

1 **Automated Derivation of Diagnostic Criteria for Lung Cancer using Natural Language Processing on**  
2 **Electronic Health Records: A pilot study.**

3 Andrew Houston<sup>1,2</sup>, Sophie Williams<sup>1,2</sup>, William Ricketts<sup>3</sup>, Charles Gutteridge<sup>1</sup>, Chris Tackaberry<sup>4</sup>, John  
4 Conibear<sup>5</sup>

5 <sup>1</sup> Barts Life Sciences, Barts Health NHS Trust, London, UK

6 <sup>2</sup> Digital Environment Research Institute, Queen Mary University of London, London, UK

7 <sup>3</sup> Respiratory Medicine, Barts Health NHS Trust, London, UK

8 <sup>4</sup> Clinithink Ltd., London, UK

9 <sup>5</sup> Barts Cancer Centre, Barts Health NHS Trust, London, UK

10 **Abstract**

11 **Background:** The digitisation of healthcare records has generated vast amounts of unstructured data,  
12 presenting opportunities for improvements in disease diagnosis when clinical coding falls short, such  
13 as in the recording of patient symptoms. This study presents an approach using natural language  
14 processing to extract clinical concepts from free-text which are used to automatically form diagnostic  
15 criteria for lung cancer from unstructured secondary-care data.

16 **Methods:** Patients aged 40 and above who underwent a chest x-ray (CXR) between 2016-2022 were  
17 included. ICD-10 and unstructured data were pulled from their electronic health records (EHRs) over  
18 the preceding 12 months to the CXR. The unstructured data were processed using named entity  
19 recognition to extract symptoms, which were mapped to SNOMED-CT codes. Subsumption of  
20 features up the SNOMED-CT hierarchy was used to mitigate against sparse features and a frequency-  
21 based criteria, combined with univariate logarithmic probabilities, was applied to select candidate  
22 features to take forward to the model development phase. A genetic algorithm was employed to  
23 identify the most discriminating features to form the diagnostic criteria.

24 **Results:** 75002 patients were included, with 1012 lung cancer diagnoses made within 12 months of  
25 the CXR. The best-performing model achieved an AUROC of 0.72. Results showed that an existing  
26 ‘disorder of the lung’, such as pneumonia, and a ‘cough’ increased the probability of a lung cancer  
27 diagnosis. ‘Anomalies of great vessel’, ‘disorder of the retroperitoneal compartment’ and ‘context-  
28 dependent findings’, such as pain, statistically reduced the risk of lung cancer, making other  
29 diagnoses more likely. The performance of the developed model was compared to the existing  
30 cancer risk scores, demonstrating superior performance.

31 **Conclusions:** The proposed methods demonstrated success in leveraging unstructured secondary-  
32 care data to derive diagnostic criteria for lung cancer, outperforming existing risk tools. These  
33 advancements show potential for enhancing patient care and results. However, it is essential to  
34 tackle specific limitations by integrating primary care data to ensure a more thorough and unbiased

35 development of diagnostic criteria. Moreover, the study highlights the importance of contextualising  
36 SNOMED-CT concepts into meaningful terminology that resonates with clinicians, facilitating a  
37 clearer and more tangible understanding of the criteria applied.

38 **Keywords:** Electronic Health Records; Natural Language Processing; Cancer; Diagnostics; SNOMED-  
39 CT; Machine Learning; Genetic Optimisation

## 40 **Background**

41 Lung cancer stands as one of the most common and serious types of cancer, ranking 2<sup>nd</sup> in terms of  
42 new cases and 1<sup>st</sup> in terms of mortalities, according to global statistics from 2020 [1]. The most  
43 recent statistics show that, in England, only 29.4% of lung cancer cases are identified at stages 1 and  
44 2 [2], underscoring the critical need for improved diagnostic criteria and detection methods to  
45 enhance the chances of successful treatment and reduce the burden of this disease on patients and  
46 healthcare systems. Recognising this urgency, the NHS has a long-term plan to diagnose 75% of all  
47 lung cancers at stages 1 and 2 by 2028, aiming to significantly improve early detection rates and  
48 patient outcomes.

49 Early diagnosis is imperative given the aggressive nature of lung cancer, with delays in detection  
50 resulting in patients presenting with more advanced stages of the disease. Recent data published by  
51 the Office for National Statistics and Public Health England showed the 5-year survival rate among  
52 patients diagnosed with stage 1 lung cancer was 56.6%, with this figure reducing to only 2.9% among  
53 those diagnosed with stage 4 disease [3]. Additionally, precise diagnostic criteria play a pivotal role in  
54 distinguishing lung cancer from a spectrum of cardiothoracic and respiratory conditions that may  
55 exhibit similar symptoms. With that said, to ensure the cost-effectiveness and cost-benefit of  
56 targeted interventions aimed at improving the diagnosis of lung cancer, judicious allocation of  
57 resources is required [4].

58 Electronic health records (EHRs) have revolutionized clinical research by offering a vast and  
59 comprehensive repository of patient information. Records encompass a range of data, including  
60 patient demographics, medical history, laboratory results, medication prescriptions, and procedure  
61 information. Such extensive and structured data enable researchers to conduct large-scale,  
62 population-based studies, aiding in the identification of trends [5], risk factors [6–8], and treatment  
63 outcomes [9,10]. However, a significant limitation of EHRs pertains to accuracy and completeness,  
64 particularly among symptoms and diagnosis data. Symptoms and diagnoses are often documented in

65 unstructured free-text clinical notes, requiring manual coding into clinical ontologies such as ICD-10  
66 and SNOMED-CT. The process of clinical coding can introduce inaccuracies and 'missingness' in the  
67 data, posing considerable challenges for clinical research [11,12]. Considering these challenges,  
68 techniques such as natural language processing (NLP) offer valuable solutions for not only mitigating  
69 the limitations of structured data but also unlocking valuable insights that may be exclusive to free-  
70 text narratives of patient encounters.

