

1 **TITLE PAGE**

2 **Towards Medical Billing Automation: NLP for Outpatient Clinician Note**  
3 **Classification**

4  
5 **Authors:**

6 Matthew G. Crowson MD, MPA, MASc FRCSC<sup>1,2</sup>

7 Emily Alsentzer, MS, PhD<sup>3</sup>

8 Julie Fiskio, BS<sup>3</sup>

9 David W. Bates, MD, MSc<sup>3,4</sup>

10

11 1 Department of Otolaryngology-Head & Neck Surgery, Massachusetts Eye & Ear, Boston,  
12 Massachusetts, USA

13 2 Department of Otolaryngology-Head & Neck Surgery, Harvard Medical School,  
14 Massachusetts, USA

15 3 Division of General Internal Medicine and Primary Care, Brigham and Women's Hospital,  
16 Boston, MA, USA

17 4 Department of Health Policy and Management, Harvard T. H. Chan School of Public Health,  
18 Boston, MA, USA

19

20 **Word Count:** 2,964

21

22 **Keywords:** natural language processing, provider billing, level of service, outpatient  
23 care

24

25 **Manuscript Submission:**

26

27 (1) The authors indicated above have contributed to, read, and approved this  
28 manuscript.

29

30 (2) **FINANCIAL DISCLOSURE:** Dr. Bates reports grants and personal fees from EarlySense,  
31 personal fees from CDI Negev, equity from ValeraHealth, equity from Clew, equity from  
32 MDCclone, personal fees and equity from AESOP, personal fees and equity from Feelbetter,  
33 equity from Guided Clinical Solutions, and grants from IBM Watson Health, outside the  
34 submitted work. Dr. Bates has a patent pending (PHC-028564 US PCT), on intraoperative  
35 clinical decision support.

36

37 (3) **CONFLICT DISCLOSURE:** no authors have conflicts related to this manuscript.

38 (4) In consideration of the journal reviewing and editing my submission, the authors  
39 undersigned transfer, assign and otherwise convey all copyright ownership if such  
40 work is published.

41

42 (4) **AUTHOR CONTRIBUTION:** The authors confirm contribution to the paper as  
43 follows: study conception and design: MC, DB, EA, JF. Data collection: MC, JF. Analysis

44 and interpretation of results: MC, DB, EA. Draft manuscript preparation: MC, DB, EA.  
45 All authors reviewed the results and approved the final version of the manuscript.

46

47 **Corresponding Author:**

48 Matthew G. Crowson MD MPA MASc FRCSC

49 Massachusetts Eye & Ear

50 Department of Otolaryngology-Head & Neck Surgery

51 243 Charles Street

52 Boston, Massachusetts

53 02114 USA

54 [matthew\\_crowson@meei.harvard.edu](mailto:matthew_crowson@meei.harvard.edu), (p) 617-573-6559, (f) 617-573-3914

55 **ABSTRACT**

56  
57 **Objectives:** Our primary objective was to develop a natural language processing  
58 approach that accurately predicts outpatient Evaluation and Management (E/M)  
59 level of service (LoS) codes using clinicians' notes from a health system electronic  
60 health record. A secondary objective was to investigate the impact of clinic note  
61 de-identification on document classification performance.

62 **Methods:** We used retrospective outpatient office clinic notes from four medical  
63 and surgical specialties. Classification models were fine-tuned on the clinic notes  
64 datasets and stratified by subspecialty. The success criteria for the classification  
65 tasks were the classification accuracy and F1-scores on internal test data. For the  
66 secondary objective, the dataset was de-identified using Named Entity  
67 Recognition (NER) to remove protected health information (PHI), and models  
68 were retrained.

69 **Results:** The models demonstrated similar predictive performance across  
70 different specialties, except for internal medicine, which had the lowest  
71 classification accuracy across all model architectures. The models trained on the  
72 entire note corpus achieved an E/M LoS CPT code classification accuracy of  
73 74.8% (CI 95: 74.1-75.6). However, the de-identified note corpus showed a  
74 markedly lower classification accuracy of 48.2% (CI 95: 47.7-48.6) compared to  
75 the model trained on the identified notes.

