



Type of the Paper (Article, Review, Communication, etc.)

1

# Effective Natural Language Processing Algorithms for Gout Flare Early Alert from Chief Complaints

2

3

Lucas Lopes Oliveira <sup>1+</sup>, Aryalakshmi Nellippillipathil Babu <sup>2+</sup>, Xiaorui Jiang <sup>3\*</sup>, Poonam Karajagi <sup>4</sup>, and Alireza Daneshkhan <sup>5</sup>

4

5

<sup>1</sup> School of Computing, Mathematics and Data Sciences, Coventry University; [lopesoll@uni.coventry.ac.uk](mailto:lopesoll@uni.coventry.ac.uk)

6

<sup>2</sup> School of Computing, Mathematics and Data Sciences, Coventry University; [nellippila@uni.coventry.ac.uk](mailto:nellippila@uni.coventry.ac.uk)

7

<sup>3</sup> Centre for Computational Sciences and Mathematical Modelling, Coventry University; [xiaorui.jiang@coventry.ac.uk](mailto:xiaorui.jiang@coventry.ac.uk)

8

<sup>4</sup> School of Computing, Mathematics and Data Sciences, Coventry University; [karajip@uni.coventry.ac.uk](mailto:karajip@uni.coventry.ac.uk)

9

<sup>5</sup> School of Computing, Mathematics and Data Sciences, Coventry University; [alireza.daneshkhan@coventry.ac.uk](mailto:alireza.daneshkhan@coventry.ac.uk)

10

\* Correspondence: [xiaorui.jiang@coventry.ac.uk](mailto:xiaorui.jiang@coventry.ac.uk)

11

\* Equal contributions.

12

13

14

**Abstract:** In this study, we extend the exploration of gout flare detection initiated by Osborne, J. D. et al, through the utilization of their dataset of Emergency Department (ED) triage nurse chief complaint notes. Addressing the challenge of identifying gout flares prospectively during an ED visit, where documentation is typically minimal, our research focuses on employing alternative Natural Language Processing (NLP) techniques to enhance the detection accuracy. This study investigates the application of medical domain-specific Large Language Models (LLMs), distinguishing between generative and discriminative models. Models such as BioGPT, RoBERTa-large-PubMed-M3, and BioElectra were implemented to compare their efficacy with the original implementation by Osborne, J. D. et al. The best model was Roberta-large-PM-M3 with a 0.8 F1 Score on the Gout-CC-2019 dataset followed by BioElectra with 0.76 F1 Score. We concluded that discriminative LLMs performed better for this classification task compared to generative LLMs. However, a combination of using generative models as feature extractors and employing SVM for the classification of embeddings yielded promising results comparable to those obtained with discriminative models. Nevertheless, all our implementations surpassed the results obtained in the original publication.

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

**Keywords:** keyword 1; keyword 2; keyword 3 (List three to ten pertinent keywords specific to the article yet reasonably common within the subject discipline.)

30

31

32

**Citation:** To be added by editorial staff during production.

Academic Editor: Firstname Lastname

Received: date

33

Revised: date

34

Accepted: date

35

Published: date

36



Copyright: © 2023 by the authors.

37

Submitted for possible open access publication under the terms and conditions of the Creative Commons

38

Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

39

40

41

42

43

## 1. Introduction

44

Gout affects over 9 million Americans[1,2] and is the most common form of inflammatory arthritis in men with a prevalence rate over 5%. The U.S. National Emergency Department Sample (NEDS) documents more than 200,000 annual visits where gout is identified as the primary diagnosis, constituting 0.2% of all Emergency Department visits and resulting in annual billable charges exceeding \$280 million [3].

45

Despite strides in natural language processing (NLP) techniques for detecting gout flares from textual data, the prospective identification of such instances remains a complex task, especially within the constraints of Emergency Department (ED) environments. This study addresses this critical gap by advancing the methodologies proposed by Osborne, J. D. et al [1].

