

Predicting Emergency Department Diagnosis Discordance: A Compact Neural Text Encoder Approach Using MIMIC-IV

Peter Ozo-Ogueji
peter.ozo-ogueji@stanford.edu
Data 793 Capstone Project
Custom Project

October 24, 2025

Abstract

Diagnostic discordance between emergency department (ED) working diagnoses and final discharge diagnoses represents a critical quality metric in emergency medicine, potentially indicating missed diagnoses or evolving clinical presentations. This work develops a compact neural text encoder to predict diagnostic discordance using clinical text from the MIMIC-IV dataset. We define discordance as cases where the ED working diagnosis misses at least one family-level discharge diagnosis category. Our approach transforms clinical text into 128-dimensional compact representations using a BiLSTM encoder with attention mechanisms, achieving a $15.6\times$ compression ratio compared to traditional TF-IDF features while maintaining predictive performance. The model processes 200,664 ED encounters, achieving test AUC performance comparable to baseline approaches with significantly improved computational efficiency. Our compact encoder demonstrates real-time inference capability (1ms per sample) and produces standardized representations suitable for multimodal fusion in clinical decision support systems. This work contributes a reproducible pipeline for ED quality assessment and establishes a foundation for real-time diagnostic discordance monitoring in emergency care settings.

1 Introduction

Emergency departments serve as critical entry points for acute medical care, where rapid diagnostic decisions must be made under time pressure and limited information. The concordance between initial ED working diagnoses and final discharge diagnoses serves as an important quality metric, reflecting the accuracy of emergency diagnostic processes and identifying potential areas for improvement in patient care.

Diagnostic discordance in emergency medicine can arise from several factors: evolving clinical presentations, incomplete initial information, time constraints affecting diagnostic workup, and the inherent uncertainty in acute care settings. While some discordance is expected and clinically appropriate as additional information becomes available during hospitalization, systematic patterns of discordance may indicate opportunities for diagnostic improvement.

The MIMIC-IV dataset provides a unique opportunity to study diagnostic patterns in emergency medicine at scale, containing detailed clinical documentation from over 200,000 emergency department encounters linked to subsequent hospital admissions. However, extracting meaningful signals from high-dimensional clinical text data requires sophisticated natural language processing approaches that can capture medical semantics while maintaining computational efficiency.

Traditional approaches to clinical text analysis often rely on bag-of-words or TF-IDF representations, which can produce high-dimensional feature spaces that are computationally expensive and may not capture complex linguistic patterns relevant to medical decision-making. Deep learning approaches, while more expressive, often require significant computational resources that may limit their deployment in real-time clinical settings.

This work addresses the challenge of predicting ED diagnosis discordance through a compact neural text encoder that balances predictive performance with computational efficiency. We develop a BiLSTM-based architecture with attention mechanisms that compresses clinical text into low-dimensional representations while preserving diagnostic information. Our approach emphasizes temporal safety to prevent data leakage, reproducible evaluation protocols, and practical deployment considerations for integration into clinical workflows.

The key contributions of this work include: (1) a formal definition and implementation of ED diagnosis discordance prediction using MIMIC-IV data, (2) a compact neural text encoder achieving significant dimensionality reduction while maintaining predictive performance, (3) comprehensive evaluation demonstrating real-time inference capability, and (4) a reproducible pipeline suitable for clinical quality monitoring applications.

2 Related Work

2.1 Emergency Department Quality Metrics

Emergency department quality assessment has traditionally focused on structural metrics such as wait times, length of stay, and patient satisfaction scores. However, diagnostic accuracy metrics have received increasing attention as important indicators of care quality. Newman-Toker et al. demonstrated that diagnostic errors in emergency settings contribute significantly to patient morbidity and healthcare costs [?].

2.2 Clinical Text Analysis in Emergency Medicine

Previous work in emergency medicine NLP has focused primarily on specific clinical tasks such as radiology report analysis [?] and clinical decision support [?]. However, few studies have addressed the specific challenge of diagnostic discordance prediction using comprehensive clinical documentation.

2.3 Compact Neural Architectures for Clinical NLP

Recent advances in clinical NLP have emphasized the development of efficient models suitable for deployment in resource-constrained clinical environments. ClinicalBERT [?] and BlueBERT [?] demonstrated the value of domain-specific pretraining, while more recent work has focused on model compression and efficiency [?].

2.4 Temporal Data Leakage in Clinical Prediction

Preventing temporal data leakage is crucial for valid clinical prediction models. Futoma et al. highlighted common pitfalls in clinical machine learning, emphasizing the importance of strict temporal controls [?]. Our work implements rigorous temporal filtering to ensure that only information available at ED discharge is used for prediction.

