

1 **Identifying secondary findings in PET/CT reports in oncological cases:**
2 **A quantifying study using automated Natural Language Processing**
3 **Secondary findings in PET/CT reports: A quantifying study using**
4 **Natural Language Processing**

5

6 Julia Sekler¹, Benedikt Kämpgen², Christian Philipp Reinert¹, Andreas Daul¹, Brigitte Gückel¹, Helmut
7 Dittmann³, Christina Pfannenberg¹, Sergios Gatidis¹

8 1. Department of Radiology, Diagnostic and Interventional Radiology, University Hospital Tübingen,
9 Hoppe-Seyler-Str. 3, 72076 Tübingen, Germany

10 2. Empolis Information Management GmbH, Technologiepark, Kettelerstraße 5 – 11, Pavillion 17,
11 97222 Rimpar, Germany

12 3. Department of Radiology, Nuclear Medicine and Clinical Molecular Imaging, University Hospital,
13 Tübingen, Hoppe-Seyler-Str.3, 72076, Tübingen, Germany

14

15

16 **Corresponding Author:**

17 Julia Sekler

18 Department of Diagnostic and Interventional Radiology

19 University Hospital Tübingen

20 Hoppe-Seyler-Str.3, 72076 Tübingen, Germany

21 Phone: +49-7071-2981282; Fax: +49-7071-295392;

22 **NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

23 **Abstract**

24 **Background:** Because of their accuracy, positron emission tomography/computed tomography (
25 PET/CT) examinations are ideally suited for the identification of secondary findings but there are only
26 few quantitative studies on the frequency and number of those.

27 Most radiology reports are freehand written and thus secondary findings are not presented as
28 structured evaluable information and the effort to manually extract them reliably is a challenge. Thus
29 we report on the use of natural language processing (NLP) to identify secondary findings from PET/CT
30 conclusions.

31 **Methods:** 4,680 anonymized German PET/CT radiology conclusions of five major primary tumor
32 entities were included in this study. Using a commercially available NLP tool, secondary findings were
33 annotated in an automated approach. The performance of the algorithm in classifying primary
34 diagnoses was evaluated by statistical comparison to the ground truth as recorded in the patient
35 registry. Accuracy of automated classification of secondary findings within the written conclusions was
36 assessed in comparison to a subset of manually evaluated conclusions.

37 **Results:** The NLP method was evaluated twice. First, to detect the previously known principal
38 diagnosis, with an F1 score between 0.65 and 0.95 among 5 different principal diagnoses.

39 Second, affirmed and speculated secondary diagnoses were annotated, and the error rate of false
40 positives and false negatives was evaluated. Overall, rates of false-positive findings (1.0%-5.8%) and
41 misclassification (0%-1.1%) were low compared with the overall rate of annotated diagnoses. Error
42 rates for false-negative annotations ranged from 6.1% to 24%. More often, several secondary findings
43 were not fully captured in a conclusion. This error rate ranged from 6.8% to 45.5%.

44

45 **Conclusions:** NLP technology can be used to analyze unstructured medical data efficiently and quickly
46 from radiological conclusions, despite the complexity of human language. In the given use case,
47 secondary findings were reliably found in in PET/CT conclusions from different main diagnoses.

48 **Keywords:** NLP; PET/CT; cancer; patient management; secondary findings

49

50 **Background**

51 In order to evaluate clinically relevant questions both retrospectively and prospectively within studies
52 as well as for therapy optimization, it is often necessary to evaluate radiological reports since these
53 are important sources of clinical diagnostic information. However, manual evaluation is only possible
54 with a significant effort if a large number of reports and findings are involved (1). To extract important
55 information from freehand texts, artificial intelligence applications can be helpful. However
56 standardized artificial intelligence (AI) applications are difficult to establish, since radiological texts are
57 usually freely written and language use and vocabulary are heterogeneous. Therefore, particular AI
58 solutions are needed, such as natural language processing (NLP). These can evaluate certain questions
59 quickly, effectively and error-controlled and can be adapted to the respective problem.

60 NLP describes a subfield of AI. It is used in numerous medical applications where text data has to be
61 analyzed and human writing or speech has to be understood and interpreted. For example, this
62 includes medical chatbots, in retrospective selection of data from unstructured records (2), in research
63 queries (3) in billing and coding (4), and in studies to analyze drug safety (5, 6).

64 Regarding the use of NLP for assessment of radiological reports, different clinical question have been
65 addressed in previous studies such as the detection of suspicious findings in mammography (7),
66 identification of site-specific bone fractures (8), tumor stage NLP (9) and other specified diagnoses
67 (10, 11) (12, 13). In this context, NLP is increasingly being used for the extraction of relevant
68 information from radiology reports in clinical studies (14-16).

69 PET/CT is mainly used for detection of tumor lesions and staging of tumor spread in oncological
70 patients. Major tumor entities examined by PET/CT include Melanoma, Prostate Cancer, Lung Cancer,
71 Lymphoma and Neuroendocrine Tumors (17-21). However, not only the status of the known or
72 suspected disease is crucial for therapy and patient management but also clinically relevant secondary
73 findings such as inflammation, vascular complications and unknown secondary tumors. Incidental
74 findings are quite common (22) and can be important for definition of scan protocols, reporting
75 strategies and further management, especially in oncology.

76 The purpose of this study was thus to automatically extract information about the occurrence of
77 secondary findings by automated analysis of freehand written radiological conclusions using NLP.

78 **Methods**

79 This study was based on a PET/CT registry (04/2013 – 12/2018) (23) including 7715 scans in total. The
80 study was reviewed and approved by the local institutional review board (Ethics committee of the
81 University of Tuebingen, reference number 064/2013B01). Informed consent regarding the use of data
82 for research was obtained from all patients.

83 **PET/CT protocols**

84 All PET/CT examinations were performed on a state-of-the art clinical scanner (Biograph mCT, Siemens
85 Healthineers, Knoxville, TN). using a standardized examination protocol. Different PET tracers were
86 applied: [68Ga]-HA-DOTATATE in case of neuroendocrine tumors, [68Ga]-PSMA in case of prostate
87 cancer, [11C]-Choline in case of prostate cancer, and [18F]-FDG in all other oncological indications. All
88 CTs were acquired in full-dose technique with contrast agent where appropriate.

89 **Structure of reports**

90 Free text PET/CT reports were written in German in a clinical routine setting using a standardized
91 structure described as follows:

92 **1. Clinical Information**

93 After providing an appropriate indication for the study by the referring physicians, the primary
94 clinical questions to be answered by the PET/CT examination are documented in the reports.

95 **2. Technique**

96 This section describes how the study was generated including information on the
97 radiopharmaceutical used, the administered activity and the CT technique. Also, the axial
98 coverage of the scan was documented (e.g., “skull base to mid-thigh”). In certain cases, PET/CT
99 protocols may have included additional acquisitions such as delayed imaging.