46

The importance of this research lies in the need to improve the continuity of care for gout patients, especially after an ED visit. Often, gout flares treated in the ED lack optimal follow-up care, necessitating the development of methods for identifying and referring

patients with gout flares during an ED visit [1]. While retrospective studies have leveraged NLP for gout flare detection, the prospective identification of patients in real-time ED settings presents a unique challenge. The study builds upon the groundwork laid by Osborne, J. D. et al [1], who annotated a corpus of ED triage nurse chief complaint notes for gout flares, paving the way for our exploration.

### 1.1 Rationale for Using Large Language Models

Large language models, such as BERT (Bidirectional Encoder Representations from Transformers), GPT-3 (Generative Pre-trained Transformer 3), and their variants, have demonstrated remarkable success in a wide range of natural language processing tasks. The use of large language models in text classification offers several compelling reasons:

**Contextual Understanding:** Large language models leverage deep learning techniques to encode contextual information and relationships between words in a sentence. This contextual understanding allows them to capture subtle nuances and semantics, which is especially relevant in the medical domain where precise interpretation of clinical text is vital.

**Transfer Learning:** Pre-training on vast corpora of textual data enables large language models to learn general language patterns. This pre-trained knowledge can be fine-tuned on domain-specific datasets, making them adaptable and effective for text classification tasks in the medical field with relatively limited labelled data.

### 1.2 Natural Language Processing and Large Language Models in Healthcare

In recent years, the domain of healthcare has witnessed a revolutionary transformation due to the rapid advancements in Natural Language Processing (NLP) and the emergence of Large Language Models (LLMs). These technologies have the potential to revolutionize the healthcare industry by enhancing medical decision-making, patient care, and biomedical research.

Some tasks in NLP could be automated using LLM such as text classification [4, 5, 6], keyword Extraction [7, 8, 9], machine translation [10], and text summarization [11]. Furthermore, NLP and LLM can assist in the early detection and diagnosis of diseases by sifting through vast datasets to identify patterns, symptoms, and risk factors.

### 1.3 Gaps and Limitations of Current Literature

Insufficient Comparative Studies Between Domain-Specific Generative LLMs and Discriminative LLMs

While some studies have compared a single generative LLM (GPT) with discriminative LLMs, a comprehensive comparison between multiple domain-specific generative LLMs and discriminative LLMs for medical intent classification is lacking. Such comparisons are essential to determine the performance disparities between different LLM types and guide the selection of the most suitable model for our specific medical intent classification task.

In light of these gaps, our research aims to bridge these deficiencies in the current literature. We specifically focus on intent classification of medical letters by leveraging domain-specific generative LLMs as feature extractors. Additionally, our study includes comparative analyses of multiple domain specific generative LLMs and discriminative LLMs to gain comprehensive insights into their performance on this particular medical classification task. By addressing these gaps, we hope to contribute novel findings and enrich the existing literature.

In the current research landscape, the use of Large Language Models (LLMs) in the medical domain has demonstrated remarkable success. LLMs, such as Roberta-large-Pm-m3-voc, BioElectra, and BioBart, have shown promise in their ability to comprehend and process medical text [...]. The integration of these advanced models in gout flare detection

within chief nurse complaints presents an exciting avenue for exploration. Furthermore, the study introduces a novel approach of using LLMs for feature extraction, followed by classification with a Support Vector Machine (SVM), contributing to the evolving methodologies in this field [...].

## 2. Materials and Methods

### Data Collection

We utilized the dataset curated by Osborne, J. D. et al, consisting of Emergency Department (ED) triage nurse chief complaint notes. This dataset, annotated for the presence of gout flares, served as the foundation for our investigation. Each Chief Complaint (CC) in the dataset was annotated to determine its indication of a gout flare, non-indication of a gout flare, or remained unknown in terms of gout flare status. Following this, a manual chart review was conducted by a rheumatologist (MID) and a post-doctoral fellow (GR) to ascertain the gout flare status for 197 out of the 300 Emergency Department (ED) encounters. The following table, extracted from the publication by Osborne, J. D. et al., illustrates the data structure.

Table 1: GOUT-CC-2019-CORPUS Examples (Osborne, J. D. et al.)