71 Natural Language Processing (NLP) has gained utility in extracting and analysing information in  
72 Electronic Health Records (EHRs). Koleck et al. (2019) conducted a literature review, finding 27  
73 relevant studies using NLP to analyse symptoms in EHR narratives [13]. NLP has been used for  
74 auditing discharge reports [14], predicting readmissions [15], and aiding in diagnosis [16–18].  
75 Weissman et al. (2016) used NLP to classify discharge documents based on critical illness-related  
76 keywords with high accuracy [14]. Greenwald et al. (2017) developed an NLP tool to extract  
77 readmission-related concepts and achieved comparable performance to existing prediction models  
78 [15]. In oncology, NLP extracted features from CT reports for predicting lymph node metastasis in  
79 non-small cell lung cancer (NSCLC) with competitive performance [16]. Despite its potential, there's a  
80 gap in applying NLP to oncology symptoms, highlighting an opportunity for further research [13].

81 While NLP has demonstrated its effectiveness in various healthcare applications, there is a growing  
82 recognition of the advantages of extracting ontological concepts rather than use-case-specific  
83 concepts [19,20]. This approach provides a more generalised framework for understanding and  
84 organising medical information, contributing to interoperability [21,22] and facilitating the linkage  
85 with already coded, structured, clinical data found in the EHR. This transition to ontological concept  
86 extraction aligns with the broader adoption of standardised medical terminologies like SNOMED CT,  
87 which play an important role in structuring and organizing clinical data for improved healthcare  
88 decision-making and research.

89 From a machine learning perspective, extraction of concepts from a hierarchical ontology offers a  
90 crucial advantage, enabling the retention of valuable information, even when a patient reports rarer  
91 or more specific symptoms. For instance, when a patient mentions a symptom like a 'chesty cough',  
92 machine learning systems can link it to a higher-level concept in the ontology, such as "cough." This  
93 hierarchical relationship allows the model to preserve the broader context and meaning of the  
94 symptom, preventing the loss of nuanced information that might occur in non-hierarchical concept  
95 lists, where rare features might otherwise be removed. Failure to account for such sparsity could  
96 result in poor or unreliable classification performance [23,24].

97 Given the promise of NLP for the accurate extraction of relevant features, at scale, this study applies  
98 NLP to extract SNOMED-CT concepts from free-text notes, applies subsumption to elevate rarer  
99 symptoms up the ontological hierarchy, then feeds the final feature set into a machine learning  
100 framework to train a model to discriminate lung cancer from other diseases. Furthermore, this study  
101 provides an exploration into how feature weights might be affected by demographic information like  
102 age, sex and ethnicity.

103

## 104 **Methodology**

### 105 *Eligibility*

106 Data were extracted from the Barts Health NHS Trust Data Warehouse for all patients meeting the  
107 following eligibility criteria: Patients referred for a chest x-ray (CXR), aged 40 years or older at the  
108 point of referral, during two time periods between 01 Jan 2016 and 31 Dec 2019 or 01 Jan 2022 and  
109 31 Dec 2022 were eligible for inclusion. The time window of 01 Jan 2020 - 31 Dec 2021 were not  
110 considered due to deviations from the typical cancer care-pathways as a result of the COVID-19  
111 pandemic. Patients who had opted out of their data being used for research, those without medical  
112 notes beyond four years from the original x-ray, unless a second confirmatory x-ray within four years

113 ruled out lung cancer, and patients with an existing or historical diagnosis of any cancer were  
114 excluded from participation in the study.

#### 115 *Data Sources*

116 All free-text data contained in the secondary care EHR system, from one year prior to the date of the  
117 first chest X-ray, were extracted and combined with demographic information, including Age, Sex and  
118 Ethnicity, and ICD-10 data from the same time period. Additionally, diagnostic data in the form of  
119 ICD-10 codes and the Somerset Cancer Registry were extracted for the subsequent four years post-  
120 CXR, or up to the maximum available timepoint.

121 To determine the ground truth, a patient was labelled as having lung cancer if a diagnosis was  
122 recorded in the Somerset Cancer Registry, or an ICD-10 code of C34 (Malignant neoplasm of  
123 bronchus and lung) was present in the patient's EHR post-CXR. Considering the potential delays in  
124 diagnoses, post-CXR, and the delays in uploading this information onto the electronic health records  
125 system, model training was performed iteratively, each time re-labelling the ground truth to consider  
126 an additional month of diagnoses. The iterative process was performed first considering only  
127 patients diagnosed with lung cancer within the subsequent month following their CXR, continuing to  
128 add more patients until 12-months post-CXR. Instances of lung cancer diagnoses over time the  
129 respective model performance is presented in Fig. 1a.

#### 130 *Feature Extraction*

131 To extract structured information from the free-text data, named entity recognition (NER) was  
132 performed using the NLP software, CLIX (Clinithink Ltd., London, UK). The free-text was queried  
133 against two resource sets, a 'Core-Problems' list containing common clinical symptoms and  
134 diagnoses, and the Human Phenotype Ontology. The top 100 clinical features for each resource set  
135 are presented in the supplementary file.

#### 136 *Feature Engineering*

137 To handle missing data, sex and ethnicity were imputed using the most common category. Symptom  
138 data were binary, and an assumption was made that if a diagnosis or symptom was not found in

139 either the structured ICD-10 data or identified by the NLP algorithm, the patient did not have the  
140 diagnosis or symptom.

141 To ensure a harmonised dataset, all features were mapped to the SNOMED-CT ontology. To address  
142 the sparseness of features in the lower levels of the SNOMED-CT hierarchy, we employed a  
143 subsumption process to generate and maintain features at higher levels of the hierarchy, ensuring  
144 the inclusion of all subordinate features.

#### 145 *Feature Selection*

146 Given the high dimensionality, with the number of symptom features exceeding 12,000, the dataset  
147 could not be analysed statistically. Instead, a genetic approach was taken. First, symptom features  
148 were removed where less than 0.5% of all patients or less than 5% of lung cancer patients had the  
149 symptom documented in their notes. Thereafter, the remaining features were ranked according to  
150 their Bayesian importance value, calculated as:

$$\text{IMPT}_{\text{NB}} = |\log(p(x_i = 1|y_j = 1)) - \log(p(x_i = 1|y_j = 0))|$$

151 where  $x_i$  and  $y_j$  are 1-dimensional binary arrays indicating the presence of feature  $i$  and diagnosis  $j$   
152 for each patient.