100 **3. Previous Studies**

101 All reports included information on prior studies which are used for comparison or correlation.
102 If no previous imaging studies are available, this was also stated.

103 **4. Findings**

104 Findings were organized by anatomic region describing both PET and CT findings relevant to
105 the clinical question within each anatomic subsection. This part also included a description of
106 incidental PET and CT findings unrelated to the primary cancer being studied. The intensity of
107 radiotracer uptake was reported using both qualitative (e.g. moderate or intense) terminology
108 as well as semiquantitative measures such as the SUV.

109 **5. Conclusion**

110 All reports concluded with a summarizing evaluation of the findings answering the specific
111 clinical questions raised by the referring physician and providing a diagnosis or a brief list of
112 differential diagnoses. In addition, potentially clinically relevant secondary findings were
113 summarized in this section.

114 **NLP**

115 The annotation of diagnoses in the report sections were automatically generated using a proprietary
116 NLP tool, Empolis Knowledge Express by Empolis Information Management GmbH (Kaiserslautern,
117 Germany; <https://knowledge.express/>). The Empolis NLP system (24, 25) implements a common NLP
118 pipeline consisting of cleansing (e.g., replacement of abbreviations), contextualization (e.g. into

119 segments "clinical information", "findings", and "conclusion"), concept recognition using common
120 terminologies such as the Radiological Lexicon (RadLex) and the International Classification of Diseases
121 (ICD), and negation detection (e.g., "affirmed", "negated", and "speculated"). The NLP system uses a
122 neural language model and word embeddings trained with fastText (26) on a medical corpus of more
123 than 100.000 German radiological reports and other medical literature (457 MB of text data). The
124 language model computes for every word a 128-dimensional vector. For concept recognition, a full
125 text index and morpho-syntactic operations such as tokenization, lemmatization, part of speech
126 tagging, compounding, noun phrase extraction and sentence detection were used. The index was
127 populated with synonyms for all entities (both from terminologies and by manual extensions). For
128 negation detection, typically, a rule-based approach is used (27); however, the heterogeneity in which
129 pathological findings are affirmed, negated or speculated require a more elaborate learning approach.
130 Therefore, the NLP system uses a bidirectional recurrent neural network based on two stacked Gated
131 Recurrent Unit (GRU) layers (28) trained and validated on more than 2.000 manually labelled reports
132 with negation information using the NLP library spaCy (29). Every input was a 50-word window, the
133 output returned a negation status for each word. The validation dataset showed 0.93 accuracy. For the
134 analysis by the Empolis NLP system, no pre-processing of the annotated radiological reports was
135 necessary.

136 Findings identified by the NLP system were classified in two categories: Unconfirmed secondary
137 findings, such as those given as differential diagnoses or as suspicions, were annotated as *speculated*,
138 whereas confirmed diagnoses are annotated as *affirmed*.

139 For automated detection of the primary patient diagnosis, the *Clinical Information* field of the
140 radiology report was used, for automated detection of secondary findings identified in the PET/CT
141 examination, the *Evaluation* field was used as input to the NLP system.

142 **The Radiological Lexicon (RadLex)**

143 In order to interpret radiological findings by NPL in a standardized way, a uniform representation of
 144 the radiological terms is required. The Radiological Lexicon (RadLex) was developed to standardize
 145 radiological terms (30). RadLex consists of a uniform vocabulary of radiological terminology that is
 146 organized hierarchically so that relationships between terms are maintained (31). In RadLex
 147 terminology there are very detailed terms for anatomy, pathology and radiological diagnoses. Some of
 148 these concepts, such as the diagnosis "neuroendocrine tumors", are therefore much easier to map
 149 with the RadLex system compared to other coding systems, such as the ICD system.

150

151 **Annotation of radiological evaluations of PET/CT scans**

152 **Selection of scans**

153 A total of 4680 scans in patients with the 5 most frequent tumor entities from the registry was
 154 annotated in this study (melanoma, non-hodgkin-lymphoma (NHL), lung cancer (lung-CA), prostate
 155 cancer (prostate-CA) and neuroendocrine tumors (NET)). Only scans from patients investigated for
 156 staging in either histologically affirmed or speculated malignancy of the above-mentioned entities
 157 were allowed. Reports were anonymized to remove patient identifiers. All characteristics of chosen
 158 scans are listed in Table .

159 **Table 1** List of all scans, divided into the five tumor types with detailed characteristics.

	Melanoma	Prostate-CA	Lung-CA	NHL	NET
Total	1178	1255	928	533	786
(H. aff./ H. spec.)	(1178/0)	(1255/0)	(901/27)	(533/0)	(719/67)
Gender					
Male	687	1255	627	314	417
Female	491	0	301	219	369
Age	63 (17 – 95)	70 (44 – 88)	66 (28 – 89)	58 (6 – 87)	62 (14 – 91)

160 **H. aff. = histologically affirmed H. spec. = histologically speculated**

161 **Age is presented as mean of years and range in parenthesis.**

162

163 **Annotation of radiological conclusions**

164 **Annotation of clinical information**

165 In order to estimate the performance of the NLP system in a setting with available ground truth, the
166 primary diagnosis was annotated first. The system was supposed to find out the main or tentative
167 diagnosis which is, in most cases, noted in the clinical information.

168 Since the principal diagnoses may be indicated with different synonyms or paraphrases within the
169 clinical information, synonyms or paraphrases were introduced into the NLP system. Subsequently, the
170 F1-score, positive predictive value and sensitivity were calculated.

171 **Annotation of secondary findings**

172 Only the conclusion and not the entire report was used for annotation of the main and secondary
173 diagnoses.

174 All radiological evaluations were uploaded onto a healthcare-analytics database provided by Empolis
175 Information Management GmbH. In this database all secondary findings, that were automatically
176 annotated were presented in a structure analogous to RadLex (31) in which supersets were in turn
177 subdivided into further specific subgroups. This categorization provides a hierarchical representation
178 of diagnoses with more general supersets such as "infectious or inflammatory disease" as well as more
179 specific subgroups such as "sinusitis". Most secondary findings were categorized within these specific
180 subgroups; remaining (rare) findings among the supersets were subsumed into the more general
181 categories, such as "infectious or inflammatory disease" or "mechanical disorder" and will be referred
182 to as "others" in the following. All affirmed or speculated secondary tumors, are subsumed as a
183 separate category of supersets. These are not further divided subgroups.

184 A list of all annotated diagnoses and their division into supersets and subgroups with the corresponding
185 RadLex codes can be found in the supplementary material (S1 Fig).

186 **Assessment of algorithm performance for classification of primary diagnosis**

187 To assess algorithm performance for classification of the primary diagnosis, algorithm output derived
188 from the clinical information field was compared to the actual clinical diagnosis of each patient.
189 Accuracy, positive predictive value and sensitivity were computed.