76 **Conclusion:** The study demonstrates the potential of NLP-based document  
77 classifiers to accurately predict E/M LoS CPT codes using clinical notes from  
78 various medical and procedural specialties. The models' performance suggests  
79 that the classification task's complexity merits further investigation. The de-  
80 identification experiment demonstrated that de-identification may negatively  
81 impact classifier performance. Further research is needed to validate the  
82 performance of our NLP classifiers in different healthcare settings and patient  
83 populations and to investigate the potential implications of de-identification on  
84 model performance.

85

## 86 INTRODUCTION

87

88 The administrative burden of clinical billing activities for health insurance  
89 reimbursement is substantial and contributes to rising healthcare costs in the  
90 United States and other insurance-based systems worldwide <sup>1,2</sup>. Compared to  
91 other western countries, the United States' proportion of total hospital costs  
92 devoted to administrative tasks is higher--it has exceeded 25% and has been  
93 rising <sup>2</sup>. Administrative tasks related specifically to billing and insurance cost the  
94 United States healthcare system an estimated \$471 billion in 2012, comprising  
95 nearly 15% of all healthcare spending <sup>3</sup>. At the institutional level, the estimated  
96 billing and insurance-related administration costs range widely from \$20 for a  
97 primary care visit to \$215 for an inpatient surgical procedure <sup>4</sup>.

98 Several reasons for the administrative cost burden have been proposed, but  
99 the complexity of the United States' healthcare reimbursement scheme appears  
100 to be the leading cause <sup>2,5</sup>. At the individual practice level, physicians and clinic  
101 team members devote considerable time interacting with health plans daily <sup>6,7</sup>.  
102 Proposed cost mitigation efforts have ranged from reforming healthcare payment  
103 processes via standardizing payment rules, claim forms, and other innovations to  
104 computer-assisted claims processing systems that automate clinicians' receipt,  
105 validation, formatting, and sending of individual claims in an attempt to reduce  
106 the time for completion <sup>8,9</sup>.

107 Natural language processing (NLP) is a subfield of artificial intelligence  
108 devoted to understanding and analyzing human language. It has seen  
109 exponential interest in applications across all domains and subspecialties of  
110 medicine<sup>10,11</sup>. NLP has been used for a variety of tasks in medicine spanning  
111 pharmaceutical and biological knowledge discovery, adverse event detection,  
112 prognosis modeling, clinical document classification, decision support systems,  
113 and point-of-care assistive technologies such as patient-facing chatbots and triage  
114 tools<sup>10,12-17</sup>. The advantages of NLP approaches include unlocking unstructured  
115 data types in clinical narratives and incorporating alternative data sources in  
116 diagnostic or prognostic predictive models. Transformer-based NLP models  
117 have demonstrated large improvement gains over several tasks, and open-source  
118 NLP libraries have made their implementation practical. Recent work has  
119 demonstrated the utility of transformer NLP models for classifying medical  
120 entities<sup>18,19</sup>, including diagnosis codes (e.g., International Classification of  
121 Diseases (ICD) codes)<sup>20,21</sup>, as well as Current Procedural Terminology (CPT)  
122 codes from clinical pathology reports<sup>22</sup>. However, to the best of our knowledge,  
123 no work has leveraged NLP models to classify Evaluation and Management  
124 (E/M) level of service (LoS) codes. Healthcare providers use outpatient E/M  
125 LoS codes to report the complexity and intensity of patient visits for billing and  
126 reimbursement purposes. An accurate point-of-care clinical note LoS decision  
127 support tool may help increase administrative efficiency and reduce variation in

128 clinical encounter coding.

129 In this work, we developed an NLP-based classifier to predict the Evaluation  
130 and Management (E/M) level of service (LoS) codes from outpatient clinician  
131 notes. We compared the performance of several general-domain and clinical  
132 language models and assessed the impact of fine-tuning on specialty-specific  
133 notes. Furthermore, we de-identified the clinical notes and studied the impact of  
134 note de-identification on classifier importance. This study demonstrates that  
135 NLP-based document classifiers can accurately predict E/M LoS CPT codes  
136 using clinical notes from various medical and procedural specialties. Further  
137 development of this autonomous classification approach might contribute to  
138 redesigning administrative billing processes to reduce resource and cost burden.