153 Following the ranking of all features, starting from the lowest ranking symptom, symptoms were  
154 removed should they have a Jaccard coefficient greater than 0.8. Thereafter, the remaining  
155 symptoms were input into a tabu asexual genetic algorithm (TAGA) [25], configured to select the  
156 feature set which maximises the area under the receiver operating characteristic curve (AUROC).  
157 TAGA was tasked with returning  $\lambda$  features, where  $\lambda$  is a number between 5 and 20. The rationale for  
158 capping the number of features included in a model at 20 was to ensure the interpretability of the  
159 final diagnostic criteria and to prevent overfitting.

#### 160 *Model Development and Evaluation*



161 20% of the data was held out of the model development process and used as a test set, with the  
162 remaining 80% being used for training and validation. For each model, a 5-fold cross-validation  
163 process was applied to select the most relevant features and identify the appropriate  
164 hyperparameters for the model, following which performance was examined using the test set. The  
165 performance of the trained models was assessed using the following diagnostic test characteristics;  
166 accuracy, sensitivity, and specificity, in addition to the calculation of the AUROC, with AUROC acting  
167 as the primary evaluation measure.

168 This study considered the following classification models: Logistic Regression, Mixed Naïve Bayes,  
169 and Decision Trees. The rationale for the selection of these models lies in their interpretability and  
170 ease of application.

#### 171 *Comparison with existing risk tools*

172 To determine whether the proposed method improves the diagnosis of lung cancer beyond that of  
173 existing methods that make use of similar features, a comparison with existing risk tools was  
174 performed. The proposed method was compared against the lung cancer component of the QCancer  
175 score [26,27] and the lung cancer-related risk assessment tools listed on the Cancer Research UK  
176 website [28].

177 In applying the QCancer score, the publicly available weights were used, and the score calculated on  
178 the same test set used for all previous comparisons. The risk assessment tools (RATs) of Hamilton et  
179 al. (2005) [28] are solely a set of feature combinations and their associated positive predictive values.  
180 Therefore, to apply the RATs to the data used in this body of work, a logistic regression model was  
181 trained for each feature combination, using the training set used for all previous experiments,  
182 returning the probability of lung cancer for each patient in the test set. Thereafter, the highest  
183 probability of all feature combinations was regarded as the final prediction for each patient. As  
184 before, the AUROC was used to compare each model.

185

## 186 **Results**

### 187 *Demographic Information*

188 In total, 75002 patients (35628 female) were included in this study. The study population had a mean  
189 age of 63 years  $\pm$  14 years. 36123 identified as 'White', 20219 identified as 'Asian', 7851 identified as  
190 'Black', 3330 identified as 'Other' and 835 identified as 'Mixed Ethnicity'. Two and 6644 patients were  
191 missing sex and ethnicity data, respectively, which were imputed.

192 In total, over the 12-month observation period after the first CXR, a total of 1012 lung cancer  
193 diagnoses were made. The occurrence of lung cancer at each monthly increment are shown in Fig.  
194 1a. Also, plotted are the number of diagnoses made following a repeat scan. The total number of  
195 diagnoses following the first scan plateaued four months post-CXR, with additional diagnoses after  
196 which time being made only after a further CXR. Aside from lung cancer, other common respiratory  
197 diagnoses in the dataset included: COPD (n=1883), atelectasis (n=2432) and pneumonia (n=398).

### 198 *Risk-Score Performance Characteristics*

199 Fig. 1b shows the performance of each of the three models, in terms of AUROC, across all 12 time  
200 intervals. The performance of the logistic regression model significantly outperformed the other two  
201 models, in terms of absolute performance but also model stability, denoted by the reduced standard  
202 deviation of AUROC. Of note, the performance of all models was less stable in the first five months,  
203 highlighting the likelihood of poorer class labelling resulting from a delay in diagnoses being  
204 uploaded to the EHRs. Considering the stabilisation in performance at five months, coupled with the  
205 plateau in diagnoses without additional scans, to strike the balance between the highest quality  
206 labelling and stable model performance, the ground truth labels established at 5-months were used  
207 for all future experiments.

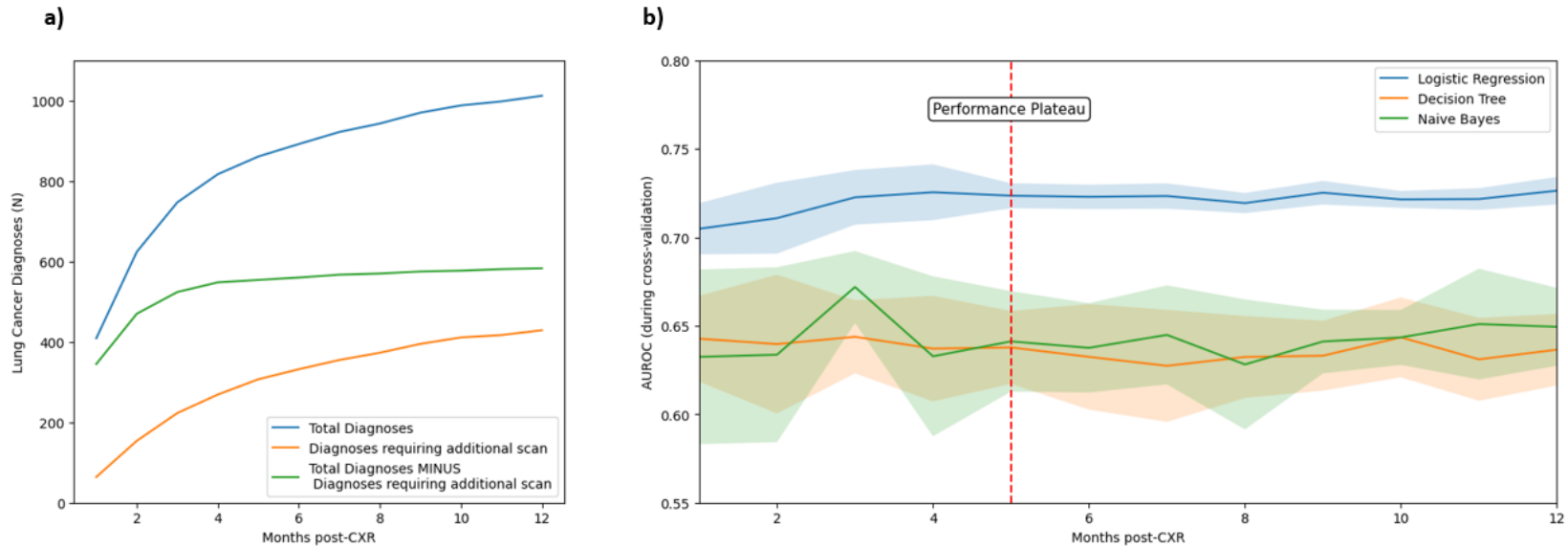


Figure 1: a) The number of diagnoses occurring at each monthly interval, for the subsequent 12 months post-CXR. b) The mean AUROC of the three tested models across each monthly interval, demonstrating the performance stabilisation and plateau from month five onwards. The shaded area indicates the standard deviation of the AUROC across the cross-validation folds.

