# Simplifying ICD-10 Coding: Evaluating Pre-trained Language Models and Traditional Methods for Clinical Notes Mapping

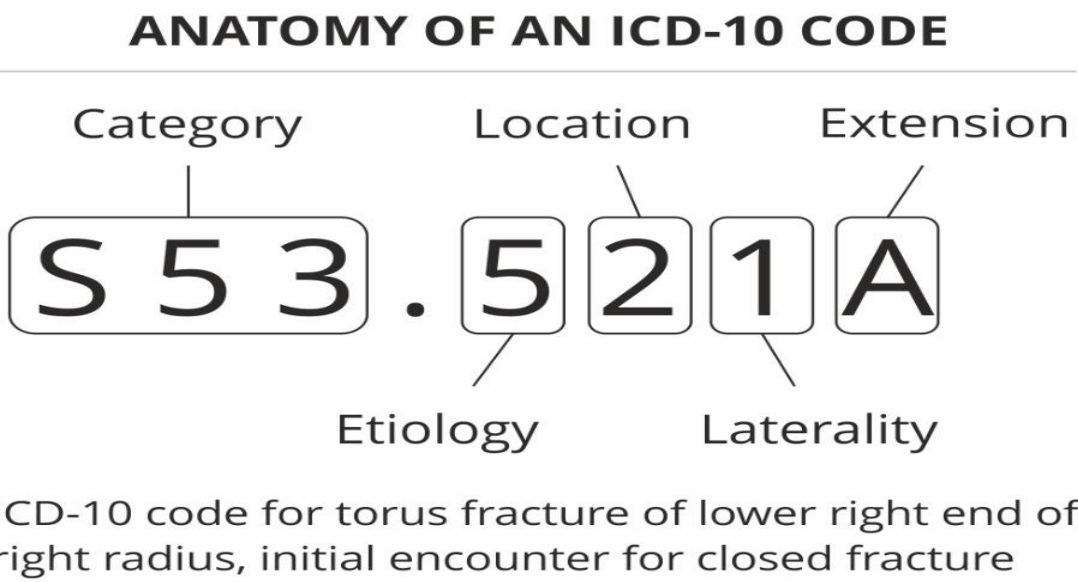## Syed Bilal Rizwan[1], Prof Dr. Sajjad Haider[1]*

[1]Department of Computer Science, Institute of Business Administration, Karachi, Pakistan