190 **Assessment of algorithm performance for classification of secondary findings**

191 For automated classification of secondary findings, algorithm output was compared to the content of
192 the conclusion section of each radiological conclusion. To this end all findings generated by the
193 algorithm were re-evaluated by two experts in medical imaging identifying correct and false positive
194 findings. For the evaluation of false-positive findings, the number of false-positive findings was
195 counted by manual verification by two experts in medical imaging. False positive findings were divided
196 in two categories: Non-annotated finding or wrong level of uncertainty (speculated vs. affirmed).

197 All secondary findings in total were summarized and the percentage of false positives was calculated
198 as a result. The number of false positives in which affirmed and speculated are interchanged was also
199 analyzed.

200 In order to estimate the frequency of false negative findings, a random sample of 500 radiological
201 conclusions (100 per cancer entity) were manually evaluated by two experts in medical imaging
202 identifying secondary findings that were not captured by the NLP system. Subsequently, all manually
203 recorded secondary findings were matched with those found by the NLP system.

204 **Statistical analysis:**

205 To evaluate the performance of the NLP system in detecting the principal diagnosis from the clinical
206 information, we calculated the overall correlation between the proposed NLP algorithm and the gold
207 standard. Three metrics, being sensitivity, specificity, and F1-score, were used for this purpose.

208 For the evaluation of the NLP system for annotation of secondary findings, false-positive and false-
209 negative cases were counted and correlated to the total number of annotations.

210 Results

211 Quality of automated annotation

212 Classification of main diagnoses

213 The NLP system's performance was first tested regarding the classification of the primary diagnosis.
 214 The system achieved an F1-score of 0.95 for the diagnosis of melanoma, 0.65 for the diagnosis of lung-
 215 CA, 0.90 for the diagnosis of prostate-CA, and 0.90 for the principal diagnosis of NHL showing the
 216 efficacy of the NLP system for identifying primary diagnoses from clinical information. The lowest F1-
 217 score with 0.65 was achieved for lung-CA. We achieved a perfect positive predictive value in
 218 melanoma, NHL and prostate-CA demonstrating that the NLP algorithm has high precision in
 219 identifying primary diagnoses from clinical information. The best sensitivity was in melanoma with 0.91
 220 whereas we got the lowest sensitivity with 0.49 in cases with lung-CA meaning that the system was
 221 able to identify between 49% and 91% of the cases. All primary diagnoses and the number of
 222 histologically affirmed and speculated cases with the respective F1-scores of the clinical information
 223 annotation are listed in Table 2.

224 **Table 2** Results of the annotation of the primary diagnosis in the clinical information of all scans.

	RadLex Code	<i>n</i> H. aff.	<i>n</i> H. spec.	Pos. pred. value	Sensitivity	F1- Score
Lung-CA	RIDE2220 Mass or nodule	901	27	0.98	0.49	0.65
Melanoma	RID34617 Melanoma	1178	0	1	0.91	0.95
NET	RID4483 Neuroendocrine neoplasm	719	67	0.96	0.52	0.67
NHL	RID3840 Lymphoma or RID3843 non-Hodgkin lymphoma	533	0	1	0,82	0,90
Prostate- CA	RID45689 Prostate cancer	1255	0	1	0.82	0.90

225 **Pos. pred. value = Positive predictive value H. aff. = histologically affirmed H. spec. = histologically**
 226 **speculated**

227 **Pos. pred. value, sensitivity and F1-Score of the annotation of the h. aff. and spec. main diagnosis in**
 228 **the clinical information of all scans.**

229

230 **Classification and distribution of secondary findings**

231 First, all secondary findings were combined into supersets to determine their distribution. Although
232 distributions were quite similar within the main diagnoses, there were obvious differences (Figure 1).
233 In general, the rate of "mechanical disorders" was highest in all cohorts but patients with lung CA had
234 a very high rate of "mechanical disorders, comparatively." This superset included subgroups such as
235 atelectasis, thrombosis, and pleural effusion "Infectious or inflammatory disorders" such as
236 pneumonitis, diverticulitis, and sinusitis occurred most frequently in patients with melanoma.

237 **Figure 1** Distribution of affirmed supersets of secondary findings of all cohorts as identified by the NLP-
238 System.

239

240 Second, supersets were divided into more specific subgroups (SG), secondary tumors (ST) and "others",
241 and their respective numbers were determined. "Others" included all secondary findings in supersets
242 that were not specifically divided into further subgroups based on the RadLex hierarchy (Table 3).

243 **Table 3** List of all superset categories with their included diagnoses.

Superset category	Included diagnoses (=subgroups)
Body-system-specific disorder	Dissection, aneurysm, pseudoaneurysm, pneumonia, pneumothorax
Degenerative disorder	Degeneration* ¹ , necrosis, deposition, ossification, resorption
Mechanical disorder	Atelectasis, thrombosis, embolism, pleural effusion, pericardial effusion, ascites, lymphocele, hernia, obstructive uropathy, hydronephrosis, hemorrhage, gallstone, urolithiasis, thrombus
Injury	Fracture
Iatrogenic disorder	Postoperative complication
Growth disorder	Cirrhosis
Infectious or inflammatory disease	Cholangitis, cholecystitis, diverticulitis, colitis, sarcoidosis, Crohn disease, ulcerative colitis, pancreatitis, pneumonitis, abscess, sinusitis
Cyst	Cyst* ² , epidermoid, mucocele

244 **List according to the RadLex hierarchy** *¹<http://radlex.org/RID/RID5043>,
 245 *²<http://radlex.org/RID/RID3890>)

246

247 The most affirmed specific subgroups in total were found in the cohort with the principal diagnosis of
 248 lung-CA (244 SG and 53 ST). This was followed in descending order by the cohorts with melanoma (124
 249 SG and 49 ST), prostate-CA (127 SG and 37 ST), NET (93 SG and 38 ST), and the lowest number of
 250 subgroups was observed in patients with the principal diagnosis of NHL (61 SG and 18 ST). A
 251 differentiated analysis of the individual subgroups showed that this distribution occurred for almost
 252 all main diagnoses. Only “infectious or inflammatory diseases” occurred more frequently in melanoma
 253 patients than in all other. In particular, the secondary diagnosis “sinusitis” was found very often in this
 254 cohort. The greatest amount of "others" was identified in the cohort with lung CA (351). The lowest
 255 number was observed in the NET cohort (91). All results of the analysis are presented in Table 4 and
 256 Figure 2. The detailed distribution of all secondary findings can be found in the supplementary material
 257 S1 Fig.

