

1 Article

2 Effective Natural Language Processing Algorithms for Gout 3 Flare Early Alert from Chief Complaints

4 Lucas Lopes Oliveira ¹, Xiaorui Jiang ^{2*}, Aryalakshmi Nellippillipathil Babu ³, Poonam Karajagi ⁴, Alireza
5 Daneshkhah ^{5*}

6 ¹ School of Computing, Mathematics and Data Sciences, Coventry University; lopesoll@uni.coventry.ac.uk
7 ² Centre for Computational Sciences and Mathematical Modelling, Coventry University;
8 xiaorui.jiang@coventry.ac.uk
9 ³ School of Computing, Mathematics and Data Sciences, Coventry University; nellippila@uni.coventry.ac.uk
10 ⁴ School of Computing, Mathematics and Data Sciences, Coventry University; karajagip@uni.coventry.ac.uk
11 ⁵ School of Computing, Mathematics and Data Sciences, and Centre for Computational Sciences and Mathe-
12 matical Modelling, Coventry University; alireza.daneshkhan@coventry.ac.uk
13 * Correspondence author.

14 **Abstract:** Early identification of acute gout is crucial, enabling healthcare professionals to imple-
15 ment targeted interventions for rapid pain relief and preventing disease progression, ensuring
16 improved long-term joint function. In this study, we comprehensively explored the potential of
17 gout flare (GF) early detection based on nurse chief complaint notes in the Emergency Department
18 (ED). Addressing the challenge of identifying GFs prospectively during an ED visit, where docu-
19 mentation is typically minimal, our research focuses on employing alternative Natural Language
20 Processing (NLP) techniques to enhance the detection accuracy. We investigate GF detection algo-
21 rithms using both sparse representations by traditional NLP methods and dense encodings by
22 medical domain-specific Large Language Models (LLMs), distinguishing between generative and
23 discriminative models. Three methods are used to alleviate the issue of severe data imbalance, in-
24 cluding oversampling, class weights, and focal loss. Extensive empirical studies are done on the
25 Gout Emergency Department Chief Complaint Corpora. Sparse text representations like tf-idf
26 proved to produce strong performance, achieving higher than 0.75 F1 Score. The best deep learning
27 models are RoBERTa-Large-PM-M3-Voc and BioGPT, with the best F1 Scores on each dataset with
28 a 0.8 on the 2019 dataset and a 0.85 F1 Score the 2020 dataset. We concluded that although dis-
29 criminative LLMs performed better for this classification task, compared to generative LLMs, a
30 combination of using generative models as feature extractors and employing support vector ma-
31 chine for classification yields promising results comparable to those obtained with discriminative
32 models.

Citation: To be added by editorial
staff during production.

Academic Editor: Firstname
Lastname

Received: date
Revised: date
Accepted: date
Published: date



Copyright: © 2024 by the author
Submitted for possible open access
publication under the terms and
conditions of the Creative Commons
Attribution (CC BY) license
(<https://creativecommons.org/licenses/by/4.0/>)

33 **Keywords:** Gout Flare; Chief Complaint; Natural Language Processing; Deep Learning; Large
34 Language Models

36 1. Introduction

37 More than 9 million Americans suffer from gout [1], which is the most prevalent
38 type of inflammatory arthritis among men, affecting over 5% of them. According to the
39 U.S. National Emergency Department Sample (NEDS), gout accounts for more than
200,000 visits to the Emergency Department (ED) every year, making up 0.2% of all ED
visits and costing more than \$280 million in annual charges [2]. It is important to improve
the continuity of care for gout patients, especially after an ED visit. Often, gout flares
(GF) treated in the ED lack optimal follow-up care, necessitating the development of
methods for identifying and referring patients with GFs during an ED visit [3]. While
retrospective studies have leveraged NLP for GF detection, the prospective identification

46 of patients in real time ED settings presents a unique challenge, especially within the
47 constraints of Emergency Department (ED) environments.

48 Despite of the success of natural language processing (NLP) techniques in healthcare
49 [4], NLP-based Gout Flare Early Detection (GFED) is in severe lack of study. Only a few
50 were identified, like Zheng et al [5], which however worked on Electronic Medical Rec-
51 ords. The problem of early warning of acute GFs becomes more challenging in the ED
52 setting where only chief complaints of patients are taken by nurses in an extremely suc-
53 cinct format. It is of paramount challenge to develop an effective GFED algorithm using
54 such limited amount of information. The current study tries to address this critical gap by
55 advancing the methodologies proposed by Osborne et al [3]. Our study builds upon the
56 groundwork laid by Osborne et al., who annotated two corpora of ED chief complaint
57 notes for GFs and paves the way for our exploration of effective text representation
58 methods and state-of-the-art medical/clinical Large Language Models (LLM).

59 *1.1 Rationale for Using Large Language Models*

60 Large language models, such as BERT [6] (Bidirectional Encoder Representations
61 from Transformers), [7] (Generative Pre-trained Transformer 3), and their variants, have
62 demonstrated remarkable success in a wide range of natural language processing tasks.
63 The use of large language models in text classification offers several compelling reasons:

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

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

73 These technologies have the potential to revolutionize the healthcare industry by
74 enhancing medical decision-making, patient care, and biomedical research. Some tasks in
75 NLP could be automated using LLM such as text classification [8-9], keyword Extraction
76 [10-11], machine translation [12], and text summarization [13]. Furthermore, NLP and
77 LLM can assist in the early detection and diagnosis of diseases by sifting through vast
78 datasets to identify patterns, symptoms, and risk factors.