## BACKGROUND & MOTIVATION

The healthcare industry is dynamic and data-driven, evolving constantly to improve patient care, operational efficiency, and service quality. The usage of standardized coding systems dates to the 1800s, when the importance of such systems became clear due to nomenclature irregularities and the need for better statistical data.
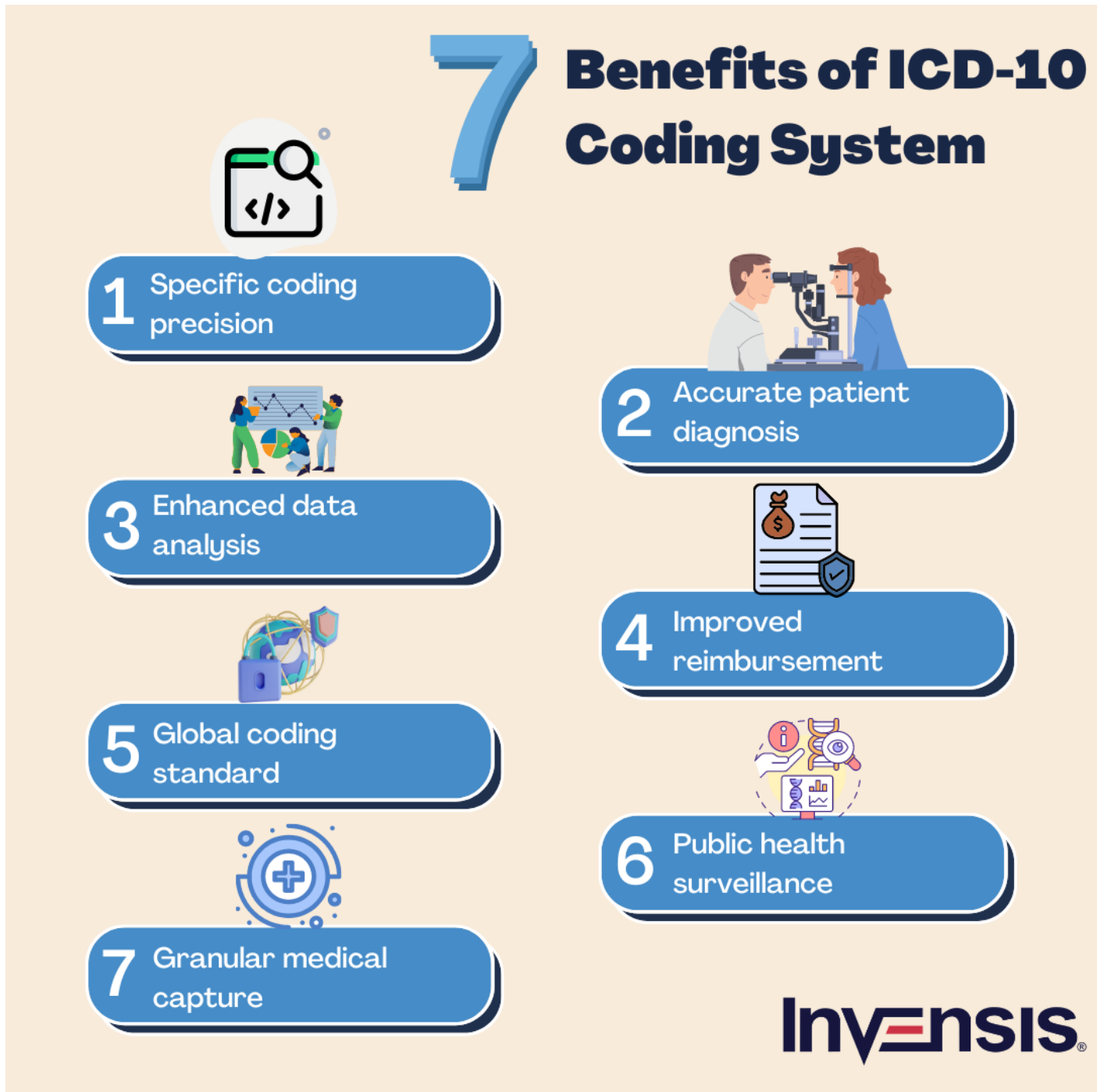
### What are ICD-10 Codes?

The International Classification of Diseases, Tenth Revision (ICD-10) is a worldwide recognized and standardized system of alphanumeric numbers used to identify diseases, disorders, and other health-related issues.
The first 3 digits represent the category of disease while the latter digits represents further specifications such as location, etiology



ICD-10 code for torus fracture of lower right end of right radius, initial encounter for closed fracture

### Challenges and Importance in Medical Coding

There are over 73,000 ICD-10 codes which makes manual coding by medical practitioners intricate and time-consuming. This requires healthcare institutes to hire additional experts to carry out the task.
However, despite the challenges, the benefits of medical coding compel healthcare institutes to carry out the tedious task.
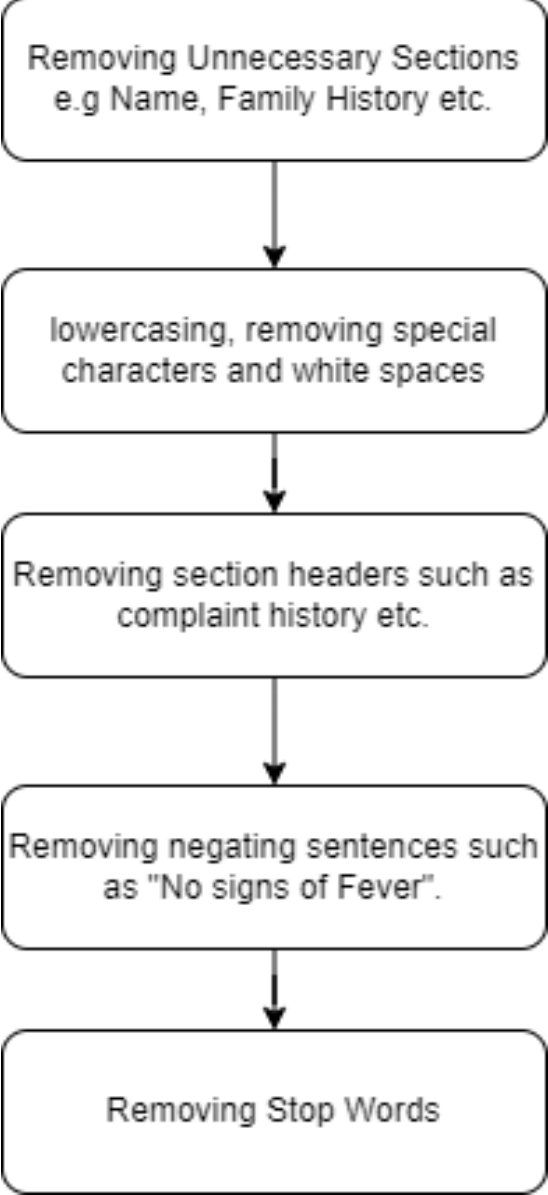


