.08	0.10	0.08	994
weighted avg	0.23	0.37	0.27	994

```

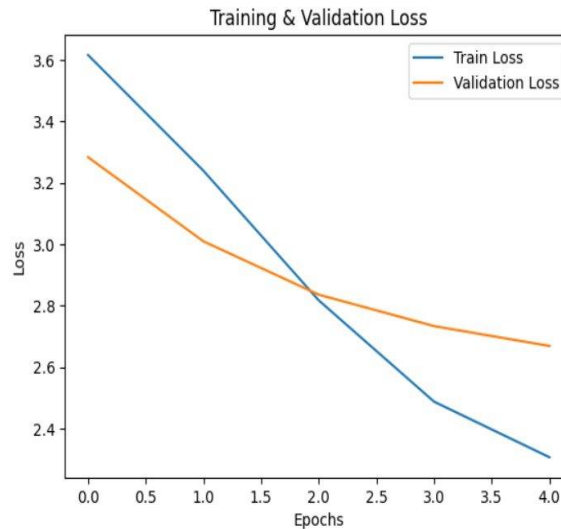
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior
_warn_prf(average, modifier, f'{metric.capitalize()} is', len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior
_warn_prf(average, modifier, f'{metric.capitalize()} is', len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior
_warn_prf(average, modifier, f'{metric.capitalize()} is', len(result))
1/1 [=====] - 2s 2s/step

Text: The patient presented with chest pain and was diagnosed with a heart attack.
Predicted Specialty: Cardiovascular / Pulmonary

Text: MRI scan showed a brain lesion consistent with glioblastoma, surgery scheduled.
Predicted Specialty: Radiology

Text: Patient shows signs of chronic kidney disease and requires dialysis.
Predicted Specialty: SOAP / Chart / Progress Notes

```



## BiLSTM :

```

Epoch 1/5
125/125 — 103s 781ms/step - accuracy: 0.1726 - loss: 3.3561 - val_accuracy: 0.3310 - val_loss: 2.6510 - learning_rate: 0.0010
Epoch 2/5
125/125 — 84s 676ms/step - accuracy: 0.3247 - loss: 2.7105 - val_accuracy: 0.3481 - val_loss: 2.4947 - learning_rate: 0.0010
Epoch 3/5
125/125 — 142s 680ms/step - accuracy: 0.3398 - loss: 2.4676 - val_accuracy: 0.3431 - val_loss: 2.4052 - learning_rate: 0.0010
Epoch 4/5
125/125 — 148s 726ms/step - accuracy: 0.3681 - loss: 2.2418 - val_accuracy: 0.3441 - val_loss: 2.3359 - learning_rate: 0.0010
Epoch 5/5
125/125 — 142s 724ms/step - accuracy: 0.3648 - loss: 2.1148 - val_accuracy: 0.3169 - val_loss: 2.3417 - learning_rate: 0.0010

```

This defines the The BiLSTM model trained on the mtsamples dataset shows a steady improvement in training performance across the five epochs. The training accuracy increases from 17% to 36%, while the training loss decreases significantly from 3.35 to 2.11, indicating that the model is learning from the data. However, the validation accuracy fluctuates between 31% and 34%, with the validation loss showing only a slight reduction from 2.65 to 2.34. This suggests that while the model is fitting the training data, it struggles to generalize well on unseen data, hinting at possible overfitting or limitations in model capacity. Additionally, the learning rate remains constant at 0.001 throughout training.

```

32/32 ----- 4s 121ms/step - accuracy: 0.3131 - loss: 2.3968
Test Accuracy: 34.41%
32/32 ----- 6s 188ms/step
              precision    recall  f1-score   support

Allergy / Immunology      0.00      0.00      0.00         1
Autopsy                   0.00      0.00      0.00         1
Bariatrics                0.00      0.00      0.00         4
Cardiovascular / Pulmonary 0.18      0.08      0.11        74
Chiropractic              0.00      0.00      0.00         3
Consult - History and Phy. 0.27      0.87      0.41       103
Cosmetic / Plastic Surgery 0.00      0.00      0.00         5
Dentistry                 0.00      0.00      0.00         5
Dermatology               0.00      0.00      0.00         6
Diets and Nutritions      0.00      0.00      0.00         2
Discharge Summary        0.00      0.00      0.00        22
ENT - Otolaryngology      0.00      0.00      0.00        19
Emergency Room Reports    0.00      0.00      0.00        15
Endocrinology             0.00      0.00      0.00         4
Gastroenterology         0.18      0.04      0.07        45
General Medicine          0.12      0.02      0.03        52
Hematology - Oncology     0.00      0.00      0.00        18
Hospice - Palliative Care 0.00      0.00      0.00         1
IME-QME-Work Comp etc.   0.00      0.00      0.00         3
Lab Medicine - Pathology  0.00      0.00      0.00         2
Letters                   0.00      0.00      0.00         5
Nephrology                0.00      0.00      0.00        16
Neurology                 0.23      0.16      0.19        45
Neurosurgery              0.00      0.00      0.00        19
Obstetrics / Gynecology   0.00      0.00      0.00        31
Office Notes              0.00      0.00      0.00        10
Ophthalmology             0.00      0.00      0.00        17
Orthopedic                0.06      0.01      0.02        71
Pain Management           0.00      0.00      0.00        12
Pediatrics - Neonatal     0.00      0.00      0.00        14
Physical Medicine - Rehab 0.00      0.00      0.00         4
Podiatry                  0.00      0.00      0.00         9
Psychiatry / Psychology   0.00      0.00      0.00        11
Radiology                 0.31      0.69      0.43        55
Rheumatology              0.00      0.00      0.00         2
SOAP / Chart / Progress Notes 0.00      0.00      0.00        33
Sleep Medicine            0.00      0.00      0.00         4
Speech - Language         0.00      0.00      0.00         2
Surgery                   0.45      0.90      0.60       218
Urology                   0.00      0.00      0.00        31

accuracy                  0.34       994
macro avg                 0.05       994
weighted avg              0.19       994