209 *Influence of Age, Sex and Ethnicity on Risk-Score Performance*

210 Table 1 shows the performance of the model solely using the symptoms found in the EHR of the  
 211 patient, then with the inclusion of demographic data. The inclusion of age and ethnicity was shown  
 212 to improve the diagnostic performance of the model, increasing AUROC to 0.69 and 0.67,  
 213 respectively. Gender did not improve model performance in isolation. The inclusion of age, gender  
 214 and ethnicity improved model performance across all metrics resulting in an AUROC of 0.72, with an  
 215 associated sensitivity and specificity of 0.69 and 0.67, respectively.

Table 1: Performance characteristics of the logistic regression model on the test set, when each combination of the demographic features is incorporated. Values in brackets indicate the mean and standard deviation of the cross-validation performed on the training set.

Input Features	Accuracy	Balanced Accuracy	AUROC	Sensitivity	Specificity
Symptoms Only	0.78 (0.77 ± 0.02)	0.59 (0.6 ± 0.01)	0.63 (0.63 ± 0.02)	0.41 (0.44 ± 0.03)	0.78 (0.77 ± 0.02)
Symptoms and Age	0.66 (0.66 ± 0.00)	0.64 (0.66 ± 0.01)	0.69 (0.71 ± 0.01)	0.62 (0.66 ± 0.03)	0.66 (0.66 ± 0.00)
Symptoms and Gender	0.78 (0.73 ± 0.08)	0.59 (0.6 ± 0.02)	0.64 (0.64 ± 0.03)	0.41 (0.46 ± 0.07)	0.78 (0.74 ± 0.08)
Symptoms and Ethnicity	0.54 (0.55 ± 0.02)	0.61 (0.61 ± 0.02)	0.67 (0.66 ± 0.02)	0.69 (0.68 ± 0.06)	0.54 (0.55 ± 0.02)
Symptoms, Age and Gender	0.66 (0.66 ± 0.00)	0.66 (0.67 ± 0.01)	0.7 (0.72 ± 0.01)	0.66 (0.67 ± 0.02)	0.66 (0.66 ± 0.00)
Symptoms, Age and Ethnicity	0.66 (0.66 ± 0.00)	0.67 (0.66 ± 0.01)	0.71 (0.72 ± 0.01)	0.68 (0.67 ± 0.02)	0.66 (0.66 ± 0.00)
Symptoms, Gender and Ethnicity	0.66 (0.63 ± 0.04)	0.6 (0.61 ± 0.02)	0.68 (0.67 ± 0.02)	0.54 (0.59 ± 0.04)	0.66 (0.63 ± 0.04)
Symptoms, Age, Gender and Ethnicity	0.66 (0.66 ± 0.00)	0.67 (0.67 ± 0.02)	0.72 (0.72 ± 0.01)	0.69 (0.69 ± 0.03)	0.66 (0.66 ± 0.00)

216

## 217 *Feature Importance*

218 To understand how each predictor influences the prediction of lung cancer, SHAP (Shapley Additive  
219 Explanations) values were calculated (Fig. 2). The most influential feature was age, with older  
220 individuals exhibiting a significant increase in the model's output towards predicting lung cancer.  
221 Additionally, the presence of an 'existing disorder of the lung' was found to positively impact the  
222 prediction. Notably, individuals of white ethnicity had the greatest influence on the model outputs,  
223 increasing the SHAP value towards the prediction of lung cancer, although all ethnicities displayed  
224 varying degrees of impact toward a positive diagnosis. Males had an increased SHAP value,  
225 contributing to the prediction. Conversely, the presence of a 'congenital anomaly of a great vessel'  
226 and 'disorders of the retroperitoneal compartment' reduced the SHAP value. Context-dependent  
227 factors, such as pain, bleeding, and arthropathy, also reduced the SHAP value, making a prediction of  
228 lung cancer less likely. Finally, the presence of a cough was found to increase the SHAP value, further  
229 emphasising its relevance in the prediction of lung cancer.

230 Given the high-level nature of several the features, due to the subsumption process applied, an  
231 exploration into what symptoms or co-morbidities comprised such features was performed. Fig. 3  
232 shows each of the selected features, and some of the most prominent features which comprise  
233 them.

## 234 *Comparison of the proposed approach with other cancer risk tools*

235 Fig. 4 shows the receiver operating characteristic curve of the model produced using the methods  
236 described in this paper, the QCancer score [26,27], and the lung cancer related risk assessment tools  
237 listed on the Cancer Research UK website [28]. As previously reported, the proposed methods  
238 resulted in an AUROC of 0.72. The application of the QCancer calculator to the test set used  
239 throughout this paper resulted in an AUROC of 0.67 and the methods of Hamilton et al. (2005)  
240 achieved an AUROC of 0.55.