### Project Objective

This project aims to improve how clinical notes are mapped to ICD-10 codes by exploring pretrained language models and refining natural language processing techniques. An essential step involves creating a strong preprocessing technique for complex clinical notes. The main goal is to compare these new methods with existing ones like SNOMED CT with MEDCAT library, with the ultimate goal of finding an accurate and efficient existing model for clinical coding
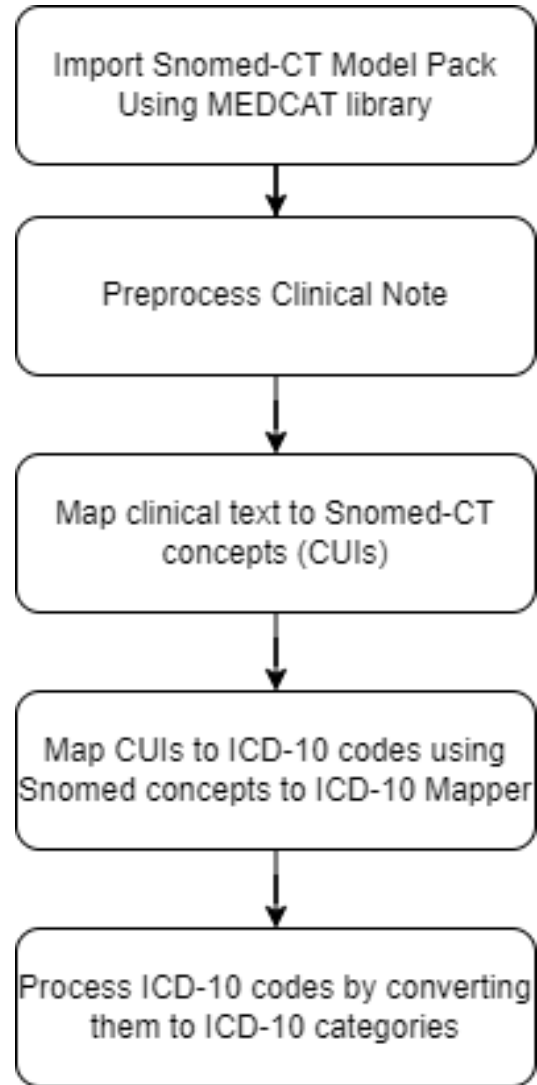
### Methodology