Chief Complaint Text	Predicted*	Actual**
AMS, lethargy, increasing generalized weakness over 2 weeks. Hx: ESRD on hemodialysis at home, HTN, DM, gout, neuropathy	No	No
I started breathing hard" hx-htn, gout, anxiety,	No	No
R knee pain x 8 years. pmh: gout, arthritis	Unknown	No
Doc N Box DX pt w/ R hip FX on sat. Pt states no falls or injuries. PMH: gout	Unknown	No
out of gout medicine	Yes	Yes
sent from boarding home for increase BP and bilateral knee pain for 1 week. Hx of HTN, gout.	Yes	Yes

\*Consensus predicted gout flare status determined by annotator examination of CC

\*\*Gout flare status determined by chart review.

### Large Language Models

In this study, we harnessed the power of Large Language Models (LLMs) and transfer learning for the task of gout flare detection within Emergency Department (ED) triage nurse chief complaint notes. LLMs are state-of-the-art natural language processing models, designed to comprehend and generate human-like text trained on vast amounts of pre-existing linguistic data.

### Model Selection

We employed several LLMs tailored for the medical domain, for their ability to capture intricate patterns within medical text, making them well-suited for discerning nuances in chief complaints related to gout flares.

#### Discriminative models

In the domain of discriminative Large Language Models (LLMs), we strategically incorporated robust models renowned for their discriminative prowess—Roberta-PM-M3-Voc and BioElectra.

Model	Roberta-PM-M3-Voc	BioElectra	BioBart
Model Size	355M Parameters	---	139M Parameters

<b>Hidden Size</b>	1024	768	768
<b>Model Size</b>	24 Layers, 16 attention heads	12 Layers, 12 attention heads	12 Layers, 12 attention heads
<b>Base Model</b>	RoBERTa-large	Electra Base	Bart Base
<b>Training Data</b>	PubMed and MIMIC-III corpora	PubMed and PubMed central (millions of articles)	PubMed abstracts, PMC articles

### Generative models

In the realm of generative Large Language Models (LLMs), we strategically chose BioGPT, BioMedLM, and PMC\_LLaMA\_7B for their renowned scale and exceptional performance in natural language processing tasks. These models represent the forefront of generative language understanding, and their comprehensive specifications, training data, and architectural features are elucidated below.

Table 2: Description of Generative LLMs implemented

Model	BioGPT	BioMedLM	PMC_LLaMA_7B
<b>Model Size</b>	347M Parameters	2.7B Parameters	7B Parameters
<b>Hidden Size</b>	1024	2560	4096
<b>Model Size</b>	24 Layers, 16 attention heads	32 Layers, 20 attention heads	32 Layers, 32 attention heads
<b>Base Model</b>	GPT2-medium	GPT2	LLaMA_7B
<b>Training Data</b>	15M PubMed abstracts from scratch	All the PubMed abstracts and full documents from The Pile.	4.8 million Biomedical academic papers from the S2ORC dataset.

### Benchmark methods

To facilitate a comprehensive benchmarking analysis, we incorporated benchmark methods for comparison with our Large Language Models (LLMs). The benchmark methods involved the transformation of textual data into numerical vectors, a crucial step for machine learning algorithms that inherently require numerical input.

#### Textual Data Transformation:

Given that machine learning algorithms cannot interpret textual data directly, we employed Sklearn's 'TfidfVectorizer' algorithm to translate textual information into numerical vectors. This algorithm transforms documents into a matrix of tf-idf (term frequency-inverse document frequency) characteristics, capturing the significance of words within the corpus.

#### N-gram Exploration:

The tf-idf vectorizer, in its default setting, considers single-word tokens (unigrams) from sentences. In our research, we expanded this exploration by incorporating and evaluating various n-gram combinations of words. N-grams represent consecutive sequences of n words in a sentence. After experimentation, we opted for the (1, 2) ngram setting, utilizing both unigrams and bigrams to capture a more comprehensive contextual understanding of chief nurse complaints.

## Performance Evaluation

The performance of each model was evaluated using standard metrics, including precision, recall, and Macro F1-score. We compared our results with the original algorithm proposed by Osborne et al., ensuring a comprehensive assessment of the advancements achieved.