```

The BiLSTM model achieved a test accuracy of 34.41% on the mtsamples dataset, which is consistent with the validation performance observed during training. However, the detailed classification report reveals that the model struggles significantly across most individual classes. Many specialties, such as Allergy/Immunology, Dentistry, Dermatology, and Endocrinology, have precision, recall, and F1-scores close to zero, indicating that the model is unable to correctly predict these categories. Only a few larger classes, like Speech - Language and Radiology, show relatively better precision and recall, but even here, the performance is modest. The overall macro average F1-score is extremely low at 0.05, and the weighted average F1-score stands at 0.22, confirming a heavy imbalance in prediction quality across different medical specialties.



	BiLSTM	DistilBERT	Bio_ClinicalBERT
Validation Accuracy	34.41%	36.82%	38.37%
Training Accuracy	38.37%	42.35%	31.31%
Loss	2.3968	2.044	2.0491

## Conclusion:

This study explored various deep learning approaches for medical text classification using the mtsamples dataset, focusing on distinguishing between multiple medical specialties. Among the models evaluated, DistilBERT and Bio\_ClinicalBERT, both transformer-based architectures, demonstrated notable potential in capturing complex semantic patterns inherent in clinical narratives. DistilBERT achieved a validation accuracy of 36.8%, while Bio\_ClinicalBERT showed comparable performance, both benefiting from their pretrained contextual embeddings and ability to process long-range dependencies within text.

The BiLSTM model, designed to capture sequential dependencies, exhibited consistent learning during training with a rise in accuracy from 17% to 36% and a steady decline in training loss. Though its validation accuracy stabilized between 31% and 34%, this reflects a reliable ability to recognize broader patterns in clinical data, especially given the diversity and length of medical transcriptions.

Across all models, key insights emerged regarding the importance of domain-specific pretraining (as seen in Bio\_ClinicalBERT) and the effectiveness of transformer-based architectures for nuanced text. The classification outputs highlight that models are adept at predicting common medical categories such as Cardiovascular / Pulmonary and Surgery, and they provide a solid foundation for further refinement.

Future improvements could explore enhanced data balancing strategies, domain-adapted embeddings, and additional fine-tuning to bolster generalization across less frequent medical

classes. The promising results underline the potential of AI-driven models in automating and supporting medical documentation and specialty categorization.