139

## 140 **RESULTS**

141

142 ***Clinic Notes Included in Modeling.*** Individual notes lacking a ground-truth  
143 label were removed, and notes with a E/M LoS CPT ‘Level I’ code were removed  
144 before model training as these codes represented less than 1% of all records  
145 (**Table 2**). Most notes had either Level III or Level IV coding. After data  
146 preprocessing, 31,115 patient notes were available for model development across  
147 all subspecialties (Cardiology, n =13,279; Gastroenterology, n = 5,299; Internal  
148 Medicine, n =10,178; Otolaryngology-Head & Neck Surgery, n = 2,197). The E/M

149 LoS CPT ‘new patient’ codes distributions differed across specialties (**Table 2**).

150

151 **Model Performance.** The models trained on the entire note corpus achieved an  
152 E/M LoS CPT code classification accuracy of 75.6% (classification accuracy CI 95:  
153 73.0-78.3%; weighted F1-score CI-95: 0.73-0.78). The predictive performance of  
154 the models trained on specialty-specific data was similar (**Table 3; Table 4**).  
155 Internal medicine had the lowest classification accuracy (CI-95: 66.0 to 68.0%;  
156 weighted F1-score CI-95: 0.57-0.68), whereas Otolaryngology-Head & Neck  
157 surgery had the highest (classification accuracy CI-95: 82.3 to 84.6%; weighted  
158 F1-score CI-95: 0.82-0.85). Within each specialty, the models do not substantially  
159 differ in predictive performance (**Table 3; Table 4**). However, *Clinical-Longformer*  
160 generally had a higher predictive performance than the other models. The  
161 models demonstrated variation in performance across E/M LoS classes, likely  
162 due to class imbalance in the data (**Supplemental Tables S1-S5**). The models  
163 trained on the de-identified note corpus achieved an E/M LoS CPT code  
164 classification accuracy of 48.3% on average (classification accuracy CI 95: 48.0-  
165 48.5; **Table 5**).

166

## 167 **DISCUSSION**

168

169 In this study, we developed natural language processing (NLP) document

170 classifier models to assign Evaluation and Management (E/M) level of service  
171 (LoS) codes using clinical notes from an electronic health record. Our models that  
172 were trained on the entire clinical note corpus produced reasonable classification  
173 accuracy. The four models trained on the subspecialty notes showed similar  
174 predictive performance across different specialties, except for internal medicine,  
175 which had the lowest classification accuracy across all model architectures. Our  
176 results suggest that an NLP-based decision support tool for outpatient clinic  
177 notes might be a feasible approach with further development. Implementing  
178 such a tool might save healthcare system costs by reducing administrative  
179 burden and minimizing the potential for billing errors and fraud. Additionally,  
180 the tool could also alleviate the workload of physicians and clinic staff, allowing  
181 them to focus on patient care.

182 The overall classification accuracy of the model trained on the entire note  
183 corpus was good, indicating that our NLP-based approach might be useful in  
184 predicting E/M LoS CPT codes. However, we observed variability in classifier  
185 performance across models fine-tuned separately on notes from different medical  
186 specialties and across LoS CPT code levels. There could be several reasons for  
187 this observation. First, medical language is often characterized by complex  
188 terminology, abbreviations, and jargon that can be challenging for an NLP model  
189 to classify accurately. It is plausible that some medical specialties may have more  
190 complex or specialized language. Second, variability in documentation practices

191 such as differences in documentation styles, structure, and terminology between  
192 specialties, could contribute to variation across specialties. This may be due, in  
193 part, to the diverse patient populations and conditions typically encountered  
194 within each specialty. Last, clinicians within medical specialties may have  
195 idiosyncratic standards for assigning E/M LoS codes. These possibilities have  
196 support in prior work that has demonstrated considerable variability in both the  
197 electronic health record utilization<sup>23</sup> and variation in electronic health record  
198 documentation between clinicians belonging to different medical specialties and  
199 health systems<sup>24</sup>. Finally, there was class imbalance in the training data, which  
200 likely also contributed to performance variability. Taken together, the classifiers  
201 may struggle to generalize to these sources of variability, resulting in lower  
202 performance for some specialties versus others.

203  
204 Across the clinical specialties, we observed that the ‘short-form’ models (i.e.,  
205 models that accept 512 tokens) tended to underperform the *Clinical-Longformer*  
206 model in classification accuracy, which accepts inputs with a longer token length  
207 (i.e., a maximum token length of 4,096). The performance gap was not large. This  
208 finding is counterintuitive as we would expect longer notes to contain a richer  
209 representation and more informative features derived from the clinical  
210 encounter. One possible explanation is that the input sequences in clinical notes  
211 may contain redundant and irrelevant information, which might offset the

212 benefits of having longer input sequences. Additionally, the pertinent and  
213 predictive information may be more likely to be contained in the first portion of  
214 the note’s body. This portion of a standard clinical note tends to be occupied by  
215 the patient’s chief complaint, history, and physical examination.