258 **Figure 2** Chart illustrating the pattern of affirmed and speculated subgroups (SG) of secondary findings
 259 and “others” (conglomerate of unspecific subgroups) as identified by the NLP system. Secondary
 260 tumors (ST) are a special part of subgroups.

261

262 **Table 4** Distribution of affirmed and speculated subgroups and “others” in radiological conclusions.

		Lung-CA	Prostate-CA	Melanoma	NET	NHL
Total SG/SM (%)	affirmed	244/53 (9%/2%)	127/37 (5%/1%)	124/49 (5%/2%)	93/38 (4%/1%)	61/18 (2%/1%)
	speculated	53/71 (2%/3%)	37/54 (1%/2%)	49/51 (2%/2%)	38/27 (1%/1%)	18/22 (1%/1%)
Total “others” (%)	affirmed	351 (13%)	227 (9%)	330 (13%)	91 (3%)	136 (5%)
	speculated	143 (5%)	50 (2%)	161 (6%)	29 (1%)	68 (3%)

263 **SG = subgroups SM = secondary malignancies**

264 **Secondary malignancies (SM) are a special subgroup. Percentage in relation of the total number of**
265 **secondary findings is shown in parentheses.**

266

267 **False positives and false negatives in secondary findings**

268 In order to classify the accuracy of the NLP-System, the number of false positives and false negatives
269 were also evaluated.

270 In cases with the main diagnosis NET the highest error rate of false positives was found. 17 out of 295
271 secondary findings were rated as false positives which results in an error rate of 5.8%. In contrast,
272 hardly any false positives were found in the cohort diagnosed with NHL. There only 1% of all secondary
273 findings were false positives and there was no incorrect assignment to affirmed or speculated. Overall,
274 the rate of false positives (1.0% - 5.8%) and incorrect assignment (0% - 1.1%) was very low compared
275 to the overall rate of annotated diagnoses. The complete calculation is listed in Table .

276 **Table 5** Calculation of false-positives.

FP	Lung-CA	Prostate-CA	Melanoma	NET	NHL
Total SF (aff./spec.)	855 (648/207)	531 (391/140)	695 (503/192)	295 (222/73)	309 (215/94)
subgroups/"others" <i>n</i>	425/430	267/264	274/421	176/119	131/178
FP incorrect assignment aff./spec. <i>n</i> (%)	3 (0.4%)	6 (1.1%)	1 (0.1%)	1 (0.3%)	0 (0.0%)
FP incorrect SF <i>n</i> (%)	14 (1.6%)	27 (5.1%)	8 (1.2%)	17 (5.8%)	3 (1.0%)

277 **FP = false-positives SF = secondary findings aff. = affirmed spec. = speculated SF = secondary findings**

278 **Results are presented as number and percentage (%) of total in parenthesis. Distinction was made**

279 **between incorrect assignment of affirmed and speculated and incorrect SF. The percentage of total**

280 **was calculated.**

281 Error rates of false-negative secondary findings calculated using a random sample of 100 cases per
 282 principal diagnosis ranged from 6.1% (NET) - 24% (prostate-CA) meaning that there were up to 24% of
 283 conclusions with a secondary diagnosis that was not found. More frequently, in conclusions with
 284 multiple secondary findings not all of them were recorded. This error rate varied from 6.8% for NHL to
 285 45.5% for NET. In part, this high number can be attributed to the fact that multiple secondary diagnoses
 286 were sometimes not found in a single conclusion. The complete calculation is listed in Table .

287 **Table 6** Calculation of false-negatives from a sample of 100 random radiological conclusions per
 288 principal diagnosis.

FN	Lung-CA	Prostate-CA	Melanoma	NET	NHL
Conclusions with SF	45	25	19	33	44
Non-annotated SF <i>n</i> (%)	5 (11.1%)	6 (24.0%)	2 (10.5%)	2 (6.1%)	8 (18.2%)
Deficit multiple SF <i>n</i> (%)	9 (20.0%)	7 (28.0%)	5 (26.3%)	15 (45.5%)	3 (6.8%)

289 **FN = false-negatives SF = secondary findings**

290 **Results are presented as number and percentage (%) of total in parenthesis. Distinction was made**
 291 **between non-annotated conclusions although containing SF and those with a deficit between the**
 292 **number of SF present and annotated. The percentage of total conclusions with SF was calculated.**

293

294 One example each of a false positive and false negative case is shown in Table .

295 **Table 7** Example of a typical case of a false positive and false negative radiological conclusion *¹,
296 respectively.

Example of a false positive (main diagnosis melanoma, incorrect annotated secondary finding “embolism”):

No evidence of metastatic lesions or recurrence in the left arch of the foot. Further increased metabolic activity of the thyroid gland, with hypothyroidism in need of substitution and TSH 7.4 mU/l consistent with Hashimoto's thyroiditis.

Example of a false negative (main diagnosis lung-CA, missing annotated secondary finding “aneurysm”):

...Increased metabolism sternal as well as small but metabolically active lymph node precarinally, compatible with reactive in post-thoracotomy condition. Long-segment infrarenal abdominal aortic aneurysm. Occlusion of the left iliac artery.

297 ***¹Translation from German**

298

299 **Discussion**

300 This study evaluated the applicability of an NLP technique for the annotation of secondary findings
301 from free-text written German radiology conclusions. Furthermore, the annotation results were
302 interpreted to discuss them in the context of five main oncological diagnoses.

303 The gold standard for the evaluation from radiology reports currently still is the manual selection of
304 information by experts. However, this is very time-consuming. Our data show that NLP technology is a
305 useful tool to efficiently extract secondary diagnoses from German freehand written radiology texts in
306 a time-saving manner. Since the clinical significance of secondary diagnoses varies considerably
307 between different patient groups, being able to extract them quickly and reliably from radiology
308 reports is important for quality management.

309 In identifying the main diagnosis within the clinical information, we achieved excellent F1-scores
310 between 0.65 and 0.95 without specific training, demonstrating the efficacy of the NLP algorithm. The
311 positive predictive value was between 0.96 and 1, indicating that all diagnoses found were correct.
312 Merely the sensitivity could be improved in cohorts with NET and lung-CA by training the NLP-System
313 (24), since currently up to 50% of the diagnoses are still hidden for the algorithm. However, the
314 complexity of the German language also plays a significant role here, as there are a large number of