## 3. Results

In this section, we meticulously analyze and compare the performance of three distinct models—Roberta-large-Pm-m3-voc, BioElectra, and BioBart—on two separate datasets: Gout-cc-2019 and Gout-cc-2020. The comprehensive assessment involves a thorough examination of overall recall and F1-score metrics, providing insights into the models' respective capabilities in capturing and identifying instances of gout flares within chief nurse complaints.

### 3.1. Direct LLMs Classification

This subcategory encompasses results obtained by directly employing Large Language Models (LLMs) for the classification of Chief Complaints (CCs). Analyze and present the performance metrics, such as recall and F1-score, achieved by each LLM (Roberta-large-Pm-m3-voc, BioElectra, and BioBart, BioGPT, BioMedLM) when used independently for gout flare prediction within CCs.

Table 3: Direct LLM Classification

Model	Gout-CC-2019			Gout-CC-2020		
	Precision	Recall	F1-score	Precision	Recall	F1-score
<b>Roberta-large- PM-M3</b>	0.80	0.79	0.80	0.62	0.72	0.63
<b>BioElectra</b>	0.76	0.76	0.76	0.63	0.68	0.65
<b>BioBart</b>	0.74	0.73	0.73	0.65	0.70	0.67
<b>BioGPT</b>	0.62	0.59	0.60	0.45	0.50	0.48
<b>BioMedLM</b>	0.49	0.49	0.47	0.51	0.52	0.52

### 3.2. LLMs Embedding Extraction and Classification with SVM

In this subcategory, we explore the outcomes derived from using LLMs to extract embeddings from Chief Complaints, followed by a secondary classification using a Support Vector Machine (SVM).

Table 4: LLMs Embedding Extraction and Classification with SVM

Algorithm	Gout-CC-2019			Gout-CC-2020		
	Precision	Recall	F1-score	Precision	Recall	F1-score
<b>SVM with BioGPT Embeddings</b>	0.79	0.79	0.75	0.69	0.73	0.71*
<b>SVM with Bio- MedLM Embed- dings</b>	0.70	0.72	0.70	0.59	0.70	0.61

<b>SVM with PMC_LLaMA_7B Embeddings</b>	0.64	0.64	0.64	0.60	0.60	0.60
---	------	------	------	------	------	------

### 3.3. Benchmark Methods

This subcategory involves benchmarking the performance of traditional methods for textual data transformation, specifically focusing on the Tf-idf vectorizer with different n-gram settings. Contrast and compare the results obtained with these benchmark methods against the outcomes achieved by the LLMs, providing valuable insights into the effectiveness of each approach for gout flare prediction. In this section we have also included the results from the original publication (shaded), the results will be discussed further in the discussion section.

<b>Algorithm</b>	<b>Gout-CC-2019</b>			<b>Gout-CC-2020</b>		
	Precision	Recall	F1-score	Precision	Recall	F1-score
<b>SVM with TF-IDF</b>	0.75	0.75	0.75	0.82	0.74	0.77
<b>NAIVE-GF</b>	0.23	1.00	0.38	0.28	0.56	0.37
<b>SIMPLE-GF</b>	0.44	0.84	0.58	0.37	0.40	0.38
<b>BERT-GF</b>	0.71	0.48	0.56	0.79	0.47	0.57

## 4. Discussion

### 4.1. General Analysis

The results on the GOUT-CC-2019-CORPUS and GOUT-CC-2020-CORPUS datasets were unsatisfactory in relation to machine learning standards. The highest performance on these datasets was the SVM with BioGPT embeddings and oversampled data on the merge of the datasets with 70% accuracy but after further analysis of the results its clear the model is not able to predict the positive label as well as the negative label, and the high results of the negative class indicate a bias of the model towards the negative class, even after oversampling.