216 Clinicians contribute to the rising administrative cost of billing and insurance  
217 activities through improper billing and coding—some of which constitute  
218 medical fraud. While the rules are sufficiently complex that errors are inevitable,  
219 prior work has shown that some clinicians and healthcare institutions actively  
220 engage in ‘upcoding’ of patients’ diagnoses or severity to receive higher  
221 payments<sup>25-28</sup>. Other healthcare entities, such as skilled nursing facilities, have  
222 also been observed to be engaged in this practice through “padding” of therapy  
223 times to increase revenues<sup>29</sup>. The substantial national economic burden of billing  
224 and insurance activities and the specific contribution of variation in coding  
225 practices represents a pressing need for improvement. Recent work has  
226 attempted to take a surveillance approach through screening for outlier behavior  
227 in coding submissions at the institutional level<sup>30</sup>. Still, we have not identified an  
228 implemented automated decision-support system that provides clinicians with  
229 coding guidance at the point of care. An NLP-based decision support tool for  
230 outpatient clinic notes might be feasible for mitigating clinician-driven  
231 administrative errors and costs associated with billing and insurance processes.  
232 However, careful attention is needed to ensure that such models are used to

233 improve billing efficiency without further assisting upcoding

234

235 De-identification of clinical data is an active area of research given the risks of  
236 inadvertent release of PHI with data sharing in clinical and research contexts<sup>31</sup>. A

237 trade-off exists in supplying sufficient informative data versus revealing

238 compromising PHI. The secondary objective of this study was to investigate the

239 impact of clinic note de-identification on document classification performance.

240 The de-identified note corpus showed a markedly lower classification accuracy

241 compared to the model trained on the identified notes. This suggests that

242 removing protected health information (PHI) elements from the clinical notes

243 significantly affects the performance of the classifier model in assigning E/M LoS

244 CPT codes. This was an unexpected finding as prior work has demonstrated that

245 de-identifying clinical notes minimally reduces information<sup>32</sup>. One possible

246 explanation for this performance drop is a loss of contextual information and

247 discriminative features. De-identification processes may inadvertently remove or

248 obscure relevant contextual information that contains predictive features. It is

249 possible that the identified data model was relying on “shortcut” features (i.e.,

250 learning to associate specific clinicians, specialty designation, or departments

251 with specific coding patterns)<sup>33</sup>. Such shortcut features could relate to the

252 prediction target (i.e, the LoS E/M code) through one or more causal paths<sup>33</sup>.

253 Similarly, the de-identification process might introduce alterations to the natural

254 flow and structure of the text that the model interprets. Another study trained on  
255 clinical notes from one emergency department setting found that the  
256 performance of word-embedding (WE)-based deep learning models did not  
257 differ when trained with identified and deidentified notes<sup>34</sup>. It is plausible that  
258 the transformer-based models used in this study may rely more heavily on  
259 contextual information. Differing approaches in de-identification may also be a  
260 factor. Various methods and algorithms are available to de-identify clinical data  
261<sup>35,36</sup>. Further research is needed to understand the potential implications of de-  
262 identification strategies on model performance across different NLP model  
263 architectures and de-identification approaches.

264  
265 This study has several limitations which should be considered. First, the  
266 study was conducted using retrospective data from a single health system, which  
267 may limit the generalizability of our findings. Further research is needed to  
268 validate the performance of our NLP classifiers in different healthcare settings  
269 and patient populations, and it should be prospectively validated in other  
270 settings. Second, our classifiers were trained on specific specialties and  
271 subspecialties, and their applicability to other medical domains remains to be  
272 investigated. Third, the performance of the NLP classifiers is influenced by the  
273 quality and structure of the clinical notes, as observed by the variability in the  
274 model performance.

275

276 An additional validation step would be to assess the model's performance  
277 against human billing auditors. Additionally, integrating the NLP-based decision  
278 support tool into electronic health record systems for seamless use by clinicians  
279 and billing staff might maximize its utility once validated. Last, estimating the  
280 economic and operational impact of the tool on healthcare costs, administrative  
281 workload, and potential reduction of billing errors and fraud may provide  
282 additional value to encourage implementation.