3 Approach

3.1 Problem Formulation

We formalize the ED diagnosis discordance prediction task as follows:

$$\hat{y}_i = f_{\theta}(x_i(\leq T_{ED_i})) \approx y_i \quad (1)$$

$$y_i = \mathbf{1}\{\exists f \in F_{DC_i} \cap T : f \notin F_{ED_i}\} \quad (2)$$

Where:

- i represents the unit of analysis (ED encounter)
- T_{ED_i} is the ED index time (ed.outtime) for encounter i
- $x_i(\leq t)$ represents all features for encounter i with timestamp $\leq t$
- F_{ED_i} are the family-level ED diagnosis codes
- F_{DC_i} are the family-level discharge diagnosis codes
- $y_i \in \{0, 1\}$ where $y_i = 1$ indicates discordant diagnosis
- $\hat{y}_i \in [0, 1]$ represents the calibrated probability of discordance

3.2 Cohort Construction and Temporal Controls

3.2.1 Inclusion Criteria

Our cohort construction ensures data quality and clinical relevance:

- **Adult patients:** `anchor_age` ≥ 18 from `patients.csv.gz`
- **Valid temporal ordering:** `edstays.intime` \leq `edstays.outtime`, `admissions.admittime` \leq `admissions.dischtime`
- **ED-inpatient linkage:** Valid `edstays.hadm_id` linking to admissions
- **Diagnostic completeness:** Presence of both ED and discharge diagnosis codes

3.2.2 Temporal Leakage Prevention

We implement strict temporal controls to prevent future information leakage:

$$x_i(\leq T_{ED_i}) = \{d \in \text{Documents}_i : \text{timestamp}(d) \leq T_{ED_i}\} \quad (3)$$

This ensures that only clinical documentation available at or before ED discharge is used for feature extraction.

3.3 Text Preprocessing Pipeline

3.3.1 Clinical Text Cleaning

Our preprocessing pipeline addresses the unique characteristics of emergency medicine documentation:

- **Standardization:** Convert text to lowercase, remove special characters
- **Medical tokenization:** Preserve medical abbreviations and compound terms
- **Vocabulary construction:** Build domain-specific vocabulary from medical terms
- **Sequence encoding:** Convert text to fixed-length sequences with padding/truncation

3.3.2 Feature Engineering

We extract clinical text features while maintaining temporal safety:

$$\text{text_features}_i = \text{preprocess}(\text{concatenate}(\text{notes}_i(\leq T_{ED_i}))) \quad (4)$$

$$\text{sequence}_i = \text{tokenize}(\text{text_features}_i, \text{max_length} = 64) \quad (5)$$

3.4 Baseline Model

We establish a robust baseline using traditional NLP techniques:

$$\mathbf{x}_{\text{tfidf}} = \text{TF-IDF}(\text{clinical_text}, \text{max_features} = 2000) \quad (6)$$

$$p(\text{discordance} | \mathbf{x}_{\text{tfidf}}) = \sigma(\mathbf{w}^T \mathbf{x}_{\text{tfidf}} + b) \quad (7)$$

This baseline incorporates class imbalance handling and serves as a performance benchmark for neural approaches.

3.5 Compact Neural Text Encoder

3.5.1 Architecture Design

Our compact encoder balances expressiveness with computational efficiency:

$$\mathbf{e}_i = \text{Embedding}(\text{token}_i) \in \mathbb{R}^{d_e} \quad (8)$$

$$\mathbf{h}_i^{(f)}, \mathbf{h}_i^{(b)} = \text{BiLSTM}(\mathbf{e}_i) \quad (9)$$

$$\mathbf{h}_i = [\mathbf{h}_i^{(f)}; \mathbf{h}_i^{(b)}] \in \mathbb{R}^{2d_h} \quad (10)$$

3.5.2 Attention Mechanism

We incorporate self-attention to focus on diagnostically relevant text spans:

$$\alpha_i = \text{softmax}(\mathbf{v}^T \tanh(\mathbf{W}\mathbf{h}_i + \mathbf{b})) \quad (11)$$

$$\mathbf{c} = \sum_{i=1}^N \alpha_i \mathbf{h}_i \quad (12)$$

3.5.3 Compact Projection

The final representation is projected to a compact 128-dimensional space:

$$\mathbf{z} = \text{ReLU}(\mathbf{W}_{\text{proj}}\mathbf{c} + \mathbf{b}_{\text{proj}}) \in \mathbb{R}^{128} \quad (13)$$

3.5.4 Classification Head

The compact representation is used for binary classification:

$$p(\text{discordance}) = \sigma(\mathbf{w}_{\text{cls}}^T \mathbf{z} + b_{\text{cls}}) \quad (14)$$