315 paraphrases and synonyms in our freehand written clinical information for these two types of cancer.
316 It has already been shown in other studies that complex non-English texts from the German language
317 family can achieve very good scores in all three metrics by training the NLP-algorithm (32, 33).
318 However, achieving perfect quality is often challenging and may not be necessary for large data sets.
319 Annotation of the secondary diagnoses was done in three steps. First, all annotated secondary
320 diagnoses were grouped into supersets, then these were subdivided into subgroups and "others". As
321 a final step, the false positives and false negatives were identified.
322 Among the supersets the most affirmed and speculated secondary findings were found in patients with
323 a principal diagnosis of lung-CA. The frequency and classification of the clinical relevance of secondary
324 findings in lung-CA is very heterogeneous in the literature and ranges from 7 - 27% (34). In the
325 evaluation of this current study, the high rate of mechanical disorder was particularly striking. This
326 includes, for example, secondary findings such as atelectasis and pleural effusion, which are typically
327 more common in lung-CA and its treatment (35) than in other oncologic diseases. In all other cohorts,
328 the most common secondary findings were also found in the superset of mechanical disorder.
329 Secondary tumors are a special part of superset which was not further divided into smaller subgroups.
330 These have been affirmed second most frequently in all cohorts. Again, the number was highest in the
331 cohort with lung-CA. In a previous study, a secondary tumor was found in 12.6% of patients with a
332 primary diagnosis of lung-CA (36). Secondary malignancies are rather rare incidental findings (37), but
333 can have a significant impact on therapy if confirmed. In our study, mainly benign secondary tumors
334 like adenomas of the adrenal gland were identified as secondary tumors.
335 The largest number of "infectious or inflammatory disorders" was found in the cohort of melanoma
336 patients. A more detailed classification of this subset into subgroups shows that these are mainly cases
337 of sinusitis.
338 The lowest number of secondary findings was found in the cohort with NHL and NET. In part, the
339 number of secondary findings may be explained by the type of therapy. Since many neuroendocrine

340 tumors are treated primarily with surgery and specific drugs (38), the full-body impact and thus
341 secondary findings are comparatively less than in patients with melanoma or lung-CA.

342 Melanoma patients in contrast often receive immunotherapy, which increases the risk of infectious
343 diseases and patients with NHL are receiving immunochemotherapy, which weakens the immune
344 system (39). Patients with prostate-CA are the oldest cohort with an average age of 70 years. At this
345 age, people frequently have other concomitant diseases by nature and therefore some secondary
346 findings were also found in further studies (40). Earlier studies have shown that some secondary
347 findings can have a significant impact on therapy (41, 42). Therefore, it is very important to be able to
348 extract this information reliably and quickly in order to adapt patient management if necessary.

349 The rate of false-positives was very low. Some false positives could be prevented by training the system
350 slightly more (24). Sometimes related terms were recognized as diseases by the NLP system (e.g. lymph
351 node metastases as lymphoma) or confused (e.g. ectasia of the aorta as hydronephrosis).
352 Abbreviations and their ambiguity can also be a problem. For example, by partially interpreting the
353 abbreviation "ALL" (acute lymphoblastic leukemia) as "all" some false positives were generated .

354 The matching of the secondary findings to the concepts affirmed and speculated succeeded almost
355 without error. Any confusion occurred only due to linguistic inaccuracies or hints hidden in sentences
356 within the conclusions. Concepts in radiological reports that are interpreted differently even by
357 clinicians have already been identified in a previous study (43).

358 The rate of false negatives was actually higher than the number of false positives. This can also be
359 attributed to the lack of training on the one side. Some false negatives are due to language diversity in
360 the conclusions. Besides many synonyms, there are also many expressions in the German language
361 that have the same meaning. Some errors are due to ambiguity or false negation detection (44). In a
362 previous study (45), the number of false negatives was also higher than the number of false positives.
363 Here, language recognition errors, syntax errors, or the inability to recognize the plural of a word,
364 among others, were identified as sources of error. Many false-negative errors could be resolved by

365 standardizing radiology reports (46). Another study (9) also recognized that shorter reports lead to
366 fewer errors in NLP recognition. The higher amount of information in more detailed reports could
367 negatively affect NLP detection.

368 In summary, NLP is a useful tool for extracting clinically relevant data, such as secondary findings, from
369 radiology reports. This is important because no statistics are available yet regarding the most common
370 secondary diagnoses in patients with particular oncologic diseases. Furthermore, an NLP tool can help
371 to prevent clinicians from missing important information and to save time in the evaluation process.
372 This can also be used to extract important information from medical reports that otherwise would
373 require tedious re-reading. Since most NLP systems are specialized for English texts or certain text
374 types, they have to be trained for other applications (47). However, free-text written radiology reports
375 are in some ways also a challenge for NLP, since natural language also uses ambiguous terms that are
376 difficult to classify by an automated system, but which an expert may easily infer by understanding the
377 context. Therefore, free texts are supported by further machine learning processes in some studies
378 (48). On the other hand, even experienced investigators might misinterpreted free-language reports
379 authored by colleagues (49). Thus, there is a need for standardization. NLP technology could be helpful
380 to develop improved imaging reporting in radiology and nuclear medicine.

381

382 **Conclusion**

383 NLP technology can be used to efficiently and easily extract important data retrospectively from
384 radiology texts. Thus, NLP is a helpful tool for research and patient management. The complexity of
385 human language and the resulting difficulties for NLP technology should be considered when writing
386 the respective reports.

387

388 **Abbreviations**

389 **PET/CT:** Positron emission tomography/computed tomography

390 **NLP:** Natural language processing

391 **AI:** Artificial intelligence

392 **RadLex:** Radiological Lexicon

393 **ICD:** International Classification of Diseases

394 **NHL:** Non-hodgkin-lymphoma

395 **lung-CA:** Lung cancer

396 **prostate-CA:** Prostate cancer

397 **NET:** Neuroendocrine tumor

398 **SF:** Secondary findings

399 **SG:** Subgroups

400 **ST:** Secondary tumors

401 **FP:** False-positive

402 **FN:** False-negative

403 **Declarations**

404 **Ethics approval and consent to participate**

405 This study was approved by the Ethics committee of the University of
406 Tuebingen, reference number 064/2013B01. Written informed consents were waived due to
407 retrospective nature.

408 **Consent for publication**

409 Not applicable.

410 **Availability of data and material**

411 The datasets generated and analyzed during the current study are not publicly available due to
412 sensitive information but are available in anonymous form from the corresponding author on
413 reasonable request.

414 **Competing interests**

415 BK is an employee of Empolis Information Management GmbH (Kaiserslautern, Germany). The other
416 authors declare no conflict of interest.

417 **Funding**

418 No funding has been received for this publication.

419 **Authors' contributions**

420 SG originated the idea for the project. JS extracted the conclusions needed. AD took care of data
421 privacy issues in exporting data. BK performed the annotation and assisted in interpretation. JS and SG
422 analyzed the data. CR summarized methods for generating radiological reports. JS prepared the figures
423 and wrote the first draft of the manuscript. BG, CP, and HD co-wrote the final version of the manuscript
424 and were significantly involved in advising on clinical aspects. All authors read and approved the final
425 manuscript.