283

284 In conclusion, we found that NLP-based document classifiers could  
285 accurately predict E/M LoS CPT codes using clinical notes from various medical  
286 and procedural specialties. This could reduce the costs associated with this  
287 process. The models' observed accuracy suggests that the classification task's  
288 complexity merits further investigation. Our de-identification experiment  
289 demonstrated markedly lower classifier performance, suggesting that de-  
290 identification may negatively impact clinical documentation processing  
291 performance. This an important finding since de-identification is an emerging  
292 method for de-risked data sharing, and collaborative research is likely to be used  
293 more widely.

294

295

## 296 **METHODS**

297

298 This study protocol was reviewed by our institutional review board and  
299 deemed exempt from formal review (Protocol #2021P002787). The development  
300 and reporting of this predictive model were completed following published  
301 guidelines from a multidisciplinary panel on the predictive model reporting<sup>37</sup>.

302

303 ***Setting and Prediction Goal.*** We developed natural language processing  
304 (NLP) document classifiers that assign evaluation and management (E/M) level  
305 of service (LoS) Current Procedural Terminology (CPT) billing codes to  
306 outpatient clinic notes. Our NLP classifiers were trained on retrospective clinical  
307 notes from a quaternary healthcare system. The success criteria for the multi-  
308 class classification tasks were the classification accuracy and weighted F1-score  
309 on internal test data. A secondary prediction goal was determining if clinic note  
310 de-identification impacted document classification performance.

311

312 ***Dataset Development.*** To account for variations in coding practices among  
313 medicine, medicine-procedural, and surgical subspecialties, retrospective  
314 outpatient office clinic notes were selected from different medical specialties and  
315 subspecialties within Cardiology, Internal Medicine, Gastroenterology, and

316 Otolaryngology-Head & Neck Surgery. Notes were obtained from clinic  
317 encounters spanning January 1, 2021 through December 31, 2021 to incorporate  
318 the latest reformed E/M LoS CPT coding definitions and criteria implemented by  
319 U.S. Centers for Medicare & Medicaid Services effective January 1, 2021<sup>38</sup>.

320 Clinic notes were included for ‘new’ patient encounters as defined by the  
321 usage of a new patient CPT code (i.e., CPT codes 99201–99205). The notes  
322 represent a wide range of patient (e.g., chief complaint, diagnosis, age) and  
323 clinician (e.g., providers, provider types (physician, physician-extender, nurse  
324 practitioners), hospital-based or community clinic sites) contexts. The LoS CPT  
325 code submitted previously for billing served as the ground truth label for each  
326 note. CPT codes corresponding to the five LoS strata were included (i.e., 99201–  
327 99205 for new patients). Individual notes were excluded if a ground-truth LoS  
328 CPT code label was unavailable or if the note text was missing. Pre-processing  
329 the clinic notes dataset included removing infrequent labels.

330  
331 ***Clinical Note De-Identification.*** To remove protected health information  
332 (PHI) from each clinical note, we leveraged a pipeline that utilized a Named  
333 Entity Recognition (NER) model to annotate clinical notes for PHI elements  
334 (Spark NLP, John Snow Labs; Lewes, Delaware). Detected PHI elements  
335 included every instance of mentioned patient age, city, country, date,  
336 doctor/clinician name, hospital name, identification number, medical record

337 number, health system organization/entity name, patient name, phone number,  
338 profession, state, street address number and name, username, and/or zip code.  
339 After identification, the PHI elements were masked in place using the type of  
340 element (**Table 1**).

341  
342 **Model Development.** We identified several state-of-the-art text classification  
343 models for these tasks. *Bio\_ClinicalBERT*, a clinical language model pre-trained  
344 on EHR notes, was selected due to its demonstrated performance on clinical  
345 tasks<sup>39</sup>. We also fine-tuned two general domain text classification models,  
346 *DistilBERT*<sup>40</sup> and *XLNet*<sup>21</sup>. *DistilBERT* was chosen because it represents a smaller  
347 model that can run computationally constrained environments. *XLNet* was  
348 chosen as it has performance improvements over the *BERT* architecture<sup>21</sup>. As  
349 clinic notes can be longer than 512 tokens, we also fine-tuned the *Clinical-*  
350 *Longformer*<sup>41</sup> model, which can handle an input sequence length of up to 4,096  
351 tokens.