## **Appendix:**

### **Code**

```
!pip install -q transformers datasets

import pandas as pd

from transformers import AutoTokenizer, TFAutoModelForSequenceClassification,
DataCollatorWithPadding

from transformers import create_optimizer

from datasets import Dataset

import tensorflow as tf

from sklearn.preprocessing import LabelEncoder

df = pd.read_csv('mtnsamples.csv') # Adjust path

df = df.dropna(subset=['transcription', 'medical_specialty'])

df['text'] = df['transcription']

label_encoder = LabelEncoder()

df['label'] = label_encoder.fit_transform(df['medical_specialty'])

num_labels = df['label'].nunique()

model_name = "emilyalsentzer/Bio_ClinicalBERT"

tokenizer = AutoTokenizer.from_pretrained(model_name)

dataset_hf = Dataset.from_pandas(df[['text', 'label']])

def tokenize_function(example):

    return tokenizer(example['text'], truncation=True, padding='max_length', max_length=256)
```

```

dataset_hf = dataset_hf.map(tokenize_function, batched=True)

train_test_split = dataset_hf.train_test_split(test_size=0.2, seed=42)

train_dataset = train_test_split['train']

val_dataset = train_test_split['test']

data_collator = DataCollatorWithPadding(tokenizer=tokenizer, return_tensors="tf")

tf_train = train_dataset.to_tf_dataset(
    columns=["input_ids", "attention_mask"],
    label_cols="label",
    shuffle=True,
    batch_size=16,
    collate_fn=data_collator,
)

tf_val = val_dataset.to_tf_dataset(
    columns=["input_ids", "attention_mask"],
    label_cols="label",
    shuffle=False,
    batch_size=16,
    collate_fn=data_collator,
)

model = TFAutoModelForSequenceClassification.from_pretrained(model_name,
num_labels=num_labels)

steps_per_epoch = len(tf_train)

num_train_steps = steps_per_epoch * 3

optimizer, schedule = create_optimizer(

```

```

init_lr=2e-5,

num_warmup_steps=0,

num_train_steps=num_train_steps

)

model.compile(optimizer=optimizer,
loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True), metrics=['accuracy'])

history = model.fit(tf_train, validation_data=tf_val, epochs=3)

results = model.evaluate(tf_val)

print(f"Final Validation Accuracy: {results[1]:.4f}")

new_texts = [

    "The patient was admitted for chest pain and diagnosed with myocardial infarction.",

    "The patient underwent a total knee replacement and showed good post-operative recovery.",

    "MRI results indicate a brain tumor in the left frontal lobe, scheduled for neurosurgery.",

]

encoded = tokenizer(new_texts, padding=True, truncation=True, max_length=512,
return_tensors='tf')

preds = model.predict(encoded)

import numpy as np

predicted_class_indices = np.argmax(preds.logits, axis=1)

label_map = {i: label for i, label in enumerate(label_encoder.classes_)}

predicted_labels = [label_map[idx] for idx in predicted_class_indices]

for text, label in zip(new_texts, predicted_labels):

    print(f"Text: {text}\nPredicted Specialty: {label}\n")

```

### **DistilBert Code:**



```

!pip install transformers datasets scikit-learn --quiet

import pandas as pd

import numpy as np

from sklearn.preprocessing import LabelEncoder

from sklearn.model_selection import train_test_split

from sklearn.metrics import classification_report

from sklearn.utils.class_weight import compute_class_weight

from transformers import DistilBertTokenizerFast, TFDistilBertForSequenceClassification,
create_optimizer

import tensorflow as tf

df = pd.read_csv('metsamples.csv') # Upload your file in Colab

df = df.dropna(subset=['transcription'])

label_encoder = LabelEncoder()

df['label'] = label_encoder.fit_transform(df['medical_specialty'])

num_classes = len(label_encoder.classes_)

train_texts, val_texts, train_labels, val_labels = train_test_split(

    df['transcription'].tolist(),

    df['label'].tolist(),

    test_size=0.2,

    random_state=42,

    stratify=df['label']

)

tokenizer = DistilBertTokenizerFast.from_pretrained('distilbert-base-uncased')

train_encodings = tokenizer(train_texts, truncation=True, padding=True, max_length=512)

```

```

val_encodings = tokenizer(val_texts, truncation=True, padding=True, max_length=512)

train_dataset = tf.data.Dataset.from_tensor_slices((dict(train_encodings), train_labels)).batch(16)

val_dataset = tf.data.Dataset.from_tensor_slices((dict(val_encodings), val_labels)).batch(16)

class_weights = compute_class_weight(class_weight="balanced", classes=np.unique(df['label']),
y=df['label'])

class_weights_dict = {i: class_weights[i] for i in range(len(class_weights))}

model = TFDistilBertForSequenceClassification.from_pretrained('distilbert-base-uncased',
num_labels=num_classes)

num_train_steps = len(train_dataset) * 5 # for 5 epochs

optimizer, schedule = create_optimizer(init_lr=2e-5, num_warmup_steps=0,
num_train_steps=num_train_steps)

model.compile(optimizer=optimizer,

              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),

              metrics=['accuracy'])

model.fit(train_dataset,

          validation_data=val_dataset,

          epochs=5,

          class_weight=class_weights_dict)

preds = model.predict(val_dataset)

y_pred = tf.argmax(preds.logits, axis=1).numpy()

print(classification_report(val_labels, y_pred, target_names=label_encoder.classes_))

test_texts = [

    "The patient presented with chest pain and was diagnosed with a heart attack.",

    "MRI scan showed a brain lesion consistent with glioblastoma, surgery scheduled.",

```

"Patient shows signs of chronic kidney disease and requires dialysis."