None of the models employed in this study were able to accurately make predictions of the GOUT-CC2019-CORPUS and GOUT-CC-2020-CORPUS datasets. The unsatisfactory results are related to the nature of the dataset. All the chief nurse complaints contain the keyword “gout” and most of the nurse complaints did not contain any clear indicator of gout flare. This is proven by the analysis of the predict column, where the professional annotators attempted to predict the presence of GOUT flare bases solely on the complaint. In the test set used more than half were miss classified by the professional rheumatologists.

### 4.2. Comparative Analysis

The following table compares the results acquired from this study, with the results obtained from the paper by Osborne et al.

<b>Algorithm</b>	<b>Gout-CC-2019</b>			<b>Gout-CC-2020</b>		
	Precision	Recall	F1-score	Precision	Recall	F1-score
<b>Roberta-large-PM-M3</b>	0.80	0.79	0.80*	0.62	0.72	0.63
<b>BioElectra</b>	0.76	0.76	0.76	0.63	0.68	0.65
<b>BioBart</b>	0.74	0.73	0.73	0.65	0.70	0.67
<b>BioGPT</b>	0.62	0.59	0.60	0.45	0.50	0.48

<b>BioMedLM</b>	0.49	0.49	0.47	0.51	0.52	0.52
<b>SVM with BioGPT Embeddings</b>	0.79	0.79	0.75	0.69	0.73	0.71*
<b>SVM with Bio-MedLM Embeddings</b>	0.70	0.72	0.70	0.59	0.70	0.61
<b>SVM with PMC_LLaMA_7B Embeddings</b>	0.64	0.64	0.64	0.60	0.60	0.60
<b>SVM with TF-IDF</b>	0.75	0.75	0.75	0.82	0.74	0.77*
<b>NAIVE-GF</b>	0.23	1.00	0.38	0.28	0.56	0.37
<b>SIMPLE-GF</b>	0.44	0.84	0.58	0.37	0.40	0.38
<b>BERT-GF</b>	0.71	0.48	0.56	0.79	0.47	0.57

As shown in the above table Roberta was the best performing model on the GOUT-CC-2019-CORPUS dataset followed by BioElectra, showcasing the superiority of discriminative LLMs in classification tasks. The SVM with BioGPT embedding and TF-IDF also performed well in relation to the other models. In the GOUT-CC-2020-CORPUS dataset the best LLM was SVM with BioGPT embeddings which outperformed all the discriminative LLMs due to the use of oversampling, which was not possible using the discriminative LLMs. This result was still outperformed by SVM with TF-IDF features. All of our models outperformed the models used in the study by Osborne et al.(in grey) in both datasets.

#### 4.3. Future Directions

Some improvements can be done to enhance the results obtained in this research:

**Full Fine-Tuning and Distributed Computing:** While parameter-efficient fine-tuning, specifically LoRA, was applied in this study due to hardware constraints and the models' size, pursuing full fine-tuning would enhance the results of the models. Implementing distributed computing is necessary to apply full fine tuning, due to the very large size of the models this process requires distributing the model load across different GPUs to perform the calculations. This strategy would enable more comprehensive fine-tuning, potentially leading to an increase in model performance.

**Enhanced Dataset Quality and Size:** with such a limited number of samples the model cannot be properly trained, validated and tested. To address this more samples must be acquired or whole new datasets to test the models effectively.

**Exploring Embeddings and Discriminative LLMs:** A new direction to follow would be a similar approach to the one employed in this study where the embeddings of the discriminative LLMs are extracted and used for classification using a separate classifier, in order to test the different embeddings side by side in a similar setting.

**Ensemble Learning for Enhanced Embeddings:** A promising route is the utilization of deep learning models to create an ensemble that enhances embeddings before their application in text classification. This strategy could potentially enhance the information captured by the embeddings, thereby leading to improved classification outcomes.