352  
353 We fine-tuned each model on all the clinic notes and separately on each  
354 subspecialty note source. Model performance on the test set was determined by  
355 clinic note classification accuracy. We also computed weighted F1 scores to  
356 account for the imbalance of the classes. The F1 score combines the precision and  
357 recall of a classifier into a single metric. The weighted F1 score is the F1 score for

358 each class weighted by its proportion in the dataset. Weighted F1 is useful in  
359 settings where an assessment of overall model performance is desired while  
360 accounting for the class imbalance. This analytic approach was repeated to serve  
361 the secondary objective using de-identified clinic notes for model fine-tuning.

362 For each fine-tuning experiment, the relevant clinic notes dataset was divided  
363 into 80% for model fine-tuning and 20% for testing. Model-specific tokenizers  
364 were used, and notes were padded to the longest sequence in the batch to ensure  
365 consistent input sequence lengths. Notes that exceeded the maximum token  
366 length of the model were truncated. Default fine-tuning training arguments were  
367 used across all experiments, including a learning rate of  $2 \times 10^{-5}$ , a batch size of 3  
368 due to memory constraints, five epochs, and a 0.01 weight decay.

369

370 ***Computing Environment.*** Data processing and NLP modeling were  
371 completed in Python (vers. 3.9) and PyTorch (vers 1.13.1). We utilized the  
372 HuggingFace transformers hub to source the models and adapt the pre-trained  
373 model fine-tuning pipeline (available at: <https://huggingface.co/>). Models were  
374 trained in a Linux environment with one NVIDIA T4 GPU with 16GB of  
375 memory.

376

377 ***Data availability:***

378

379 The data supporting this study's findings are not publicly available due to the  
380 datasets containing protected health information (PHI) that could compromise  
381 research participant privacy.

382

383

384

385

386 **ACKNOWLEDGEMENTS**

387

388 Dr. Crowson's effort is partly supported by an NIH grant (Biomedical

389 Informatics and Data Science Research Training Program; T15LM007092-30; PI

390 Nils Gehlenborg). The authors thank and acknowledge John Snow Labs/Spark

391 NLP for providing an academic research license for the de-identification toolkit.

392

393

394 **TABLES**

395

396

397 **Table 1.** Mock example of de-identification process using named entity recognition of

398 protected health information and obfuscation.

399

400

<p>NAME: <b>Earl Fullness</b> MRN: <b>138582469</b> DOB: <b>01/23/1987</b> Date of service: <b>03/02/2023</b> Location: <b>Mass Eye &amp; Ear Hospital</b> PCP: Dr. <b>B. Sick</b></p> <p>ASSESSMENT:</p> <p><b>35-year-old</b> male with a history of repeated left ear infections.</p> <p>He has previously sought treatment at several hospitals in the <b>Boston</b> area including <b>Mass General Hospital</b> and <b>Brigham &amp; Women's</b>. He reports experiencing pain and discomfort in his left ear, accompanied by redness, drainage, and decreased hearing. The patient has no other significant past medical history and takes no regular medications. No allergies are reported.</p> <p>Further evaluation, including a physical examination and potentially imaging studies, is recommended to determine the cause of the persistent infections and to develop an appropriate treatment plan.</p> <p>Provider: Dr. <b>N Cerumen</b> Department: Otolaryngology-Head &amp; Neck Surgery Phone Number: <b>857-123-4567</b></p>	<p>NAME: &lt;NAME&gt; MRN: &lt;ID&gt; DOB: &lt;DATE&gt; Date of service: &lt;DATE&gt; Location: &lt;LOCATION&gt; PCP: Dr. &lt;NAME&gt;:</p> <p>ASSESSMENT:</p> <p>&lt;AGE&gt; male with a history of repeated left ear infections.</p> <p>He has previously sought treatment at several hospitals in the &lt;LOCATION&gt; area including &lt;LOCATION&gt; and &lt;LOCATION&gt;.</p> <p>He reports experiencing pain and discomfort in his left ear, accompanied by redness, drainage, and decreased hearing. The patient has no other significant past medical history and takes no regular medications. No allergies are reported.</p> <p>Further evaluation, including a physical examination and potentially imaging studies, is recommended to determine the cause of the persistent infections and to develop an appropriate treatment plan.</p> <p>Provider: Dr. &lt;NAME&gt; Department: Otolaryngology-Head &amp; Neck Surgery Phone Number: &lt;CONTACT&gt;</p>
--	--

**Table 2.** Distribution of Evaluation and Management (E/M) level of service (LoS) codes across different specialties included in modeling.