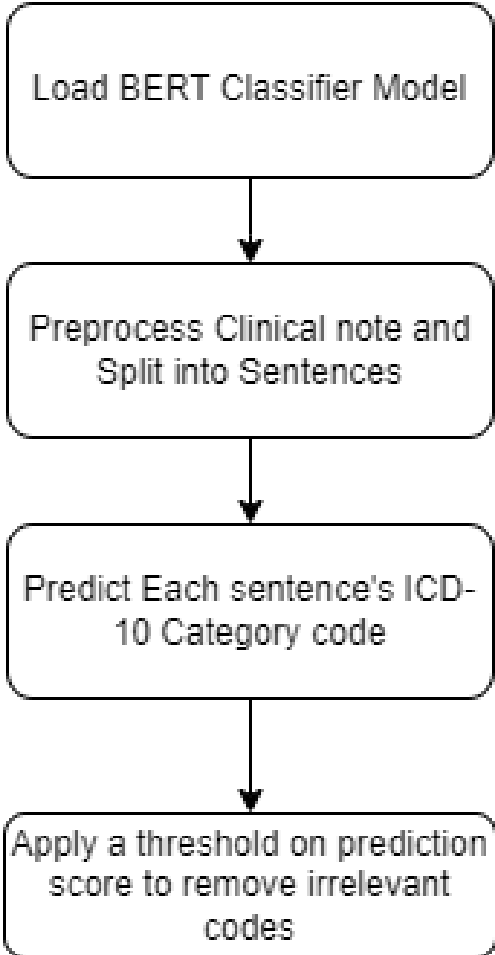There were 4 methodologies tried in this experiment with a common but extensive preprocessing technique:
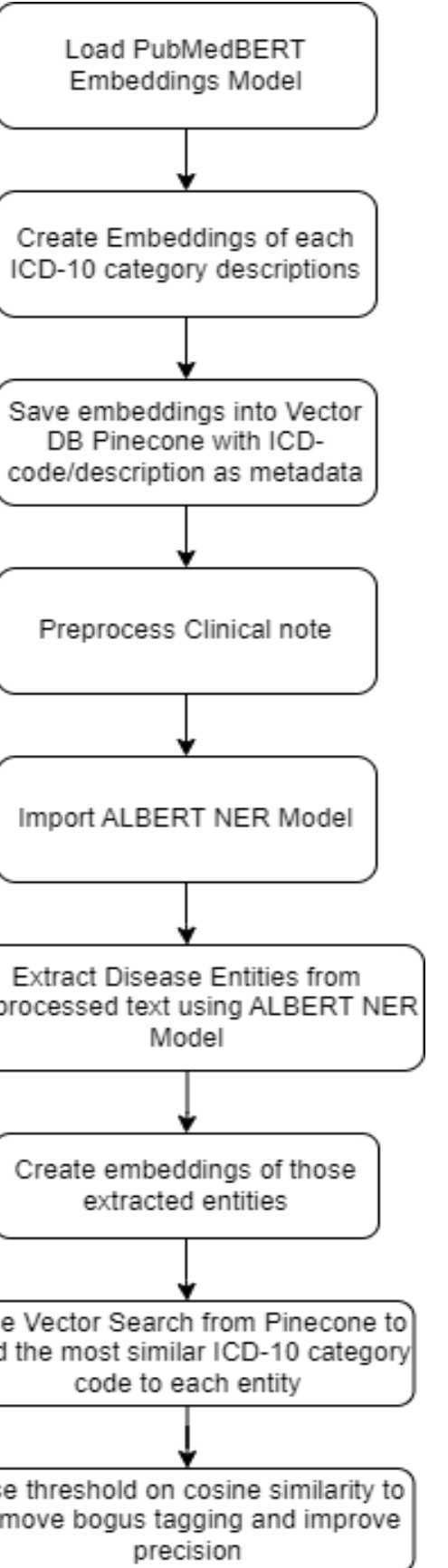
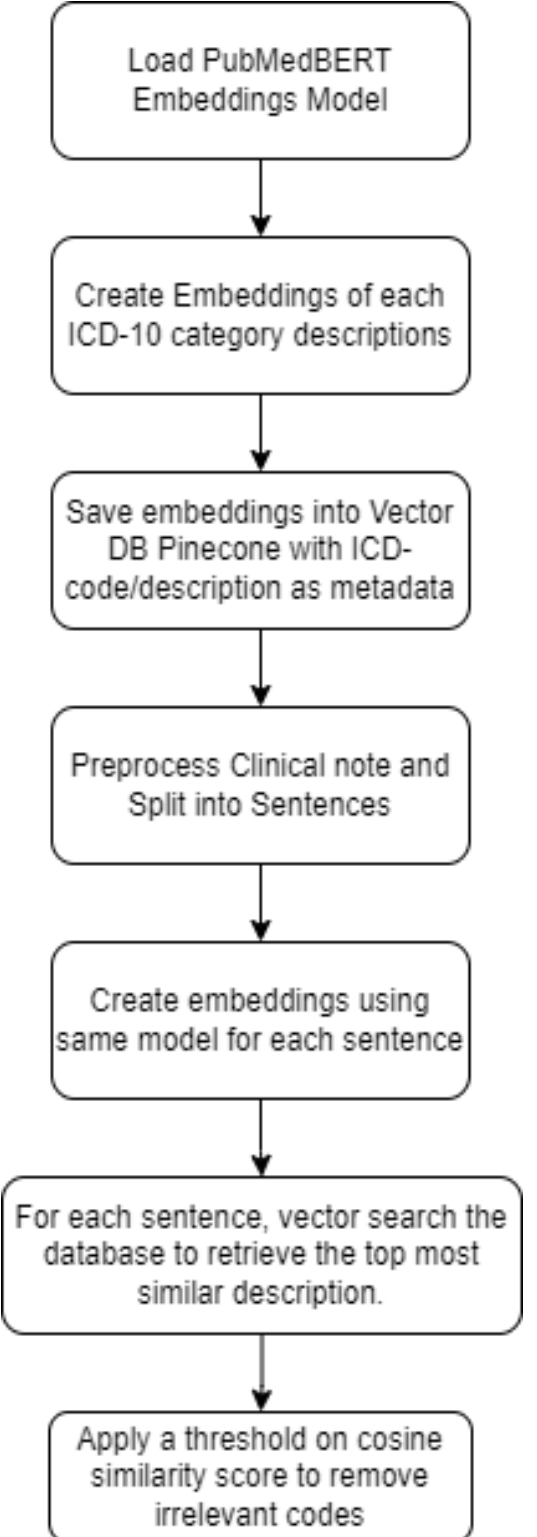#### Preprocessing Flow