79 *1.2 Gaps and Limitations of Current Literature*

80 While some studies have compared a single generative LLM (GPT) with discrimi-
81 native LLMs, a comprehensive comparison between multiple domain-specific generative
82 LLMs and discriminative LLMs for disease detection is lacking. Such comparisons are
83 essential to determine the performance disparities between different LLM types and
84 guide the selection of the most suitable model for our specific medical intent classifica-
85 tion task.

86 In light of these gaps, our research aims to bridge these deficiencies in the current
87 literature. We specifically focus on GFED by leveraging domain-specific generative LLMs
88 as feature extractors. Additionally, our study includes comparative analyses of multiple
89 domain specific generative LLMs and discriminative LLMs to gain comprehensive in-
90 sights into their performance on this particular medical classification task.

91 *1.3 Our contributions*

92 In this paper, we make three contributions to the task of gout flare detection from
93 nurse chief complaints. First, we compare the performance of domain specific discrimi-
94 native and generative models that are fine-tuned for the task. Second, we propose an al-
95 ternative approach that uses domain specific generative LLMs as feature extractors and

96 support vector machine as classifier. Third, we benchmark our methods against a base-
97 line that uses sparse text representation (tf-idf). Our results demonstrate the effectiveness
98 of using LLMs, such as RoBERTa-Large-PM-M3-Voc, BioELECTRA, and BioGPT, for
99 processing medical text and detecting GFs.

100 2. Materials and Methods

101 2.1 Data Collection

102 We utilized the dataset of ED chief complaint notes which were annotated by Os-
103 borne et al. for the presence of GFs [14]. Each CC text in the dataset was annotated to
104 determine its indication of a GF, a non-GF, or remained unknown in terms of the status of
105 GF. Following this, a manual chart review was conducted by a rheumatologist and a
106 post-doctoral fellow to ascertain the GF status for a small portion of the ED counters.
107 These were served as the gold standard annotations of the real GF status. The corpora
108 contain two datasets for the year 2019 and 202, namely GOUT-CC-2019-CORPUS and
109 GOUT-CC-2020-CORPUS respectively. Table 1 shows the annotation statistics of the two
110 datasets (from Osborne et al. [3]), while Table 2 illustrates some examples. In our ex-
111 periments, we used the human-annotated samples using Chart Review, as what Osborne
112 et al. did.

113 Table 1: Annotation Statistics of the Gout Flare Chief Complaint Datasets (Osborne et al. [3])

Dataset Name	GF-POS (Positive)	GF-NEG (Negative)	GF-UNK (Unknown)	Review	Agreement	Cohen's κ
GOUT-CC-2019-CORPUS	93	194	13	CC	0.883	0.825
GOUT-CC-2019-CORPUS*	70	118	9	Chart	0.849	0.774
GOUT-CC-2020-CORPUS	14	7992	129	CC	0.977	0.965
GOUT-CC-2020-CORPUS*	25	232	7	Chart	0.904	0.856

114 * Used for experiments as Osborne et al. [3]

115 Table 2: Examples of Chief Complaint Notes for Gout Flare (Osborne et al. [3])

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

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

117 **Gout flare status determined by chart review.

118 2.2 Feature Extraction

119 In the feature engineering approach, we extracted the n -grams ($n = 1, 2, 3$) and tested
120 different combinations of n -grams and different feature sizes. CC texts were converted
121 into sparse representations using *tf-idf* (Term Frequency-Inverse Document Frequency)
122 [15] as initial feature values. A linear support vector classifier (Linear SVC) was trained.
123 All implementations were done using the scikit-learn library¹.

¹ <https://scikit-learn.org/>

124 It was hard to extract more advanced syntactic or semantic features due to the
125 noisiness of CC texts. As can be observed from Table 2, CC texts are extremely succinct,
126 often containing a sequence of medical terms or abbreviations, which record the facts
127 reported by patients. Such CCs are not meaningful sentences for us to extract features
128 from the syntactic analysis results. Semantic analysis tools are either immature or non-
129 existent in this particular area. However, we could still observe quite good performances
130 from fine-tuning a machine learning model using the right sparse feature representation
131 of CC texts.

132 2.3 Large Language Models

133 We employed several LLMs tailored for the medical domain, for their ability to
134 capture intricate patterns within medical text, making them well-suited for discerning
135 nuances in chief complaints related to GF. All LLMs belong to the Transformers family
136 [16] because we hoped that the multi-headed self-attention mechanism of the Trans-
137 formers architecture could be able to learn the meaningful association between certain
138 words of CC texts to indicate the existence of GF.

139 2.3.1 Discriminative models

140 We strategically incorporated three robust discriminative LLMs renowned for their
141 discriminative power—RoBERTa-PM-M3-Voc², BioELECTRA³ [17], and BioBART⁴ [18].
142 These are the domain-specific versions of the RoBERTa [19], Electra [20] and BART [21]
143 models respectively. Although BART was a language model pretrained in a se-
144 quence-to-sequence fashion, it can be used equally well and in the same way as a dis-
145 criminative model [21]. So, we treated it as one representative of the discriminative cat-
146 egory. The details of the discriminative LLMs are shown in Table 3.