426 **Acknowledgements**

427 Not applicable

428

429 **References**

430 1. Mamlin BW, Heinze DT, McDonald CJ. Automated extraction and normalization of findings
431 from cancer-related free-text radiology reports. AMIA Annu Symp Proc. 2003:420-4.

- 432 2. Shah RF, Bini S, Vail T. Data for registry and quality review can be retrospectively collected
433 using natural language processing from unstructured charts of arthroplasty patients. *Bone Joint J.*
434 2020;102-B(7_Supple_B):99-104.
- 435 3. Libbus B, Rindflesch TC. NLP-based information extraction for managing the molecular biology
436 literature. *Proc AMIA Symp.* 2002:445-9.
- 437 4. Haug PJ, Ranum DL, Frederick PR. Computerized extraction of coded findings from free-text
438 radiologic reports. *Work in progress. Radiology.* 1990;174(2):543-8.
- 439 5. Jagannatha A, Liu F, Liu W, Yu H. Overview of the First Natural Language Processing Challenge
440 for Extracting Medication, Indication, and Adverse Drug Events from Electronic Health Record Notes
441 (MADE 1.0). *Drug Saf.* 2019;42(1):99-111.
- 442 6. Mohammadhassanzadeh H, Sketris I, Traynor R, Alexander S, Winqvist B, Stewart SA. Using
443 Natural Language Processing to Examine the Uptake, Content, and Readability of Media Coverage of a
444 Pan-Canadian Drug Safety Research Project: Cross-Sectional Observational Study. *JMIR Form Res.*
445 2020;4(1):e13296.
- 446 7. Jain NL, Friedman C. Identification of findings suspicious for breast cancer based on natural
447 language processing of mammogram reports. *Proc AMIA Annu Fall Symp.* 1997:829-33.
- 448 8. Wang Y, Mehrabi S, Sohn S, Atkinson EJ, Amin S, Liu H. Natural language processing of radiology
449 reports for identification of skeletal site-specific fractures. *BMC Med Inform Decis Mak.* 2019;19(Suppl
450 3):73.
- 451 9. Cheng LT, Zheng J, Savova GK, Erickson BJ. Discerning tumor status from unstructured MRI
452 reports--completeness of information in existing reports and utility of automated natural language
453 processing. *J Digit Imaging.* 2010;23(2):119-32.
- 454 10. Rink B, Roberts K, Harabagiu S, Scheuermann RH, Toomay S, Browning T, et al. Extracting
455 actionable findings of appendicitis from radiology reports using natural language processing. *AMIA Jt*
456 *Summits Transl Sci Proc.* 2013;2013:221.
- 457 11. Pham AD, Neveol A, Lavergne T, Yasunaga D, Clement O, Meyer G, et al. Natural language
458 processing of radiology reports for the detection of thromboembolic diseases and clinically relevant
459 incidental findings. *BMC Bioinformatics.* 2014;15:266.
- 460 12. Pons E, Braun LM, Hunink MG, Kors JA. Natural Language Processing in Radiology: A Systematic
461 Review. *Radiology.* 2016;279(2):329-43.
- 462 13. Cai T, Giannopoulos AA, Yu S, Kelil T, Ripley B, Kumamaru KK, et al. Natural Language Processing
463 Technologies in Radiology Research and Clinical Applications. *Radiographics.* 2016;36(1):176-91.
- 464 14. Kreimeyer K, Foster M, Pandey A, Arya N, Halford G, Jones SF, et al. Natural language
465 processing systems for capturing and standardizing unstructured clinical information: A systematic
466 review. *J Biomed Inform.* 2017;73:14-29.
- 467 15. Velupillai S, Mowery D, South BR, Kvist M, Dalianis H. Recent Advances in Clinical Natural
468 Language Processing in Support of Semantic Analysis. *Yearb Med Inform.* 2015;10(1):183-93.
- 469 16. Chen X, Xie H, Wang FL, Liu Z, Xu J, Hao T. A bibliometric analysis of natural language processing
470 in medical research. *BMC Med Inform Decis Mak.* 2018;18(Suppl 1):14.
- 471 17. Reinhardt MJ, Joe AY, Jaeger U, Huber A, Matthies A, Bucerius J, et al. Diagnostic performance
472 of whole body dual modality 18F-FDG PET/CT imaging for N- and M-staging of malignant melanoma:
473 experience with 250 consecutive patients. *J Clin Oncol.* 2006;24(7):1178-87.
- 474 18. Metser U, Dubebout J, Baetz T, Hodgson DC, Langer DL, MacCrostie P, et al. [(18) F]-FDG PET/CT
475 in the staging and management of indolent lymphoma: A prospective multicenter PET registry study.
476 *Cancer.* 2017;123(15):2860-6.
- 477 19. Kubota K, Matsuno S, Morioka N, Adachi S, Koizumi M, Seto H, et al. Impact of FDG-PET findings
478 on decisions regarding patient management strategies: a multicenter trial in patients with lung cancer
479 and other types of cancer. *Ann Nucl Med.* 2015;29(5):431-41.
- 480 20. Barrio M, Czernin J, Fanti S, Ambrosini V, Binse I, Du L, et al. The Impact of Somatostatin
481 Receptor-Directed PET/CT on the Management of Patients with Neuroendocrine Tumor: A Systematic
482 Review and Meta-Analysis. *J Nucl Med.* 2017;58(5):756-61.