## 5. Conclusions

Overall this study highlighted the potential of generative LLMs for classification tasks, achieving results comparable to the discriminative models. Additionally the models also have shown potential as feature extractors for classification tasks even without fine tuning, due to their ability to understand contextual information and produce contextual rich embeddings. Despite the results between the two types of models being comparable, the computational requirements to perform the same task is much greater using the

generative LLMs employed in this study. Similar or superior results can be obtained using much smaller discriminative models. 253

Still, this research highlights the importance of using the domain specific variants of 254  
the models when the text contains specialized and out of word vocabulary. 255

Given the considerations mentioned above, the following conclusions can be drawn. 256  
The integrations of Large language models trained on medical publications holds potential 257  
to reshape classification tasks in the medical domain for the future. 258

## References 259

1. Osborne JD, Booth JS, O'Leary T, Mudano A, Rosas G, Foster PJ, Saag KG, Danila MI. Identification of Gout Flares in Chief Complaint Text Using Natural Language Processing. *AMIA Annu Symp Proc*. 2021 Jan 25;2020:973-982. PMID: 33936473; PMCID: PMC8075438. 260
2. Michael Chen-Xu, Chio Yokose, Rai Sharan K, Pillinger Michael H, Choi Hyon K. Contemporary prevalence of gout and hyperuricemia in the united states and decadal trends: the national health and nutrition examination survey, 2007–2016. *Arthritis & Rheumatology*. 2019;71(6):991–999. 261
3. Singh Jasvinder A, Shaohua Yu. Time trends, predictors, and outcome of emergency department use for gout: a nationwide us study. *The Journal of rheumatology*. 2016;43(8):1581–1588. 262
4. Xu, B., Gil-Jardiné, C., Thiessard, F., Tellier, E., Avalos, M., & Lagarde, E. (2019). Pre-training A Neural Language Model Improves The Sample Efficiency of an Emergency Room Classification Model.10.48550/arxiv.1909.01136 263
5. Veladas, R. et al. (2021). Aiding Clinical Triage with Text Classification. In: Marreiros, G., Melo, F.S., Lau, N., Lopes Cardoso, H., Reis, L.P. (eds) *Progress in Artificial Intelligence*. EPIA 2021. Lecture Notes in Computer Science(), vol 12981. Springer, Cham. [https://doi.org/10.1007/978-3-030-86230-5\\_7](https://doi.org/10.1007/978-3-030-86230-5_7) 264
6. Noor, K., Smith, K., Bennett, J., OConnell, J., Fisk, J., Hunt, M., Philippo, G., Xu, T., Knight, S., Romao, L., Dobson, R. J., & Wong, W. K. (2022). Predicting Clinical Intent from Free Text Electronic Health Records.10.48550/arxiv.2204.09594 265
7. Ding,L.,Zhang,Z.,Liu,H.,Li,J. & Yu,G.(3921).Automatic Keyphrase Extraction from Scientific Chinese Medical Abstracts Based on Character-Level Sequence Labeling. *Journal of Data and Information Science*,6(3) 35-57. <https://doi.org/10.2478/jdis-2021-0013> 266
8. Ke, H., Lee, C. S., & Sugiyama, K. (2021). Bert-Based Chinese Medical Keyphrase Extraction Model Enhanced with External Features. *Towards Open and Trustworthy Digital Societies* (pp. 167-176). Springer International Publishing AG. 10.1007/978-3-030-91669-5\_14 267
9. Tang, M., Gandhi, P., Kabir, M. A., Zou, C., Blakey, J., & Luo, X. (2019). Progress Notes Classification and Keyword Extraction using Attention-based Deep Learning Models with BERT.10.48550/arxiv.1910.05786 268
10. Han, L., Erofeev, G., Sorokina, I., Gladkoff, S., & Nenadic, G. (2022). Investigating Massive Multilingual Pre-Trained Machine Translation Models for Clinical Domain via Transfer Learning.10.48550/arxiv.2210.06068 269
11. Tang, L., Sun, Z., Idnay, B., Nestor, J. G., Soroush, A., Elias, P. A., ... & Peng, Y. (2023). Evaluating Large Language Models on Medical Evidence Summarization. *medRxiv*, 2023-04. 270

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content. 271

253  
254

255  
256

257  
258

259  
260

261  
262

263  
264

265  
266

267  
268

269  
270

271  
272

273  
274

275  
276

277  
278

279  
280

281  
282

283  
284

285  
286

287  
288

289  
290

291