147 Table 3: Description of Discriminative LLMs Implemented

Model	RoBERTa-PM-M3-Voc	BioELECTRA	BioBART
Model Size	355M Parameters	---	139M Parameters
Hidden Size	1024	768	768
Model Size	24 Layers, 16 heads	12 Layers, 12 heads	12 Layers, 12 heads
Base Model	RoBERTa-large	Electra Base	BART Base
Training Data	PubMed articles and MIMIC-III corpora ⁵ [22]	PubMed articles	PubMed abstracts and articles

148 2.3.2 Generative models

149 In the realm of generative LLMs, we strategically chose BioGPT⁶ [23], BioMedLM⁷,
150 and PMC_LLaMA_7B⁸ [24] for their renowned scale and exceptional performance in
151 natural language processing tasks. BioGPT and PMC_LLaMA_7B are the domain-specific
152 versions of the GPT-2 [25] and LLaMA [26-27] models respectively, while BioMedLM is a
153 bespoke LLM pretrained for medical applications. These models represent the forefront
154

² <https://huggingface.co/Sedigh/RoBERTa-large-PM-M3-Voc>

³ <https://github.com/kamalkraj/BioELECTRA>

⁴ <https://github.com/GanjinZero/BioBART>

⁵ <https://www.nature.com/articles/sdata201635>

⁶ https://huggingface.co/docs/transformers/model_doc/biogpt

⁷ <https://github.com/stanford-crfm/BioMedLM>

⁸ <https://github.com/chaoyi-wu/PMC-LLaMA>

of generative language understanding, and their comprehensive specifications, training data, and architectural features are elucidated in Table 4.

Table 4: Description of Generative LLMs Implemented

Model	BioGPT	BioMedLM	PMC_LLaMA_7B
Model Size	347M Parameters	2.7B Parameters	7B Parameters
Hidden Size	1024	2560	4096
Model Size	24 Layers, 16 heads	32 Layers, 20 heads	32 Layers, 32 heads
Base Model	GPT2-medium	GPT2	LLaMA_7B
Training Data	15M PubMed abstracts from scratch	All PubMed abstracts and full texts from The Pile benchmark [28].	4.8 million Biomedical publications from the S2ORC dataset [29].

2.4 Fine-tuning

Fine tuning was implemented to improve the models' ability to understand and capture the nuances in the texts. For the discriminative models full fine tuning was implemented, but for the generative models due to the size of the models and hardware constraints full fine tuning was not possible.

2.4.1 Fine-tuning of Discriminative LLMs

All three discriminative LLMs use a bidirectional encoder as BERT [6]. The encoder part of these models was used to encode each CC text, and the "[CLS]" token was used as the dense representation. For RoBERTa-PM-M3-Voc and BioELECTRA, a further feature transformation was applied. Essentially, the classification head was a Multiple Layer Perceptron (MLP), the hidden layer of which made a nonlinear transformation (of the same size). On the contrary, BioBART used a linear classification head following the tradition of BART usage.

In the fine-tuning process, the following hyperparameters were used: learning rate = $1e-5$, epoch number = 10, batch size = 14, early stopping patience = 3. The AdamW optimiser was used for training [30].

2.4.2 Fine-tuning of Generative LLMs

Similarly, generative LLMs were used for encoding CC texts, and the "Extract" token (for all three models as they all belong to the GPT family) were used to extract the dense representation, which was then sent to a linear classification head. Due to their large sizes, the generative LLMs were not fully fine-tuned. Instead, we used LoRA (Low Rank Adaptation) to efficiently adapt LLMs to specific tasks by only modifying a small portion of the whole parameter space.

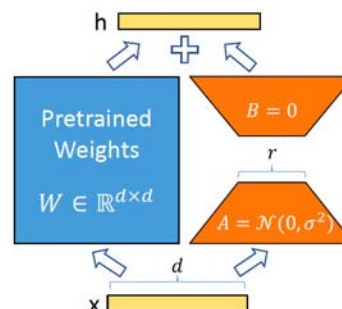
The main idea behind LoRA is to exploit the low-rank structure of the model's weight matrices during task adaptation, resulting in reduced memory usage and computational complexity [31]. The idea was inspired by Aghajanyan et al.'s finding that pre-trained language models have a low "intrinsic dimension" meaning that they can still learn efficiently when their weight matrices are randomly projected to a smaller subspace [32].

More precisely, LoRA hypothesizes that updates to model's weight matrix, W_0 , can be represented by a low-rank decomposition, which is given by $W_0 + \Delta W = W_0 + BA$, where $B \in R^{d \times r}$, $A \in R^{r \times k}$, and $\Delta W = BA$ represents weight updates. During training (i.e., fine-tuning), W_0 is frozen while A and B contain the trainable parameters.

In our fine-tuning process, we applied the following LoRA parameters:

1. The rank (r) of A and B was set to 8.
2. The LoRA regularization coefficient α was set to 16.

- 196
197
198
199
3. To prevent overfitting and enhancing model generalisation, we applied a LoRA dropout rate of 0.1.
 4. A learning rate of $3e-4$ was used, enabling efficient convergence during training.