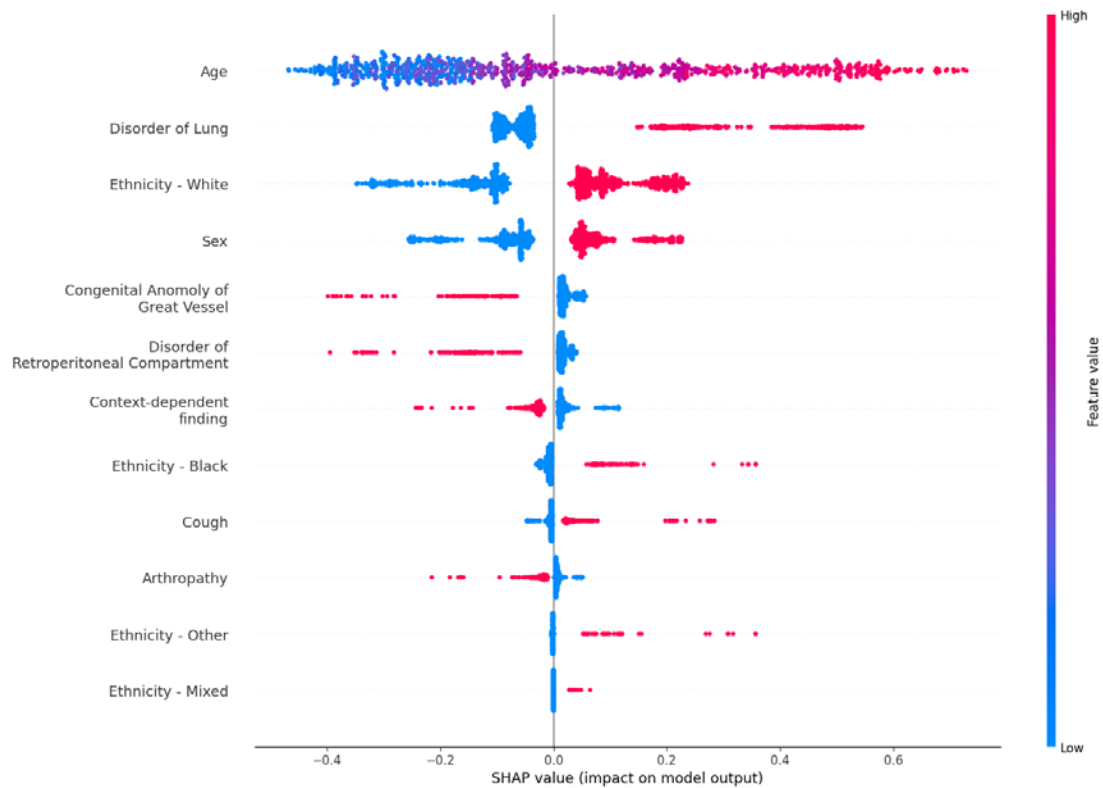


Figure 2: A summary plot of the SHAP values denoting the impact of each feature in the best performing model, on the prediction of lung cancer. Shading of each datapoint indicates the value of the feature. For all binary features, except age, a red value denotes a “true” value and blue denotes a false value. For age, the bluer a datapoint reflect a younger age, and the redder a data point, the older the patient.

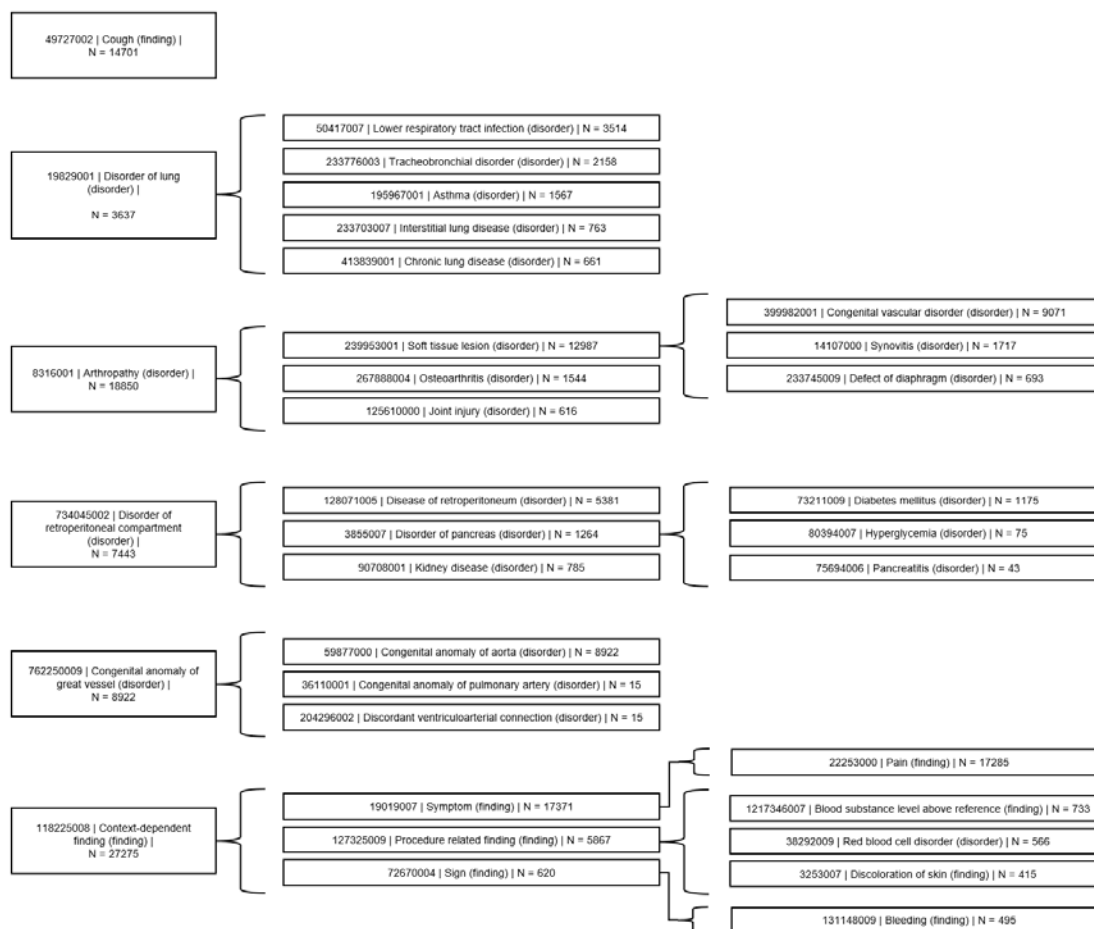


Figure 3: Visualisation of common concepts subsumed into higher level concepts in the SNOMED-CT hierarchy.