#### SNOMED-CT (MEDCAT) Flow



#### BERT Classifier



#### PubMedBERT Embeddings Similarity

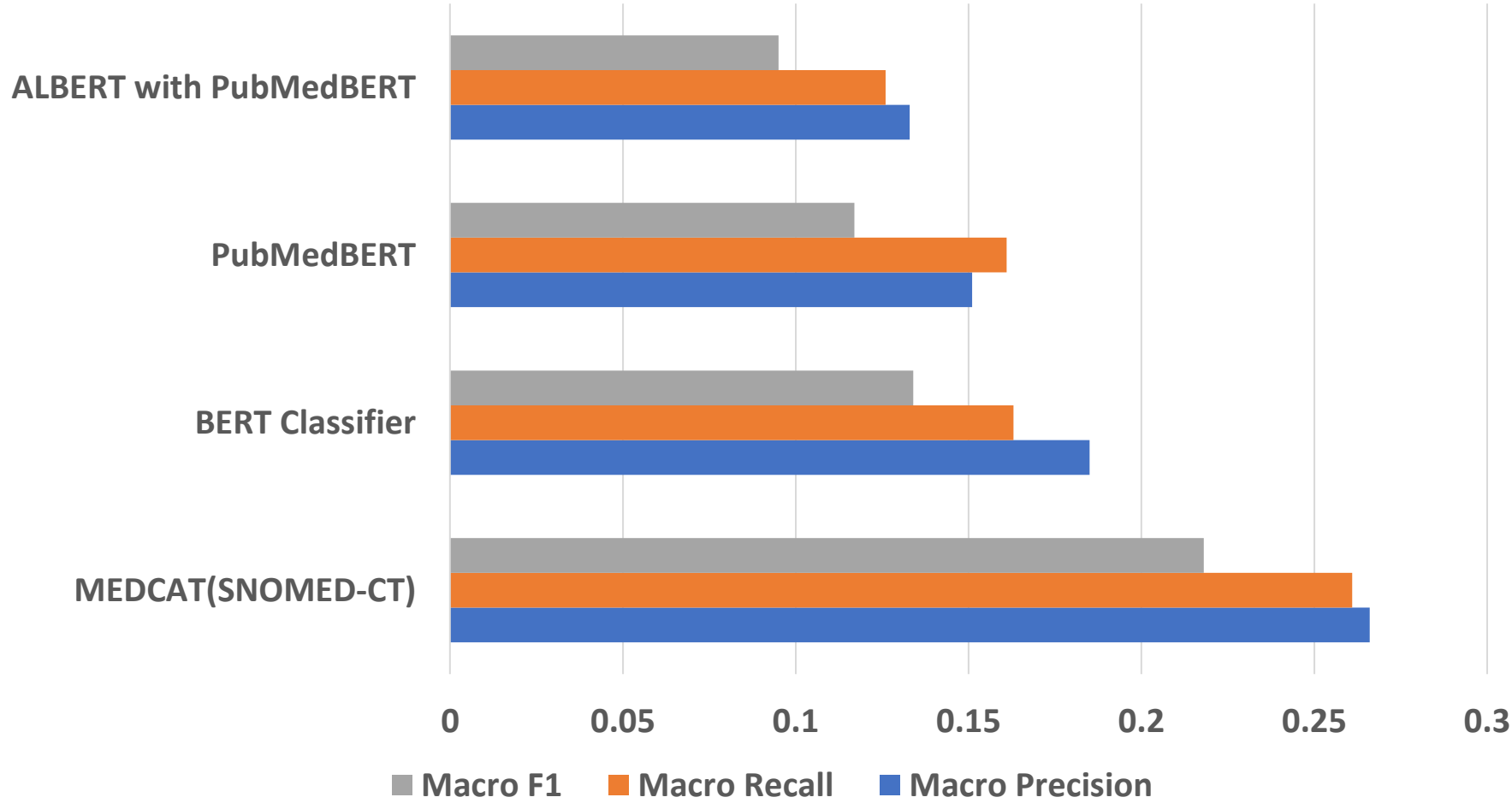

#### ALBERT with PubMedBERT



## RESULTS

- The dataset chosen for evaluation was the latest version of benchmark clinical notes dataset MIMIC IV.
- 1000 shortest notes were used from the dataset as a testing set.