200
201 Figure 1: Parametrization of LoRA. Only A and B are trained. [31]

202 2.5 Classification

203 In the feature engineering approach, a Linear SVC was trained. When finetuning
204 discriminative LLMs, either an MLP or a linear classifier was applied. Similarly, a linear
205 layer was used for classification with generative LLMs. In the experiments, we also tested
206 using generative LLMs only as the feature extractor and trained a Linear SVC for classi-
207 fication. In this alternative approach, which required significantly less computational
208 resources, generative LLMs were frozen, used to encode CC texts, and the hidden states
209 of the “Extract” token were extracted as dense representation. A Linear SVC was then
210 trained in the similar way as in the feature engineering approach. This was to demon-
211 strate LLMs’ native ability to understand and represent medical texts for the downstream
212 task.

213 2.6. Optimisation

214 2.6.1 Class weights

215 We also observed severe data balance in the corpora. The data imbalance ratio of
216 GOUT-CC-2019 is $(70 + 9) / 118 = 0.6695$, while the imbalance ratio of GOUT-CC-2020 is
217 $(25 + 7) / 232 = 0.1379$. Our first method to handle data imbalance was class weights [33],
218 which were set according to the relative sizes of each class as in Eq. (1),

$$w_j = N / (K \times N_j), \quad (1)$$

219 where w_j is the weight for the j -th class, K is the total number of classes, N is the
220 total number of samples, and N_j is the number of samples of the j -th class [34].

221 2.6.2 Oversampling

222 However, class weighting in Eq. (1) did not help improve the performances on
223 GOUT-CC-2020 much, which is 5 times more imbalanced than GOUT-CC-2019. Alt-
224 hough the discriminative LLMs performed strongly in our experiments, they were ex-
225 tremely sensitive to this severe data imbalance. Therefore, we performed random over
226 sampling on GOUT-CC-2020. The positive samples in the training split, including
227 GF-POS and GF-UNK combined, were randomly duplicated to match the size of
228 GF-NEG.

229 The second approach we used to oversample the minority class was Synthetic Mi-
230 nority Over-sampling Technique (SMOTE) [35]. SMOTE generates synthetic examples of
231 then minority class by interpolating the feature space of the existing minority samples. By
232 doing so, SMOTE effectively oversamples the minority class, thereby balancing the class
233 distribution [35]. This approach was only implemented in the method where we used the
234 LLMs as feature extractors and classified with the SVC.

235
236
237
238
239
240
241
242
243

244
245
246

247
248
249
250
251
252

253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268

269

2.6.3 Focal Loss

In the context of our classification tasks, the choice of a suitable loss function plays a pivotal role in training and optimizing our models. We employed two distinct loss functions as per dataset and model requirement, namely cross-entropy loss and focal loss [36], to effectively guide the training process and address specific challenges posed by our datasets.

In instances where class imbalance persisted even after oversampling the training data, such as in the case of GOUT-CC-2020, we employed focal loss as an alternative to cross-entropy to combat class imbalance in classification tasks, as in Eq. (2).

$$FL(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t), \quad (2)$$

where p_t is the posterior probability of each target t (here $t = 0$ or 1), $\alpha_t \in [0,1]$ is the scaling parameter, γ is the focusing parameter and $(1 - p_t)^\gamma$ is the modulating factor of the original cross-entropy loss [36].

3. Results

In this section, we meticulously analyze and compare the performances of all methods. 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. [3], ensuring a comprehensive assessment of the advancements achieved.

3.1. Fine-tuned LLM

This subcategory encompasses results obtained by directly employing LLMs for CC classification. Table 5 shows the results.

The table shows that RoBERTa-Large-PM-M3-Voc outperforms the other four models in the 2019 dataset in terms of precision, recall, and F1-score for both datasets. This suggests that this model is more effective at detecting GFs from clinical notes. Table 5 also shows that BioBERT and BioELECTRA have similar performance, while BioGPT and BioMedLM have the lowest performance among the five models.

On the 2020 dataset, the best model was by far BioGPT, outperforming others LLM competitors by large margins. Good performances were obtained due to oversampling, which improved the results from 0.67 to 0.85 macro f1 score. These results suggest that BioGPT can handle the data imbalance and the domain-specific vocabulary better than the other models, and that oversampling can boost the performance of generative LLMs for this task. On the other hand, BioMedLM did not achieve good performances, possibly due to the limitations of the LoRA adaptor, compared to BioGPT which was fully finetuned to adapt better to the special domain of gout flare CC texts.

Table 5: Performances of Gout Flare Detection using Fine-Tuned LLMs

Model	GOUT-CC-2019			GOUT-CC-2020		
	Precision	Recall	F1-score	Precision	Recall	F1-score
RoBERTa-Large-PM-M3-Voc	0.80	0.79	0.80	0.62	0.72	0.63
BioELECTRA	0.76	0.76	0.76	0.63	0.68	0.65
BioBART	0.74	0.73	0.73	0.65	0.70	0.67
BioGPT	0.62	0.59	0.60	0.82	0.88	0.85
BioMedLM	0.49	0.49	0.47	0.52	0.53	0.52

270 *3.2 Frozen LLMs as Feature Extractors*

271 In this subcategory, we used LLMs to embed CC texts to dense feature vectors and
272 use Linear SVC for classification. Table 6 shows the results.