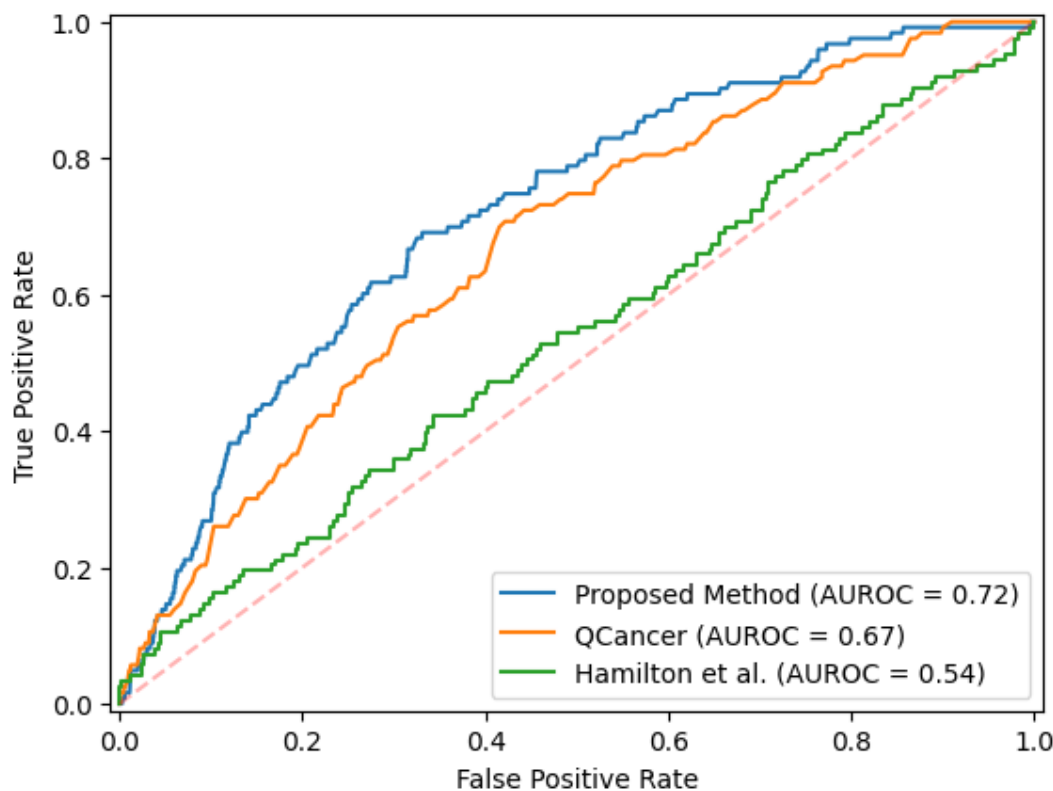


Figure 4: Receiver operating characteristic curves for the proposed method, QCancer and methods of Hamilton et al. (2005), when applied to the test set used in this study.

241

## 242 Discussion

243 This work aimed to explore the use of NLP for the extraction of SNOMED-CT concepts from  
244 unstructured clinical free-text, coupled with subsumption techniques to address the challenges  
245 posed by sparse features in high-dimensional datasets. Leveraging genetic optimisation and machine  
246 learning, the generated dataset was used to develop a predictive model for lung cancer diagnosis.  
247 Model development resulted in a classifier with stable performance characterised by low standard  
248 deviations between the cross-validation folds and an AUROC of 0.72. Additionally, the model offers a  
249 balanced trade-off between sensitivity and specificity with values of 0.69 and 0.66, respectively.  
250 Notably, our proposed methodology outperforms both the QCancer calculator [26,27] and the  
251 methods introduced by Hamilton et al. (2005) [28], highlighting the promise NLP and machine



252 learning approaches could have for the curation of rich datasets and the development of robust  
253 predictive models in the field of lung cancer risk assessment.

254 The incorporation of subsumption techniques helped mitigate the challenges posed by sparse  
255 features within our predictive model. By hierarchically organising and abstracting SNOMED-CT  
256 concepts, subsumption allowed us to identify broader, higher-level categories that encapsulate a  
257 range of related clinical terms. This not only alleviated the risk of overfitting and unreliable  
258 performance, a common concern in models trained on sparse data [23,24], but enhanced the  
259 generalisability of our model. However, the introduction of more abstract, top-level features meant  
260 that the final model was rooted in a level of granularity less commonly used in routine clinical  
261 practice. This has important implications for the practical translation and messaging of the model,  
262 highlighting the need for a clear and effective strategy to bridge the gap between the model's  
263 output, which operates at a higher conceptual level, and the clinical realities on the ground, which  
264 makes use of specific and well-established terminology.

265 The primary function of our model is to evaluate the likelihood of a positive lung cancer diagnosis  
266 when a patient enters the clinical pathway for this purpose. While this is a valuable step in enhancing  
267 early diagnosis and intervention, the success of a diagnostic tool is often measured by its ability to  
268 identify patients even before they enter the diagnostic pathway [29], ultimately achieving a  
269 significant stage shift in the diagnostic process which is associated with improved mortality rates  
270 [30]. The primary limitation of this study is its reliance on secondary care data, which did not provide  
271 sufficient longitudinal information to facilitate such an analysis. It is essential to recognise that most  
272 patient interactions with the healthcare system before a lung cancer diagnosis occur in primary care  
273 facilities, where symptoms are first reported and initial evaluations are made [31–34]. The absence  
274 of primary care data in our study thus limits the real-world applicability of the developed methods  
275 and highlights the need for future efforts to incorporate primary care data to truly impact early  
276 detection and diagnosis in clinical practice.

277 A core limitation relates to documentation bias. Although purely data-driven methods were  
278 employed to derive the features predictive of lung cancer, most patients had only one document  
279 before their CXR, a referral letter. Therefore, we must consider the possibility that the referring  
280 clinician may only include symptoms that they perceive to be relevant to the suspected diagnosis for  
281 which the scan is required, omitting other symptoms which may prove predictive. Such a limitation  
282 will often be present in such predictive modelling studies. However, if each patient were to have  
283 more clinical notes before the suspecting of lung cancer the effect of such bias may be reduced.

284 Clinically, the absence of staging data restricts our insight into the model's capacity to identify lung  
285 cancer at an early stage, which is crucial for understanding the impact of the predictions on patient  
286 outcomes. Additionally, the NER methods employed were not trained to extract genetic variants  
287 from pathology reports, specifically lung-cancer specific risk loci, which could further improve the  
288 performance of the model [35]. Future studies with access to more comprehensive and longitudinal  
289 patient data, including primary care information, genomic data and staging details, could help  
290 address these limitations and further enhance the efficacy and generalisability of the developed  
291 predictive model.