- 483 21. Kuten J, Kesler M, Even-Sapir E. [the Role of PsmA Pet/Ct in Imaging Prostate Cancer].
484 Harefuah. 2021;160(7):455-61.
- 485 22. Lumbreras B, Donat L, Hernandez-Aguado I. Incidental findings in imaging diagnostic tests: a
486 systematic review. Br J Radiol. 2010;83(988):276-89.
- 487 23. PfannenberG C, Gueckel B, Wang L, Gatidis S, Olthof SC, Vach W, et al. Practice-based evidence
488 for the clinical benefit of PET/CT-results of the first oncologic PET/CT registry in Germany. Eur J Nucl
489 Med Mol Imaging. 2019;46(1):54-64.
- 490 24. Jungmann F, Kampgen B, Mildenerger P, Tsaui I, Jorg T, Duber C, et al. Towards data-driven
491 medical imaging using natural language processing in patients with suspected urolithiasis. Int J Med
492 Inform. 2020;137:104106.
- 493 25. Jungmann F, Kuhn S, Tsaui I, Kampgen B. [Natural language processing in radiology : Neither
494 trivial nor impossible]. Radiologe. 2019;59(9):828-32.
- 495 26. Bojanowski P, Grave E, Joulin A, Mikolov T. Enriching Word Vectors with Subword Information.
496 Transactions of the Association for Computational Linguistics. 2017;5:135-46.
- 497 27. Chapman WW, Hillert D, Velupillai S, Kvist M, Skeppstedt M, Chapman BE, et al. Extending the
498 NegEx lexicon for multiple languages. Stud Health Technol Inform. 2013;192:677-81.
- 499 28. Cho K, Merrienboer Bv, Gülçehre Ç, Bahdanau D, Bougares F, Schwenk H, et al., editors.
500 Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation.
501 Conference on Empirical Methods in Natural Language Processing; 2014.
- 502 29. Honnibal M, Montani I, Van Landeghem S, Boyd A. spacy: Industrial-strength natural language
503 processing in python. spaCy <https://spacy.io/>(accessed Jun 30, 2020). 2016.
- 504 30. Langlotz CP. RadLex: a new method for indexing online educational materials. Radiographics.
505 2006;26(6):1595-7.
- 506 31. Marwede D, Daumke P, Marko K, Lobsien D, Schulz S, Kahn T. [RadLex - German version: a
507 radiological lexicon for indexing image and report information]. Rofo. 2009;181(1):38-44.
- 508 32. Dahl FA, Rama T, Hurlen P, Brekke PH, Husby H, Gundersen T, et al. Neural classification of
509 Norwegian radiology reports: using NLP to detect findings in CT-scans of children. BMC Med Inform
510 Decis Mak. 2021;21(1):84.
- 511 33. Olthof AW, Shouche P, Fennema EM, FFA IJ, Koolstra RHC, Stirlor VMA, et al. Machine learning
512 based natural language processing of radiology reports in orthopaedic trauma. Comput Methods
513 Programs Biomed. 2021;208:106304.
- 514 34. Kucharczyk MJ, Menezes RJ, McGregor A, Paul NS, Roberts HC. Assessing the impact of
515 incidental findings in a lung cancer screening study by using low-dose computed tomography. Can
516 Assoc Radiol J. 2011;62(2):141-5.
- 517 35. Moller DS, Khalil AA, Knap MM, Hoffmann L. Adaptive radiotherapy of lung cancer patients
518 with pleural effusion or atelectasis. Radiother Oncol. 2014;110(3):517-22.
- 519 36. Faehling M, Schwenk B, Kramberg S, Fallscheer S, Leschke M, Strater J, et al. Second malignancy
520 in non-small cell lung cancer (NSCLC): prevalence and overall survival (OS) in routine clinical practice. J
521 Cancer Res Clin Oncol. 2018;144(10):2059-66.
- 522 37. Sebros R, Aparici CM, Pampaloni MH. Frequency and clinical implications of incidental new
523 primary cancers detected on true whole-body 18F-FDG PET/CT studies. Nucl Med Commun.
524 2013;34(4):333-9.
- 525 38. Herrera-Martinez AD, Hofland J, Hofland LJ, Brabander T, Eskens F, Galvez Moreno MA, et al.
526 Targeted Systemic Treatment of Neuroendocrine Tumors: Current Options and Future Perspectives.
527 Drugs. 2019;79(1):21-42.
- 528 39. Tun AM, Ansell SM. Immunotherapy in Hodgkin and non-Hodgkin lymphoma: Innate, adaptive
529 and targeted immunological strategies. Cancer Treat Rev. 2020;88:102042.
- 530 40. Elmi A, Tabatabaei S, Talab SS, Hedgire SS, Cao K, Harisinghani M. Incidental findings at initial
531 imaging workup of patients with prostate cancer: clinical significance and outcomes. AJR Am J
532 Roentgenol. 2012;199(6):1305-11.

- 533 41. Stump M, Keller JR, Mott SL, Stolmeier DA, Milhem MM, Liu V. The prevalence and significance
534 of radiographic incidental findings during initial staging of melanoma: a retrospective study. *J Eur Acad*
535 *Dermatol Venereol.* 2020;34(2):e62-e4.
- 536 42. Conrad F, Winkens T, Kaatz M, Goetze S, Freesmeyer M. Retrospective chart analysis of
537 incidental findings detected by (18) F-fluorodeoxyglucose-PET/CT in patients with cutaneous
538 malignant melanoma. *J Dtsch Dermatol Ges.* 2016;14(8):807-16.
- 539 43. Hobby JL, Tom BD, Todd C, Bearcroft PW, Dixon AK. Communication of doubt and certainty in
540 radiological reports. *Br J Radiol.* 2000;73(873):999-1001.
- 541 44. Rivera Zavala R, Martinez P. The Impact of Pretrained Language Models on Negation and
542 Speculation Detection in Cross-Lingual Medical Text: Comparative Study. *JMIR Med Inform.*
543 2020;8(12):e18953.
- 544 45. Li AY, Elliot N. Natural language processing to identify ureteric stones in radiology reports. *J*
545 *Med Imaging Radiat Oncol.* 2019;63(3):307-10.
- 546 46. Vuokko R, Makela-Bengs P, Hypponen H, Lindqvist M, Doupi P. Impacts of structuring the
547 electronic health record: Results of a systematic literature review from the perspective of secondary
548 use of patient data. *Int J Med Inform.* 2017;97:293-303.
- 549 47. Weikert T, Nestic I, Cyriac J, Bremerich J, Sauter AW, Sommer G, et al. Towards automated
550 generation of curated datasets in radiology: Application of natural language processing to unstructured
551 reports exemplified on CT for pulmonary embolism. *Eur J Radiol.* 2020;125:108862.
- 552 48. Pandey M, Xu Z, Sholle E, Maliakal G, Singh G, Fatima Z, et al. Extraction of radiographic findings
553 from unstructured thoracoabdominal computed tomography reports using convolutional neural
554 network based natural language processing. *PLoS One.* 2020;15(7):e0236827.
- 555 49. Schwartz LH, Panicek DM, Berk AR, Li Y, Hricak H. Improving communication of diagnostic
556 radiology findings through structured reporting. *Radiology.* 2011;260(1):174-81.

557

558

medRxiv preprint doi: <https://doi.org/10.1101/2022.12.02.22283043>; this version posted December 5, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