273 The table shows that SVM with BioGPT Embeddings has the best performance
274 among the four algorithms on both datasets. It achieves an F1-score of 0.67 on
275 Gout-CC-2019 and 0.71 on Gout-CC-2020. This indicates that this algorithm can effec-
276 tively extract the relevant features from CC texts and classify them accurately.

277 The table also shows that SVM with BioMedLM Embeddings and SVM with
278 PMC_Llama_7B Embeddings have similar performance, but lower than SVM with
279 BigGPT Embeddings. They both have an F1-score of 0.66 on Gout-CC-2019 and 0.61 on
280 Gout-CC-2020. This suggests that these algorithms are less robust and consistent in han-
281 dling the variability and complexity of CC texts.

282 Table 6: Performances of Gout Flare Detection using LLM Embeddings

Algorithm	Gout-CC-2019			Gout-CC-2020		
	Precision	Recall	F1-score	Precision	Recall	F1-score
SVM with BioGPT Embeddings	0.68	0.67	0.67	0.69	0.73	0.71
SVM with BioMedLM Embeddings	0.69	0.66	0.66	0.59	0.70	0.61
SVM with PMC_LLaMA_7B Embeddings	0.66	0.66	0.66	0.60	0.60	0.60

283 *3.3 Sparse Text Representation*

284 This subcategory involves performance of the traditional feature engineering ap-
285 proach, which generated sparse text representations using tf-idf of n -gram features.
286 Contrast and compare these results against the outcomes achieved by the LLMs,
287 providing valuable insights into the effectiveness of each approach for GF prediction. In
288 this section we have also included the results from the original publication of Osborne et
289 al. [3], which are shaded. All results will be discussed further in the discussion section.
290 Table 7 shows the results.

291 Table 7: Performances of Gout Flare Detection using Sparse Text Representations

Algorithm	GOUT-CC-2019			GOUT-CC-2020		
	Precision	Recall	F1-score	Precision	Recall	F1-score
SVM with tf-idf	0.75	0.75	0.75	0.82	0.74	0.77
NAIVE-GF	0.23	1.00	0.38	0.28	0.56	0.37
SIMPLE-GF	0.44	0.84	0.58	0.37	0.40	0.38
BERT-GF	0.71	0.48	0.56	0.79	0.47	0.57

292 **4. Discussion**

293 *4.1 Comparative Analysis*

294 The following table compares the results acquired from this study, with the results
295 obtained from the paper by Osborne et al. As shown in Table 8, RoBERTa was the best
296 performing model on the GOUT-CC-2019-CORPUS dataset followed by BioELECTRA,
297 showcasing the superiority of discriminative LLMs in classification tasks. The SVM with
298 BioGPT embedding and tf-idf also performed well in relation to the other models. In the
299 GOUT-CC-2020-CORPUS dataset the best was BioGPT which outperformed all the dis-

300 criminative LLMs. This model responded very well to the fine tuning and oversampling.
301 This result was still outperformed by SVM with tf-idf features. All our models outper-
302 formed the models used in the study by Osborne et al. (in grey) in both datasets. Overall,
303 RoBERTa-Large-PM-M3-Voc , BioGPT and tf-idf on n -grams were more robust models
304 across datasets, particularly the latter. In addition, BioGPT was a more robust feature
305 extractor when model parameters were frozen. Finally, a promising future direction to
306 employ the strengths of different classifier to achieve better recall while at the meantime
307 keeping a better balance for precision.

308 Table 8: Comparing the Performances of All Gout Flare Detection Methods.

Algorithm	GOUT-CC-2019			GOUT-CC-2020		
	Precision	Recall	F1-score	Precision	Recall	F1-score
RoBERTa-Large-PM-M3-Voc	0.80	0.79	0.80*	0.62	0.72	0.63
BioELECTRA	0.76	0.76	0.76	0.63	0.68	0.65
BioBART	0.74	0.73	0.73	0.65	0.70	0.67
BioGPT	0.62	0.59	0.60	0.82	0.88	0.85
BioMedLM	0.49	0.49	0.47	0.52	0.53	0.52
SVM with BioGPT Embeddings	0.68	0.67	0.67	0.69	0.73	0.71
SVM with BioMedLM Embeddings	0.69	0.66	0.66	0.59	0.70	0.61
SVM with PMC_LLaMA_7B Embeddings	0.66	0.66	0.66	0.60	0.60	0.60
SVM with tf-idf	0.75	0.75	0.75	0.82	0.74	0.77
NAIVE-GF	0.23	1.00	0.38	0.28	0.56	0.37
SIMPLE-GF	0.44	0.84	0.58	0.37	0.40	0.38
BERT-GF	0.71	0.48	0.56	0.79	0.47	0.57

309 4.2 Potential and limitations

310 The best performance on these datasets was achieved by
311 RoBERTa-large-PM-M3-Voc, which outperformed other LLMs and traditional machine
312 learning algorithms. This suggests that RoBERTa-Large-PM-M3-Voc can effectively
313 capture the semantic features of CC texts and distinguish between GF and non-flares.
314 However, the results also show that there is still a large gap between the performance of
315 LLMs and the desired accuracy for GF detection.