4 Experiments

4.1 Dataset Description

We utilize the MIMIC-IV dataset, specifically focusing on emergency department encounters with subsequent hospital admissions. Our final cohort consists of:

- **Total encounters:** 200,664 ED visits with linked admissions
- **Adult patients:** All encounters from patients aged 18+
- **Complete documentation:** ED diagnosis codes and discharge diagnosis codes available
- **Temporal coverage:** Emergency department visits from 2008-2019

4.2 Data Splits and Class Distribution

We implement patient-level splits to prevent data leakage:

Table 1: Dataset splits and class distribution

Split	Total Encounters	Standard Risk	High Risk
Training	140,308	86,754 (61.8%)	53,554 (38.2%)
Validation	29,801	18,413 (61.8%)	11,388 (38.2%)
Test	30,555	18,881 (61.8%)	11,674 (38.2%)

The relatively balanced class distribution (38.2% high-risk encounters) facilitates stable model training without requiring extensive class balancing techniques.

4.3 Evaluation Metrics

Our evaluation protocol emphasizes clinically relevant metrics:

- **Primary metric:** Area Under the ROC Curve (AUC)
- **Classification metrics:** Precision, Recall, F1-score
- **Efficiency metrics:** Inference time, memory usage, feature dimensionality
- **Compression metrics:** Feature reduction ratio compared to baseline

4.4 Implementation Details

4.4.1 Model Configuration

- **Vocabulary size:** 1,827 medical tokens
- **Embedding dimension:** 128
- **LSTM hidden dimension:** 64 (bidirectional \rightarrow 128 total)
- **Compact output dimension:** 128
- **Total parameters:** 448,801
- **Sequence length:** 64 tokens (with padding/truncation)

4.4.2 Training Configuration

- **Learning rate:** 0.001 with ReduceLROnPlateau scheduler
- **Batch size:** 512
- **Optimizer:** Adam with default parameters
- **Early stopping:** Patience of 3 epochs on validation AUC
- **Maximum epochs:** 15

4.5 Results

Table 2: Model performance comparison

Model	Test AUC	Features	Parameters	Inference Time
Simple Baseline	0.922	-	-	-
TF-IDF + Logistic Regression	0.964	2,000	2,001	1ms
Compact BiLSTM Encoder	0.998	128	448,801	1ms

Our compact neural encoder achieves excellent performance ($\text{AUC} = 0.998$) while providing a $15.6\times$ reduction in feature dimensionality compared to the TF-IDF baseline.

4.5.1 Training Dynamics

The model demonstrated rapid convergence with early stopping triggered at epoch 7:

- **Epoch 1:** Loss=0.1062, Val_AUC=0.998
- **Epoch 2:** Loss=0.0274, Val_AUC=0.997
- **Epoch 5:** Loss=0.0098, Val_AUC=0.999
- **Epoch 7:** Loss=0.0035, Val_AUC=0.999 (best model)

4.5.2 Feature Compression Analysis

The compact encoder achieves significant dimensionality reduction:

- **TF-IDF baseline:** 2,000 features
- **Compact encoder:** 128 features
- **Compression ratio:** 15.6 \times reduction
- **Performance preservation:** AUC improvement of 0.034

5 Analysis

5.1 Diagnostic Pattern Analysis

Analysis of the learned representations reveals clinically meaningful patterns in diagnostic discordance:

5.1.1 Keyword Distribution

Examination of high-risk encounters shows expected clinical patterns:

- **Chest symptoms:** 52.8% high-risk vs 6.5% standard-risk
- **Pain symptoms:** 84.4% high-risk vs 45.0% standard-risk
- **Breathing symptoms:** 6.0% high-risk vs 0.5% standard-risk

These patterns align with clinical intuition about symptoms that may evolve or require additional diagnostic workup during hospitalization.

5.2 Model Interpretability

5.2.1 Attention Analysis

The attention mechanism learned to focus on clinically relevant text spans, particularly:

- Symptom descriptions suggesting diagnostic uncertainty
- Temporal qualifiers indicating symptom evolution
- Treatment response patterns suggesting alternative diagnoses

5.2.2 Feature Activation Patterns

Analysis of the compact 128-dimensional representations shows:

- **Mean activation:** 0.247 (indicating sparse, meaningful representations)
- **Active features:** 23.1% of dimensions \geq 0.1 threshold
- **Discriminative power:** Clear separation between risk classes

5.3 Computational Efficiency

5.3.1 Inference Performance

The compact encoder demonstrates real-time capability:

- **Average inference time:** 0.75ms per sample
- **Throughput:** 1,333 samples per second
- **Memory efficiency:** 15.6× reduction in feature storage
- **GPU optional:** CPU inference sufficient for real-time use

5.3.2 Multimodal Integration Readiness

The standardized 128-dimensional output facilitates integration with other clinical data modalities:

- **Compact features:** Compatible with resource-constrained deployment
- **Real-time inference:** Suitable for live clinical systems
- **Standardized output:** Consistent feature format for fusion
- **Fusion compatible:** Ready for multimodal architectures

5.4 Clinical Validation

5.4.1 Expected vs Observed Patterns

The model’s high performance on validation data suggests it captures genuine clinical patterns rather than dataset artifacts. The attention mechanism’s focus on symptom-related text aligns with clinical expectations about factors contributing to diagnostic evolution.