LoS E/M Code	Proportion of Notes (%)					Mean (CI-95)
	Clinic Note Source Stratified by Specialty					
	Full Dataset (n = 31,115)	Cardiology (n = 13,279)	Gastroenterology (n = 5,299)	Internal Medicine (n = 10,178)	Otolaryngology-HNS (n = 2,197)	
Level I New	< 0.0	0.0	0.0	< 0.0	< 0.0	0.0 (0.0-0.0)
Level II New	2.2	0.8	< 0.0	5.3	0.3	2.2 (0.4-3.9)
Level III New	25.3	8.2	14.8	43.9	69.3	32.3 (12.9-51.7)
Level IV New	48.1	45.1	73.5	42.6	29.9	47.8 (35.3-60.3)
Level V New	24.3	45.9	11.0	8.1	0.5	18.0 (4.0-31.9)

**Table 3.** LoS E/M Code classification accuracy for general domain and clinical language models finetuned on notes from different medical specialties. All metrics are reported on test set data.

Model	Note Classification Accuracy (%)				
	Full Dataset (n = 6,223)	Clinic Note Source			
		Cardiology (n = 2,656)	Gastroenterology (n = 1,060)	Internal Medicine (n = 2,036)	Otolaryngology-HNS (n = 440)
<i>Bio_ClinicalBERT</i>	74.9	76.4	79.8	67.8	<b>87.3</b>
<i>Clinical-Longformer</i>	<b>75.8</b>	<b>79.0</b>	<b>82.5</b>	<b>68.1</b>	84.3
<i>DistilBERT</i>	74.1	75.7	80.1	67.7	84.1
<i>XLNet</i>	74.7	75.5	81.1	67.1	82.3
<i>Mean (CI-95)</i>	74.8 (74.1-75.6)	76.7 (75.1-78.2)	80.9 (79.7-82.0)	67.7 (67.3-68.1)	84.5 (82.5-86.5)

**Table 4.** Weighted F1-scores for prediction of LoS E/M codes across general domain and clinical language models finetuned on notes from different medical specialties. All metrics are reported on test set data.

Model	Model Weighted F1-Score				
	Full Dataset (n = 6,223)	Clinic Note Source			
		Cardiology (n = 2,656)	Gastroenterology (n = 1,060)	Internal Medicine (n = 2,036)	Otolaryngology-HNS (n = 440)
<i>Bio_ClinicalBERT</i>	0.74	<b>0.76</b>	<b>0.79</b>	<b>0.67</b>	<b>0.85</b>
<i>Clinical-Longformer</i>	<b>0.79</b>	0.75	<b>0.79</b>	0.63	<b>0.85</b>
<i>DistilBERT</i>	0.74	0.75	0.78	0.55	0.84
<i>XLNet</i>	0.74	0.74	<b>0.79</b>	0.66	0.82
<i>Mean (CI-95)</i>	0.75 (0.73-0.78)	0.75 (0.74-0.76)	0.79 (0.78-0.79)	0.63 (0.57-0.68)	0.83 (0.82-0.85)

**Table 5.** De-identified note classification performance for models trained on all clinic notes. All metrics are reported on test set data.

<b>Model</b>	<b>LoS Classification Accuracy (%)</b>	<b>Weighted F1-Score</b>
	<i>Full Dataset (n = 6,223)</i>	
<i>Bio_ClinicalBERT</i>	48.4	0.45
<i>Clinical-Longformer</i>	47.9	0.42
<i>DistilBERT</i>	48.3	0.43
<i>XLNet</i>	48.1	0.44
<i>Mean (CI-95)</i>	48.3 (48.0-48.5)	0.43 (0.42-0.45)