316 Furthermore, the results also indicate that some models have a bias towards the
317 negative class, which may affect their ability to predict the positive label. Therefore, more
318 research is needed to address these challenges and improve the performance of LLMs for
319 GF detection. One of the main challenges is the nature of the dataset. All the chief com-
320 plaints contain the keyword “gout” and most of them did not contain any clear indicator
321 of gout flare. This makes it difficult for the models to learn the subtle differences between
322 gout flares and non-flares. Upon analysing the predict column of our test set (which
323 contains the prediction of the human annotators based solely on the CC) we found that
324 this is a challenging problem even for professional rheumatologists which achieved less
325 than 50% accuracy in our test set.

326 Although the performance on GOUT-CC-2020-CORPUS was not as good as
327 GOUT-CC-2019-CORPUS, it’s still an improvement compared to the baseline. We
328 acknowledge that the dataset is challenging due to its data imbalance and small size,

329 which contributed to the performance decline. Our approaches to tackling the data im-
330 balance did improve the performance but future work is still required to tackling this
331 issue. One potential direction is the use of semi-supervised learning do deal with the low
332 number of annotated CC's and another is to encourage the medical community to share
333 or annotate more data to create high-quality datasets.

334 *4.3 Future Directions*

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

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

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

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

350 **Task-specific continuous pre-training:** Another possible direction is to use unsu-
351 pervised learning to continuously pre-train the LLMs on the task-specific data, i.e., the
352 chief complaint texts. This could help the models to adapt to the domain and the vocabu-
353 lary, and to tackle the particular write styles of keeping CC notes in the task.

354 **5. Conclusions**

355 Overall, this study highlighted the potential of generative LLMs for classification
356 tasks, achieving results comparable to the discriminative models. Additionally, the
357 models also have shown potential as feature extractors for classification tasks even
358 without fine tuning, due to their ability to understand contextual information and pro-
359 duce contextual rich embeddings. Despite the results between the two types of models
360 being comparable, the computational requirements to perform the same task is much
361 greater using the generative LLMs employed in this study. Similar or superior results can
362 be obtained using much smaller discriminative models. Still, this research highlights the
363 importance of using the domain specific variants of the models when the text contains
364 specialized and out of word vocabulary. Our results are important because they demon-
365 strate the feasibility and effectiveness of using generative LLMs for gout flare detection
366 from chief complaints, which is a novel and challenging task that can benefit both clinical
367 practice and research. Furthermore, our approaches can potentially improve the quality
368 of care for gout patients, a large portion of them could now receive proper and in-time
369 follow-up after an ED visit.

370
371 **Author Contributions:** Conceptualization, X. Jiang, and A. Daneshkakh; methodology, X. Jiang,
372 L.L. Oliveira, A.N., Babu, P. Karajagi, and A. Daneshkakh; software, L.L. Oliveira, A.N., Babu, P.
373 Karajagi, and X. Jiang; validation, L.L. Oliveira, and X. Jiang; investigation, L.L. Oliveira, A.N.,
374 Babu, P. Karajagi, and X. Jiang; resources, X. Jiang; data curation, L.L. Oliveira; writing—original
375 draft preparation, L.L. Oliveira, and X. Jiang; writing—review and editing, L.L. Oliveira, X. Jiang,
376 and A. Daneshkakh; supervision, X. Jiang, and A. Daneshkakh; project administration, X. Jiang. All
377 authors have read and agreed to the published version of the manuscript.

378 **Data Availability Statement:** The dataset the current paper used is a public dataset, which is
379 available through PhysioNet at <https://doi.org/10.13026/96v3-dw72>.

380 **Conflicts of Interest:** The authors declare no conflict of interest.