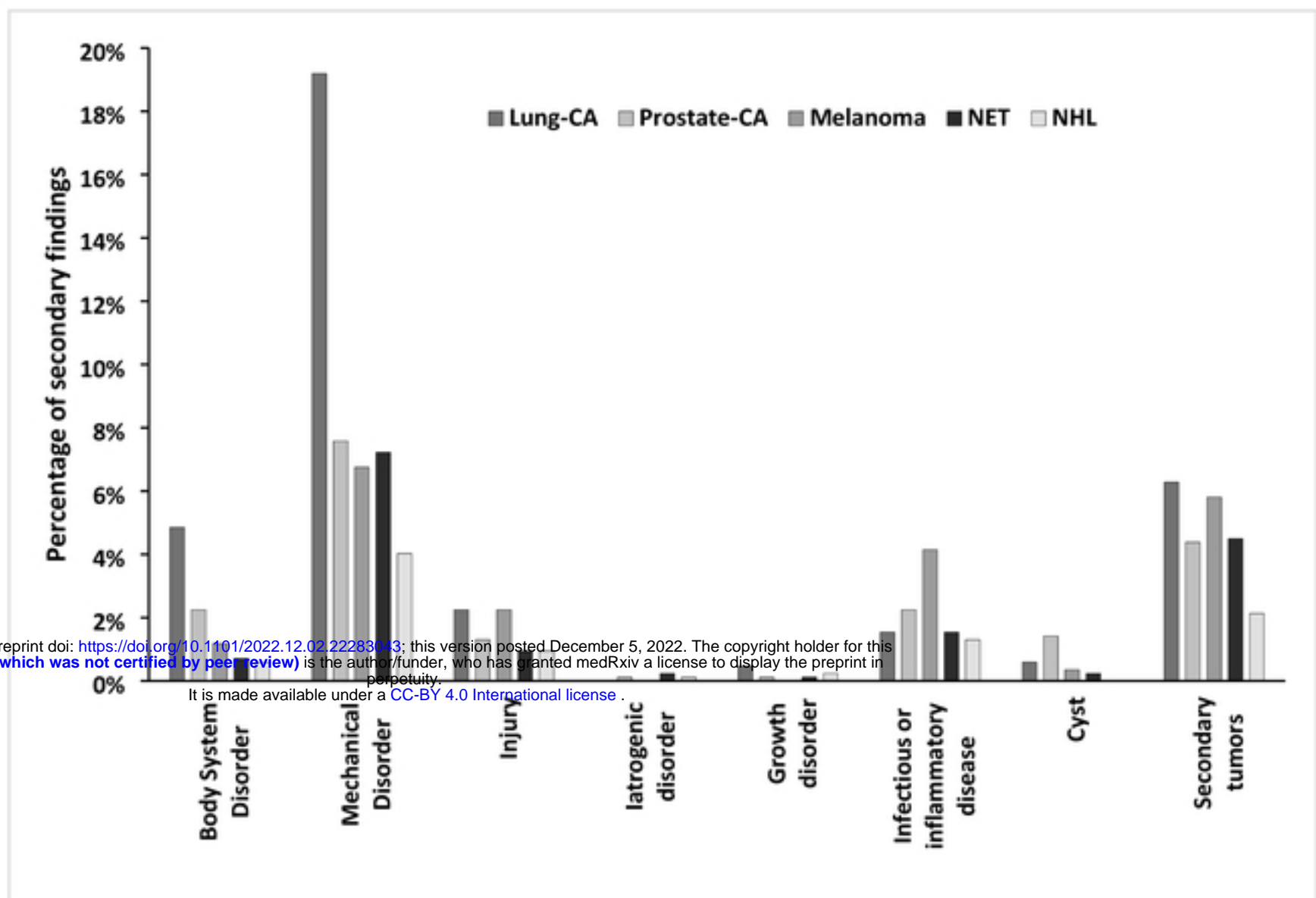


Figure 1

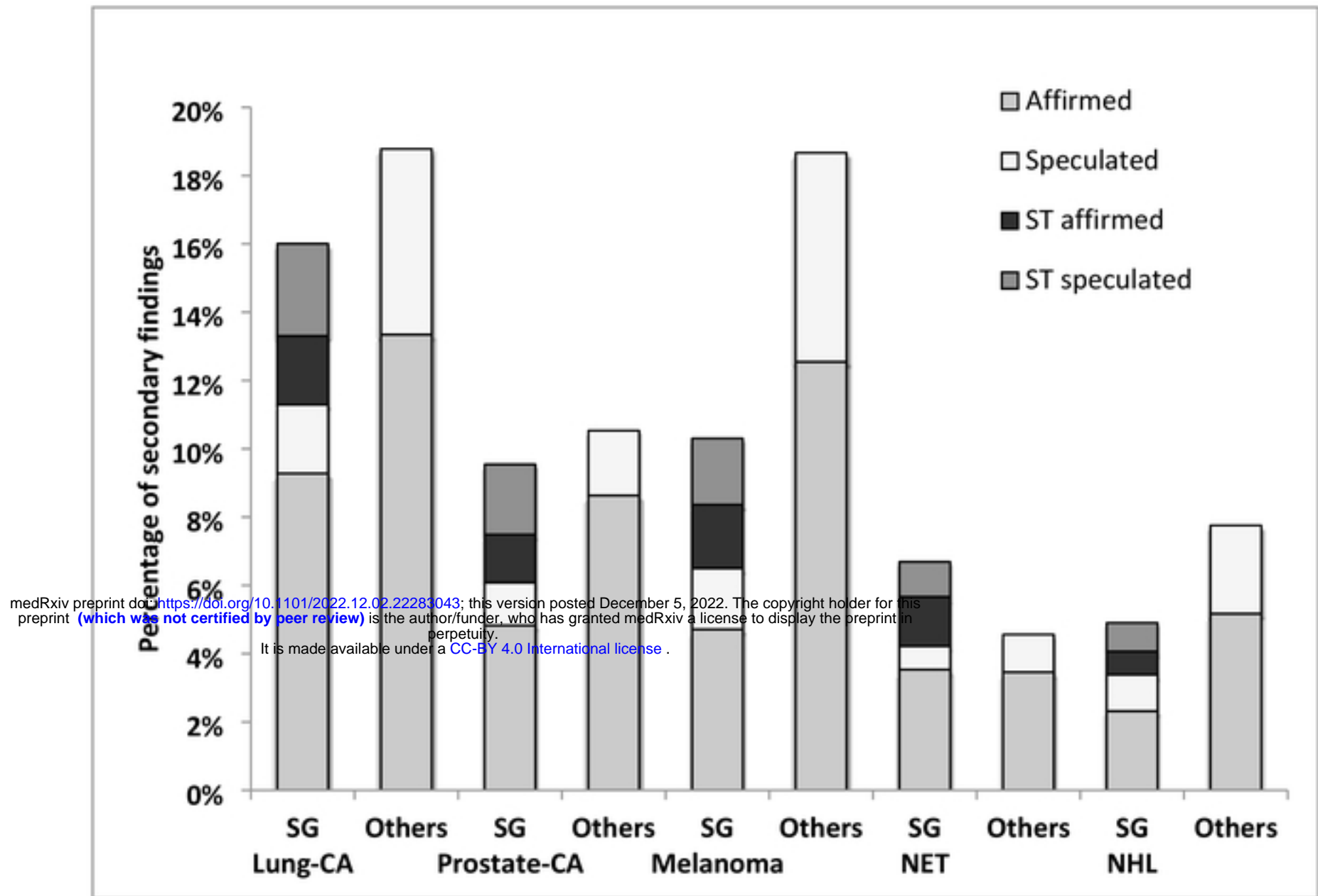


Figure 2