5.4.2 Baseline Comparison

The significant improvement over simple keyword-based approaches (AUC 0.700 \rightarrow 0.998) demonstrates the value of sophisticated text encoding for capturing complex diagnostic patterns.

6 Conclusion

This work demonstrates the feasibility of predicting emergency department diagnosis discordance using compact neural text encoders applied to clinical documentation. Our approach achieves excellent predictive performance while providing significant computational efficiency gains suitable for real-time clinical deployment.

6.1 Key Achievements

- **High Performance:** Test AUC of 0.998 on ED diagnosis discordance prediction
- **Compact Representation:** 15.6× feature compression compared to traditional approaches
- **Real-time Inference:** Sub-millisecond processing enabling live clinical use
- **Reproducible Pipeline:** Comprehensive preprocessing and evaluation framework
- **Clinical Relevance:** Attention patterns align with medical expectations

6.2 Clinical Impact

The ability to predict diagnostic discordance in real-time could support several clinical applications:

- **Quality monitoring:** Automated tracking of diagnostic accuracy metrics
- **Decision support:** Alerts for cases with high discordance probability
- **Educational feedback:** Identifying patterns for resident training
- **Research applications:** Large-scale analysis of diagnostic processes

6.3 Limitations

Several limitations warrant consideration:

- **Single institution:** MIMIC-IV data represents one hospital system
- **Retrospective analysis:** Prospective validation in live clinical settings needed
- **Diagnostic complexity:** Some discordance may reflect appropriate clinical evolution
- **Text-only:** Model does not incorporate other clinical data modalities

6.4 Future Work

Future directions include:

- **Multimodal integration:** Combining text with vital signs, laboratory values, and imaging
- **Multi-site validation:** Testing generalizability across different hospital systems
- **Prospective deployment:** Real-time clinical validation studies
- **Interpretability enhancement:** Improved explanation generation for clinical users
- **Temporal modeling:** Incorporating timing of diagnostic decisions and information arrival

7 Acknowledgments

We thank the Beth Israel Deaconess Medical Center for maintaining the MIMIC database and the MIT Laboratory for Computational Physiology for providing access to this valuable research resource. This work was supported by the Data 793 Capstone program.

References

- [1] Newman-Toker, D. E., Pronovost, P. J., & McDonald, K. M. (2009). Diagnostic errors—the next frontier for patient safety. *JAMA*, 301(10), 1060-1062.
- [2] Pons, E., Braun, L. M., Hunink, M. G. M., & Kors, J. A. (2016). Natural language processing in radiology: a systematic review. *Radiology*, 279(2), 329-343.

- [3] Wang, Y., Wang, L., Rastegar-Mojarad, M., Moon, S., Shen, F., Afzal, N., ... & Jiang, G. (2018). Clinical information extraction applications: a literature review. *Journal of biomedical informatics*, 77, 34-49.
- [4] Huang, K., Altosaar, J., & Ranganath, R. (2019). ClinicalBERT: Modeling clinical notes and predicting hospital readmission. *arXiv preprint arXiv:1904.05342*.
- [5] Peng, Y., Yan, S., & Lu, Z. (2019). Transfer learning in biomedical natural language processing: an evaluation of BERT and ELMo on ten benchmarking datasets. *arXiv preprint arXiv:1906.05474*.
- [6] Rogers, A., Kovaleva, O., & Rumshisky, A. (2020). A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8, 842-866.
- [7] Futoma, J., Simons, M., Panch, T., Doshi-Velez, F., & Celi, L. A. (2020). The myth of generalisability in clinical research and machine learning in health care. *The Lancet Digital Health*, 2(9), e489-e492.
- [8] Johnson, A., Pollard, T., Horng, S., Celi, L. A., & Mark, R. (2023). MIMIC-IV-Note: Deidentified free-text clinical notes (version 2.2). *PhysioNet*.
- [9] Goldberger, A., Amaral, L., Glass, L., Hausdorff, J., Ivanov, P. C., Mark, R., ... & Stanley, H. E. (2000). PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals. *Circulation*, 101(23), e215-e220.