381 References

- 382 1. Chen, X.M.; Yokose, C.; Rai, S.K.; Pillinger, M.H.; Choi, H.K. Contemporary Prevalence of Gout and Hyperuricemia in the
383 United States and Decadal Trends: The National Health and Nutrition Examination Survey, 2007-2016. *Arthritis Rheumatol*
384 **2019**, 71(6), 991–999. doi:[10.1002/art.40807](https://doi.org/10.1002/art.40807).
- 385 2. Singh, J.A.; Yu, S. Time Trends, Predictors, and Outcome of Emergency Department Use for Gout: A Nationwide US Study. *J*
386 *Rheumatol* **2016**, 43(8), 1581–1588. doi:[10.3899/jrheum.151419](https://doi.org/10.3899/jrheum.151419).
- 387 3. Osborne, J.D.; Booth, J.S.; O’Leary, T.; et al. Identification of Gout Flares in Chief Complaint Text Using Natural Language
388 Processing. *AMIA Annu Symp Proc.* **2020**, 973–982.
- 389 4. Hossain, E.; Rana, R.; Higgins, N.; Soar, J.; Barua, P.D.; Pisani, A.R.; Turner, K. Natural Language Processing in Electronic
390 Health Records in relation to healthcare decision-making: A systematic review. *Comput Biol Med* **2023**, 155, 106649. doi:
391 [10.1016/j.compbiomed.2023.106649](https://doi.org/10.1016/j.compbiomed.2023.106649).
- 392 5. Zheng, C.; Rashid, N.; Wu, Y.; et al. Using Natural Language Processing and Machine Learning to Identify Gout Flares From
393 Electronic Clinical Notes. *Arthritis Care Res* **2014**, 66(11), 1740–1748. doi: [10.1002/acr.22324](https://doi.org/10.1002/acr.22324).
- 394 6. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Under-
395 standing. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Ling-
396 uistics: Human Language Technologies, (NAACL-HLT’2019), 4171–4186, Minneapolis, MN, USA, 2nd-7th June 2019. doi:
397 [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423).
- 398 7. Brown, T.; Mann, B.; Ryder, N.; et al. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing*
399 *Systems*; Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M.F., Lin, H., Eds.; Curran Associates Inc., 2020; Volume 33., pp.
400 1877–1901. <https://papers.nips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- 401 8. Xu, B.; Gil-Jardiné, C.; Thiessard, F.; Tellier, E.; Avalos, M.; Lagarde, E. Pre-training A Neural Language Model Improves The
402 Sample Efficiency of an Emergency Room Classification Model. In Proceedings of The Thirty-Third International FLAIRS
403 Conference (FLAIRS-33), North Miami Beach, Florida, USA, 17-20 May 2020. <https://aaai.org/papers/264-flairs-2020-18444/>.
- 404 9. Veladas, R.; Yang, H.; Quaresma, P.; et al. Aiding Clinical Triage with Text Classification. In *Progress in Artificial Intelligence,*
405 *Lecture Notes in Computer Science*; Marreiros, G., Melo, F.S., Lau, N., Lopes Cardoso, H., Reis, L.P., Eds.; Springer International
406 Publishing; 2021; Vol 12981, pp. 83–96. doi: [10.1007/978-3-030-86230-5_7](https://doi.org/10.1007/978-3-030-86230-5_7).
- 407 10. Ding, L.; Zhang, Z.; Liu, H.; Li, J.; Yu, G. Automatic Keyphrase Extraction from Scientific Chinese Medical Abstracts Based on
408 Character-Level Sequence Labeling. *J Data Inf Sci* **2021**, 6(3), 35–57. doi: [10.2478/jdis-2021-0013](https://doi.org/10.2478/jdis-2021-0013).
- 409 11. Ding, L.; Zhang, Z.; Zhao, Y. Bert-Based Chinese Medical Keyphrase Extraction Model Enhanced with External Features. In:
410 *Towards Open and Trustworthy Digital Societies, Lecture Notes in Computer Science*; Ke, H.R., Lee, C.S., Sugiyama, K., Eds.; Springer
411 International Publishing; 2021; Volume 13133, pp. 167–176. doi: [10.1007/978-3-030-91669-5_14](https://doi.org/10.1007/978-3-030-91669-5_14).
- 412 12. Han, L.; Erofeev, G.; Sorokina, I.; Gladkoff, S.; Nenadic, G.; Investigating Massive Multilingual Pre-Trained Machine Transla-
413 tion Models for Clinical Domain via Transfer Learning. In Proceedings of the 5th Clinical Natural Language Processing
414 Workshop (ClinicalNLP’2019), 31–40, Minneapolis, MN, USA, 7th June 2019. doi: [10.18653/v1/2023.clinicalnlp-1.5](https://doi.org/10.18653/v1/2023.clinicalnlp-1.5).
- 415 13. Tang, L.; Sun, Z.; Iday, B.; et al. Evaluating Large Language Models on Medical Evidence Summarization. *npj Digit. Med.* **2003**,
416 6, Article No. 158. doi: [10.1038/s41746-023-00896-7](https://doi.org/10.1038/s41746-023-00896-7).
- 417 14. Osborne, J. D., O’Leary, T., Mudano, A., Booth, J., Rosas, G., Peramsetty, G. S., Knighton, A., Foster, J., Saag, K., & Danila, M. I.
418 Gout Emergency Department Chief Complaint Corpora (version 1.0). PhysioNet, 2020. doi: [10.13026/96v3-dw72](https://doi.org/10.13026/96v3-dw72).
- 419 15. Manning, C.D.; Raghavan, P.; Schütze, H. *Introduction to Information Retrieval*; Cambridge University Press: Cambridge, USA,
420 2008.
- 421 16. Vaswani, A.; Shazeer, N.; Parmar, N.; et al. Attention is All you Need. In *Advances in Neural Information Processing Systems*;
422 Guyon, I., Luxburg, U.V., Bengio, S., et al., Eds.; Curran Associates, Inc., 2017; Volume 30, pp. 5998–6008.
423 https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.
- 424 17. Kanakarajan, K.R.; Kundumani, B.; Sankarasubbu, M. BioELECTRA: Pretrained Biomedical Text Encoder using Discrimina-
425 tors. In Proceedings of the 20th Workshop on Biomedical Language Processing (BioNLP’2021), 143–154. Online, 16 August
426 2021. doi: [10.18653/v1/2021.bionlp-1.16](https://doi.org/10.18653/v1/2021.bionlp-1.16).
- 427 18. Yuan, H.; Yuan, Z.; Gan, R.; Zhang, J.; Xie, Y.; Yu, S. BioBART: Pretraining and Evaluation of A Biomedical Generative Lan-
428 guage Model. In Proceedings of the 21st Workshop on Biomedical Language Processing (BioNLP’2022), 97–109. Dublin, Ire-
429 land, w6th May 2022. doi: [10.18653/v1/2022.bionlp-1.9](https://doi.org/10.18653/v1/2022.bionlp-1.9).
- 430 19. Liu, Y.; Ott, M.; Goyal, N.; et al. RoBERTa: A Robustly Optimized BERT Pretraining Approach. <https://arxiv.org/abs/1907.11692>.
- 431 20. Clark, K.; Luong, M.T.; Le, Q.V.; Manning, C.D. ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Gener-
432 ators. In Proceedings of the Eighteenth International Conference on Learning Representations (ICLR’2020). Online, 27th-30th
433 April 2020. <https://openreview.net/forum?id=r1xMH1BtvB>.
- 434 21. Lewis, M.; Liu, Y.; Goyal, N.; et al. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation,
435 Translation, and Comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics
436 (ACL’2020), 7871–7880. Online, 5th-10th July 2020. doi: [10.18653/v1/2020.acl-main.703](https://doi.org/10.18653/v1/2020.acl-main.703).
- 437 22. Johnson, A.E.W.; Pollard, T.J.; Shen, L.; et al. MIMIC-III, a freely accessible critical care database. *Sci Data* **2016**, 3(1): 160035. doi:
438 [10.1038/sdata.2016.35](https://doi.org/10.1038/sdata.2016.35).