- the first three digits (category) of the ICD-10 code were predicted for this project. The evaluation metrics chosen inspired by relevant research papers were Macro F1, Precision and Recall.
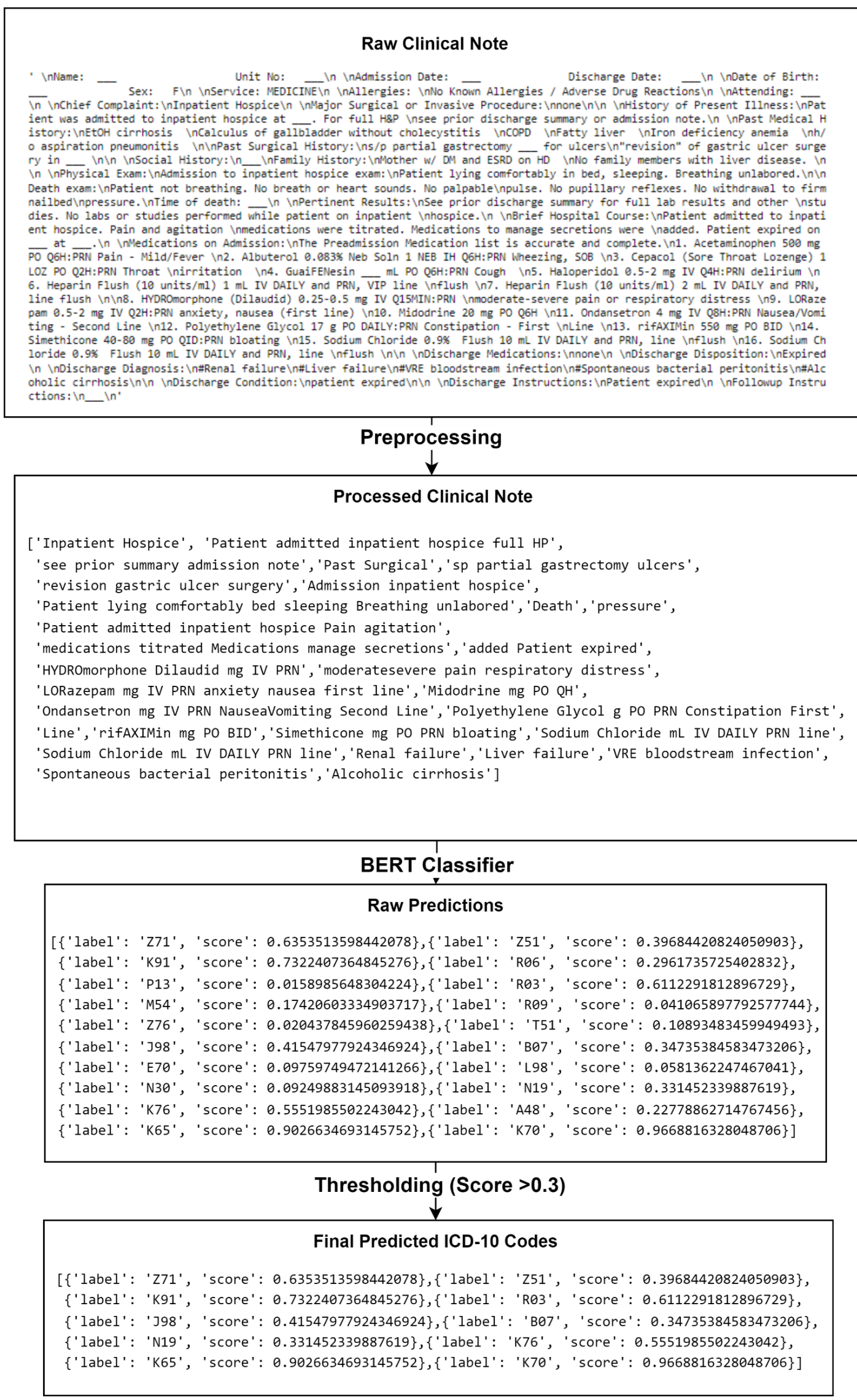- Total number of classes (ICD-10 Categories) in testing set were 769.