## 292 **Conclusions**

293 This research highlights the potential of combining natural language processing and machine  
294 learning techniques to enhance diagnostic criteria for lung cancer using unstructured healthcare  
295 data. The study's key findings include the successful identification of discriminating features  
296 associated with lung cancer diagnosis and achieving promising AUROC scores which outperform  
297 other comparable risk assessment tools. Such advancements hold promise for improving patient care  
298 and outcomes, albeit with a need to address certain limitations through the incorporation of primary  
299 care data for more comprehensive and unbiased criteria development.

300

301 **List of Abbreviations**

302 AUROC - Area Under the Receiver Operating Characteristic curve

303 COPD - Chronic Obstructive Pulmonary Disease

304 CT - Computed Tomography

305 CXR - Chest X-ray

306 EHR - Electronic Health Record

307 ICD-10 - International Classification of Diseases

308 NER - Named Entity Recognition

309 NHS - National Health Service

310 NLP - Natural Language Processing

311 NSCLC - Non-Small Cell Lung Cancer

312 SHAP - Shapley Additive Explanations

313 SNOMED-CT - Systematized Nomenclature of Medicine - Clinical Terms

314 TAGA - Tabu Asexual Genetic Algorithm

315 **Declarations**

316 *Ethics Approval and Consent to Participate*

317 This study was submitted to an NHS Research Ethics Committee, with subsequent approval being  
318 granted by the NHS Health Research Authority (IRAS ID: 320934). The requirement for informed  
319 consent to participate was waived by the NHS Health Research Authority following support, granted  
320 by the NHS Confidentiality Advisory Group, under Section 251 of the National Health Service Act  
321 2006. The study also adhered to the NHS National Data Opt-Out to respect patients' decision to opt  
322 out of the use of their data for research purposes.

323

324 *Consent for Publication*

325 Not applicable

326

327 *Availability of data and materials*

328 The datasets generated during and/or analysed during the current study are not publicly available  
329 due their identifiable, sensitive, and confidential nature. Data are however available from the  
330 authors upon reasonable request and with permission of Barts Health NHS Trust Information  
331 Governance Team and with appropriate approvals in place from the NHS Confidentiality Advisory  
332 Group.

333

334 *Competing Interests*

335 CT is the CEO of Clinithink Ltd., the entity that owns the NLP software employed for extracting clinical  
336 features from the free-text in this study. The remaining authors declare that they have no competing  
337 interests.

338

339 *Funding*

340 This work is sponsored by AstraZeneca UK Ltd. The funding body was independent of the study team  
341 and was not involved in the design of the study or collection, analysis, and interpretation of data or  
342 in writing the manuscript.

343

344 *Authors' contributions*

345 **AH:** Data curation, Software, Validation, Formal analysis, Project administration, Visualization,  
346 Writing – original draft; **SW:** Supervision, Project administration, Writing – original draft, Resources;  
347 **WR:** Writing – original draft; **CG:** Resources; **CT:** Funding acquisition, Project administration,  
348 Resources; **JC:** Supervision, Project administration; **All authors:** Conceptualisation, Methodology,  
349 Investigation, Writing – review & editing

350

351 *Acknowledgement*

352 We wish to acknowledge the technical support team at Clinithink Ltd. for their assistance in  
353 developing the NLP methods employed in this study.

354 **References**

- 355 1 Sung H, Ferlay J, Siegel RL, *et al.* Global Cancer Statistics 2020: GLOBOCAN Estimates of  
356 Incidence and Mortality Worldwide for 36 Cancers in 185 Countries. *CA Cancer J Clin.*  
357 2021;71:209–49.
- 358 2 Case-mix Adjusted Percentage of Cancers Diagnosed at Stages 1 and 2 in England - NHS  
359 Digital. [https://digital.nhs.uk/data-and-information/publications/statistical/case-mix-](https://digital.nhs.uk/data-and-information/publications/statistical/case-mix-adjusted-percentage-of-cancers-diagnosed-at-stages-1-and-2-in-england)  
360 [adjusted-percentage-of-cancers-diagnosed-at-stages-1-and-2-in-england](https://digital.nhs.uk/data-and-information/publications/statistical/case-mix-adjusted-percentage-of-cancers-diagnosed-at-stages-1-and-2-in-england) (accessed 13  
361 September 2023)
- 362 3 Cancer survival in England - Office for National Statistics.  
363 [https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsan](https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/bulletins/cancersurvivalinengland/stageatdiagnosisandchildhoodpatientsfollowedupto2018)  
364 [ddiseases/bulletins/cancersurvivalinengland/stageatdiagnosisandchildhoodpatientsfollowedu](https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/bulletins/cancersurvivalinengland/stageatdiagnosisandchildhoodpatientsfollowedupto2018)  
365 [pto2018](https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/bulletins/cancersurvivalinengland/stageatdiagnosisandchildhoodpatientsfollowedupto2018) (accessed 13 September 2023)
- 366 4 Reduced Lung-Cancer Mortality with Low-Dose Computed Tomographic Screening. *New*  
367 *England Journal of Medicine.* 2011;365:395–409.
- 368 5 Phadke NA, del Carmen MG, Goldstein SA, *et al.* Trends in Ambulatory Electronic  
369 Consultations During the COVID-19 Pandemic. *J Gen Intern Med.* 2020;35:3117.
- 370 6 Hill E, Mehta H, Sharma S, *et al.* Risk Factors Associated with Post-Acute Sequelae of SARS-  
371 CoV-2 in an EHR Cohort: A National COVID Cohort Collaborative (N3C) Analysis as part of the  
372 NIH RECOVER program. *medRxiv.* Published Online First: 17 August 2022. doi:  
373 10.1101/2022.08.15.22278603
- 374 7 Prado MG, Kessler LG, Au MA, *et al.* Symptoms and signs of lung cancer prior to diagnosis:  
375 Comparative study using electronic health records. *medRxiv.* 2022;2022.06.01.22275657.
- 376 8 Wong A, Young AT, Liang AS, *et al.* Development and Validation of an Electronic Health  
377 Record-Based Machine Learning Model to Estimate Delirium Risk in Newly Hospitalized  
378 Patients Without Known Cognitive Impairment. *JAMA Netw Open.* 2018;1:e181018.
- 379 9 van Laar SA, Gombert-Handoko KB, Guchelaar HJ, *et al.* An Electronic Health Record Text  
380 Mining Tool to Collect Real-World Drug Treatment Outcomes: A Validation Study in Patients  
381 With Metastatic Renal Cell Carcinoma. *Clin Pharmacol Ther.* 2020;108:644–52.
- 382 10 Houston A, Cosma G, Turner P, *et al.* Predicting surgical outcomes for chronic exertional  
383 compartment syndrome using a machine learning framework with embedded trust by  
384 interrogation strategies. *Scientific Reports 2021 11:1.* 2021;11:1–15.