- 439 23. Luo, R.; Sun, L.; Xia, Y.; et al. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Brief*
440 *Bioinform* **2022**, 23(6), bbac409. doi: [10.1093/bib/bbac409](https://doi.org/10.1093/bib/bbac409).
- 441 24. Wu, C.; Lin, W.; Zhang, X.; Zhang, Y.; Wang, Y.; Xie, W. PMC-LLaMA: Towards Building Open-source Language Models for
442 Medicine. Preprint, 2023. <https://arxiv.org/abs/2304.14454>.
- 443 25. Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; Sutskever, I. Language Models are Unsupervised Multitask Learners.
444 Technical Report, 2018. <https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf>.
- 445 26. Touvron, H.; Lavril, T.; Izacard, G.; et al. LLaMA: Open and Efficient Foundation Language Models. Preprint, 2023.
446 <https://arxiv.org/abs/2302.13971>.
- 447 27. Touvron, H.; Martin, L.; Stone, K.; et al. Llama 2: Open Foundation and Fine-Tuned Chat Models. Preprint, 2023.
448 <https://arxiv.org/abs/2307.09288>.
- 449 28. Gao, L.; Biderman, S.; Black, S.; et al. The Pile: An 800GB Dataset of Diverse Text for Language Modeling. Preprint, 2021.
450 <http://arxiv.org/abs/2101.00027>.
- 451 29. Lo, K.; Wang, L.L.; Neumann, M.; Kinney, R.; Weld, D. S2ORC: The Semantic Scholar Open Research Corpus. In Proceedings of
452 the 58th Annual Meeting of the Association for Computational Linguistics (*ACL 2020*), 4969–4983. Online, 5th-10th July 2020.
453 doi: [10.18653/v1/2020.acl-main.447](https://doi.org/10.18653/v1/2020.acl-main.447).
- 454 30. Loshchilov, I., Hutter, F. Decoupled Weight Decay Regularization. In Proceedings of the Seventh International Conference on
455 Learning Representations (*ICLR'2019*). New Orleans, USA, 6th-9th May 2019. <https://openreview.net/pdf?id=Bkg6RiCqY7>.
- 456 31. Hu, E.J.; Shen, Y.; Wallis, P.; et al. LoRA: Low-Rank Adaptation of Large Language Models. In Proceedings of the Ninth In-
457 ternational Conference on Learning Representations (*ICLR'2021*). Online, 3rd-7th May 2021.
458 <https://openreview.net/forum?id=nZeVKeeFYf9>.
- 459 32. Aghajanyan, A.; Zettlemoyer, L.; Gupta, S. Intrinsic Dimensionality Explains the Effectiveness of Language Model Fi-
460 ne-Tuning. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th Interna-
461 tional Joint Conference on Natural Language Processing (*ACL-IJCNLP'2020*), 7319–7328. Online, 1st-6th August 2020. Doi:
462 [10.18653/v1/2021.acl-long.568](https://doi.org/10.18653/v1/2021.acl-long.568).
- 463 33. He, J.; Cheng, X. Weighting Methods for Rare Event Identification From Imbalanced Datasets. *Front Big Data* **2021**, Volume 4,
464 Article 715320. doi: [10.3389/fdata.2021.715320](https://doi.org/10.3389/fdata.2021.715320).
- 465 34. Singh, K. How to Improve Class Imbalance using Class Weights in Machine Learning? Analytics Vidhya. Published October 6,
466 2020. <https://www.analyticsvidhya.com/blog/2020/10/improve-class-imbalance-class-weights>. Accessed on 29th January 2024.
- 467 35. Chawla, N.V.; Bowyer, K.W.; Hall, L.O.; Kegelmeyer, W.P. SMOTE: Synthetic Minority Over-sampling Technique. *J Artif Intell*
468 *Res* **2002**, 16, 321–357. doi: [10.1613/jair.953](https://doi.org/10.1613/jair.953).
- 469 36. Lin, T.Y.; Goyal, P.; Girshick, R.; He, K.; Dollar, P. Focal Loss for Dense Object Detection. *IEEE Trans Pattern Anal Mach Intell*
470 **2020**, 42(2), 318–327. doi: [10.1109/TPAMI.2018.2858826](https://doi.org/10.1109/TPAMI.2018.2858826).
- 471

472 **Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual
473 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
474 to people or property resulting from any ideas, methods, instructions or products referred to in the content.