| S.no | Method | Macro Precision | Macro Recall | Macro F1 |
|---|---|---|---|---|
| 1 | MEDCAT(SNOMED-CT) | 0.266 | 0.261 | 0.218 |
| 2 | BERT Classifier | 0.185 | 0.163 | 0.134 |
| 3 | PubMedBERT | 0.151 | 0.161 | 0.117 |
| 4 | ALBERT with PubMedBERT | 0.133 | 0.126 | 0.095 |

*Table 1: Final Metrics*



## Walkthrough Example



## Discussion and Final Remarks

**Top Performer:** SNOMED-CT from MEDCAT library excelled a Macro F1-score of 0.218, establishing a significant project milestone.

**BERT Success:** Emran's BERT model, coupled with advanced preprocessing, outperformed other transformer-based techniques with a comparable Macro-F1 score of 0.134.

**Complexity Acknowledged:** Considering the intricate multi-label classification problem involving around 24,000 classes with 769 classes in testing set, achieving such scores reflects commendable progress in addressing the inherent complexity of the task.

**Comparable Results:** Related works carried out for this problem had very high metrics compared to this project however it was seen that in those research papers the size of a clinical note used were very short which is around 20-50 words while the MIMIC IV dataset had notes ranging between 2000-10000 words. Furthermore, the reference papers mostly focused on a classification of a single ICD-10 code while our methods tried to perform a multi-label classification which came with its complexities.

**Contributions to Healthcare:** The methodologies tested present tangible solutions to challenges in clinical note mapping, making a noteworthy contribution to the healthcare information systems landscape.

## References

Nguyen, Anthony N., Donna L. Truran, Madonna Kemp, Bevan Koopman, David Conlan, John O'Dwyer, Ming Zhang, Sarvnaz Karimi, Hamed Hassanzadeh, Michael Lawley and Damian J. Green. "Computer-Assisted Diagnostic Coding: Effectiveness of an NLP-based approach using SNOMED CT to ICD-10 mappings." AMIA ... Annual Symposium proceedings. AMIA Symposium 2018 (2018): 807-816.

Al-Bashabsheh, Emran, Ahmad Alaiad, Mahmoud Al-Ayyoub, Othman Beni-Yonis, Raed Abu Zitar, and Laith Abualigah. "Improving clinical documentation: automatic inference of ICD-10 codes from patient notes using BERT model." The Journal of Supercomputing (2023): 1-25.

Huang, Chao-Wei, Shang-Chi Tsai, and Yun-Nung Chen. "PLM-ICD: automatic ICD coding with pretrained language models." arXiv preprint arXiv:2207.05289 (2022).

"International Classification of Diseases." Encyclopædia Britannica. Accessed December 29, 2023. https://www.britannica.com/topic/International-Classification-of-Diseases.

Zeljko. "MedCAT: Introduction-Analyzing Electronic Health Records." Medium, April 23, 2020. https://medium.com/@w_is_h/medcat-introduction-analyzing-electronic-health-records-e1c420afa13a.

## Acknowledgement