]

```
test_encodings = tokenizer(test_texts, truncation=True, padding=True, return_tensors='tf')
```

```
outputs = model.predict(dict(test_encodings))
```

```
predicted_class_indices = tf.argmax(outputs.logits, axis=1).numpy()
```

```
predicted_labels = label_encoder.inverse_transform(predicted_class_indices)
```

```
for text, pred in zip(test_texts, predicted_labels):
```

```
    print(f"\nText: {text}\nPredicted Specialty: {pred}")
```

### **BiLSTM Code:**

```
import tensorflow as tf
```

```
from tensorflow.keras.models import Sequential
```

```
from tensorflow.keras.layers import Embedding, Bidirectional, LSTM, Dense, Dropout
```

```
# BiLSTM Model
```

```
model = Sequential([
```

```
    Embedding(input_dim=20000, output_dim=128, input_length=512),
```

```
    Bidirectional(LSTM(64, return_sequences=False)),
```

```
    Dropout(0.5),
```

```
    Dense(64, activation='relu'),
```

```
    Dropout(0.3),
```

```
    Dense(num_classes, activation='softmax')
```

```
])
```

```

model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])

model.summary()

from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau

# Callbacks for better performance

early_stopping = EarlyStopping(monitor='val_loss', patience=2, restore_best_weights=True)

lr_scheduler = ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=1)

# Train

history = model.fit(X_train, y_train,

                    validation_data=(X_test, y_test),

                    epochs=5,

                    batch_size=32,

                    callbacks=[early_stopping, lr_scheduler])

loss, accuracy = model.evaluate(X_test, y_test)

print(f"Test Accuracy: {accuracy * 100:.2f}%")

# Predict

y_pred = model.predict(X_test)

y_pred_classes = np.argmax(y_pred, axis=1)

y_true = np.argmax(y_test, axis=1)

```

```
from sklearn.metrics import classification_report
```

```
print(classification_report(y_true, y_pred_classes, target_names=label_encoder.classes_))
```

## References:

1. Shin, D., Moon, J., & Hwang, H. (2022). Identifying the creation and impact of new technologies in patent text using NLP and machine learning. *Discover Artificial Intelligence*, 2(1), 6. <https://doi.org/10.1007/s44163-022-00006-3>
2. Kaur, H., Tyagi, V., & Garg, R. (2021). NLP Implementation: Current State, Challenges, and Perspectives. 2021 6th International Conference on Signal Processing, Computing and Control (ISPCC), 383–388. IEEE. <https://doi.org/10.1109/ISPCC53510.2021.9609373>
3. Barcelona, V., Scharp, D., Moen, H., Davoudi, A., Idnay, B. R., Cato, K., & Topaz, M. (2023). Using natural language processing to identify stigmatizing language in labor and birth clinical notes. *Maternal and Child Health Journal*, 28(3), 578–586. <https://doi.org/10.1007/s10995-023-03857-4>
4. Sheikhalishahi, S., Miotto, R., Dudley, J. T., Lavelli, A., Rinaldi, F., & Osmani, V. (2019). Natural Language Processing of Clinical Notes on Chronic Diseases: Systematic Review. *JMIR Medical Informatics*, 7(2), e12239. <https://doi.org/10.2196/12239>
5. Rojas-Carabali, W., Agrawal, R., Gutierrez-Sinisterra, L., Baxter, S. L., Cifuentes-González, C., Wei, Y. C., Abisheganaden, J., Kannapiran, P., Wong, S., Lee, B., de-la-Torre, A., & Agrawal, R. (2024). Natural Language Processing in Medicine and ophthalmology: A review for the 21st-century clinician. *Asia-Pacific Journal of Ophthalmology*, 13(4), 100084. <https://doi.org/10.1016/j.apjo.2024.100084>
6. Sezgin, E., Hussain, S.-A., Rust, S., & Huang, Y. (2023). Extracting medical information from free-text and unstructured patient-generated health data using natural language processing methods: Feasibility study with real-world data. *JMIR Formative Research*, 7, e43014. <https://doi.org/10.2196/43014>

## Project Assignment 7 (10 + 5 = 15 points)

Submit the final report after implementing all my feedback. Make sure that all the sections in the report are complete. In addition to that, the team will have an opportunity to present their project outcomes in the classroom for about 10-15 minutes. More instructions about the presentation will be provided by the instructor over the email.