- 385 11 Naran S, Hudovsky A, Antscherl J, *et al.* Audit of accuracy of clinical coding in oral surgery. *Br J*  
386 *Oral Maxillofac Surg.* 2014;52:735–9.
- 387 12 Nouraei SAR, Hudovsky A, Frampton AE, *et al.* A Study of Clinical Coding Accuracy in Surgery:  
388 Implications for the Use of Administrative Big Data for Outcomes Management. *Ann Surg.*  
389 2015;261:1096–107.
- 390 13 Koleck TA, Dreisbach C, Bourne PE, *et al.* Natural language processing of symptoms  
391 documented in free-text narratives of electronic health records: a systematic review. *J Am*  
392 *Med Inform Assoc.* 2019;26:364–79.
- 393 14 Weissman GE, Harhay MO, Lugo RM, *et al.* Natural Language Processing to Assess  
394 Documentation of Features of Critical Illness in Discharge Documents of Acute Respiratory  
395 Distress Syndrome Survivors. *Ann Am Thorac Soc.* 2016;13:1538–45.
- 396 15 Greenwald JL, Cronin PR, Carballo V, *et al.* A Novel Model for Predicting Rehospitalization Risk  
397 Incorporating Physical Function, Cognitive Status, and Psychosocial Support Using Natural  
398 Language Processing. *Med Care.* 2017;55:261–6.
- 399 16 Hu D, Li S, Zhang H, *et al.* Using Natural Language Processing and Machine Learning to  
400 Preoperatively Predict Lymph Node Metastasis for Non-Small Cell Lung Cancer With Electronic  
401 Medical Records: Development and Validation Study. *JMIR Med Inform.* 2022;10. doi:  
402 10.2196/35475
- 403 17 Chase HS, Mitrani LR, Lu GG, *et al.* Early recognition of multiple sclerosis using natural  
404 language processing of the electronic health record. *BMC Med Inform Decis Mak.* 2017;17:24.
- 405 18 Zhou L, Baughman AW, Lei VJ, *et al.* Identifying Patients with Depression Using Free-text  
406 Clinical Documents. *Stud Health Technol Inform.* 2015;216:629–33.
- 407 19 Fodeh SJ, Zirkle M, Finch D, *et al.* MedCat: A framework for high level conceptualization of  
408 medical notes. *Proceedings - IEEE 13th International Conference on Data Mining Workshops,*  
409 *ICDMW 2013.* 2013;274–80.
- 410 20 Bean DM, Kraljevic Z, Shek A, *et al.* Hospital-wide natural language processing summarising  
411 the health data of 1 million patients. *PLOS Digital Health.* 2023;2:e0000218.
- 412 21 Lee D, de Keizer N, Lau F, *et al.* Literature review of SNOMED CT use. *Journal of the American*  
413 *Medical Informatics Association.* 2014;21:e11–9.

- 414 22 Benson T, Grieve G. Principles of Health Interoperability. Published Online First: 2016. doi:  
415 10.1007/978-3-319-30370-3
- 416 23 Avanzi B, Taylor G, Wang M, *et al.* Machine Learning with High-Cardinality Categorical  
417 Features in Actuarial Applications. Published Online First: 30 January 2023.
- 418 24 Ohno-Machado L, Musen MA. Learning rare categories in backpropagation. *Lecture Notes in*  
419 *Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture*  
420 *Notes in Bioinformatics)*. 1995;991:201–9.
- 421 25 Salesi S, Cosma G, Mavrovouniotis M. TAGA: Tabu Asexual Genetic Algorithm embedded in a  
422 filter/filter feature selection approach for high-dimensional data. *Inf Sci (N Y)*. 2021;565:105–  
423 27.
- 424 26 Hippisley-Cox J, Coupland C. Symptoms and risk factors to identify women with suspected  
425 cancer in primary care: derivation and validation of an algorithm. *Br J Gen Pract*. 2013;63. doi:  
426 10.3399/BJGP13X660733
- 427 27 Hippisley-Cox J, Coupland C. Symptoms and risk factors to identify men with suspected cancer  
428 in primary care: derivation and validation of an algorithm. *Br J Gen Pract*. 2013;63. doi:  
429 10.3399/BJGP13X660724
- 430 28 Hamilton W, Peters TJ, Round A, *et al.* What are the clinical features of lung cancer before the  
431 diagnosis is made? A population based case-control study. *Thorax*. 2005;60:1059–65.
- 432 29 Balata H, Quaife SL, Craig C, *et al.* Early Diagnosis and Lung Cancer Screening. *Clin Oncol*.  
433 2022;34:708–15.
- 434 30 Flores R, Patel P, Alpert N, *et al.* Association of Stage Shift and Population Mortality Among  
435 Patients With Non–Small Cell Lung Cancer. *JAMA Netw Open*. 2021;4:e2137508–e2137508.
- 436 31 Bradley SH, Kennedy MPT, Neal RD. Recognising Lung Cancer in Primary Care. *Adv Ther*.  
437 2019;36:19–30.
- 438 32 Holtedahl K, Scheel BI, Johansen ML. General practitioners' participation in cancer treatment  
439 in Norway. *Rural Remote Health*. 2018;18. doi: 10.22605/RRH4276
- 440 33 Tørring ML, Frydenberg M, Hansen RP, *et al.* Evidence of increasing mortality with longer  
441 diagnostic intervals for five common cancers: a cohort study in primary care. *Eur J Cancer*.  
442 2013;49:2187–98.



- 443 34 Ewing M, Naredi P, Nemes S, *et al.* Increased consultation frequency in primary care, a risk  
444 marker for cancer: a case-control study. *Scand J Prim Health Care*. 2016;34:204–11.
- 445 35 Timofeeva MN, Hung RJ, Rafnar T, *et al.* Influence of common genetic variation on lung  
446 cancer risk: meta-analysis of 14 900 cases and 29 485 controls. *Human molecular*  
447 *genetics*, 2012;21:4980-4995.
- 448