## REFERENCES

- 1 Cutler, D. M. & Ly, D. P. The (paper) work of medicine: understanding international medical costs. *J Econ Perspect* **25**, 3-25 (2011). <https://doi.org:10.1257/jep.25.2.3>
- 2 Himmelstein, D. U. *et al.* A comparison of hospital administrative costs in eight nations: US costs exceed all others by far. *Health Aff (Millwood)* **33**, 1586-1594 (2014). <https://doi.org:10.1377/hlthaff.2013.1327>
- 3 Jiwani, A., Himmelstein, D., Woolhandler, S. & Kahn, J. G. Billing and insurance-related administrative costs in United States' health care: synthesis of micro-costing evidence. *BMC Health Serv Res* **14**, 556 (2014). <https://doi.org:10.1186/s12913-014-0556-7>
- 4 Tseng, P., Kaplan, R. S., Richman, B. D., Shah, M. A. & Schulman, K. A. Administrative Costs Associated With Physician Billing and Insurance-Related Activities at an Academic Health Care System. *JAMA* **319**, 691-697 (2018). <https://doi.org:10.1001/jama.2017.19148>
- 5 Himmelstein, D. U., Campbell, T. & Woolhandler, S. Health Care Administrative Costs in the United States and Canada, 2017. *Ann Intern Med* **172**, 134-142 (2020). <https://doi.org:10.7326/M19-2818>
- 6 Casalino, L. P. *et al.* What does it cost physician practices to interact with health insurance plans? *Health Aff (Millwood)* **28**, w533-543 (2009). <https://doi.org:10.1377/hlthaff.28.4.w533>
- 7 Sakowski, J. A., Kahn, J. G., Kronick, R. G., Newman, J. M. & Luft, H. S. Peering into the black box: billing and insurance activities in a medical group. *Health Aff (Millwood)* **28**, w544-554 (2009). <https://doi.org:10.1377/hlthaff.28.4.w544>
- 8 Blanchfield, B. B., Heffernan, J. L., Osgood, B., Sheehan, R. R. & Meyer, G. S. Saving billions of dollars--and physicians' time--by streamlining billing practices. *Health Aff (Millwood)* **29**, 1248-1254 (2010). <https://doi.org:10.1377/hlthaff.2009.0075>
- 9 Boranbayev, A. S. & Boranbayev, S. N. in *2010 Seventh International Conference on Information Technology: New Generations* 1282-1284 (2010).
- 10 Koleck, T. A., Dreisbach, C., Bourne, P. E. & Bakken, S. Natural language processing of symptoms documented in free-text narratives of electronic health records: a systematic review. *J Am Med Inform Assoc* **26**, 364-379 (2019). <https://doi.org:10.1093/jamia/ocy173>
- 11 Wang, J. *et al.* Systematic Evaluation of Research Progress on Natural Language Processing in Medicine Over the Past 20 Years: Bibliometric Study on PubMed. *J Med Internet Res* **22**, e16816 (2020). <https://doi.org:10.2196/16816>
- 12 Juhn, Y. & Liu, H. Artificial intelligence approaches using natural language processing to advance EHR-based clinical research. *J Allergy Clin Immunol* **145**, 463-469 (2020). <https://doi.org:10.1016/j.jaci.2019.12.897>
- 13 Locke, S. *et al.* Natural language processing in medicine: A review. *Trends in Anaesthesia and Critical Care* **38**, 4-9 (2021). <https://doi.org:10.1016/j.tacc.2021.02.007>

- 14 Marafino, B. J. *et al.* Validation of Prediction Models for Critical Care Outcomes Using Natural Language Processing of Electronic Health Record Data. *JAMA Netw Open* **1**, e185097 (2018). [https://doi.org:10.1001/jamanetworkopen.2018.5097](https://doi.org/10.1001/jamanetworkopen.2018.5097)
- 15 Patra, B. G. *et al.* Extracting social determinants of health from electronic health records using natural language processing: a systematic review. *J Am Med Inform Assoc* **28**, 2716-2727 (2021). [https://doi.org:10.1093/jamia/ocab170](https://doi.org/10.1093/jamia/ocab170)
- 16 Wu, S. *et al.* Deep learning in clinical natural language processing: a methodical review. *J Am Med Inform Assoc* **27**, 457-470 (2020). [https://doi.org:10.1093/jamia/ocz200](https://doi.org/10.1093/jamia/ocz200)
- 17 Young, I. J. B., Luz, S. & Lone, N. A systematic review of natural language processing for classification tasks in the field of incident reporting and adverse event analysis. *Int J Med Inform* **132**, 103971 (2019). [https://doi.org:10.1016/j.ijmedinf.2019.103971](https://doi.org/10.1016/j.ijmedinf.2019.103971)
- 18 Abadeer, M. in *Proceedings of the 3rd clinical natural language processing workshop*. 158-167.
- 19 Li, Y. *et al.* BEHRT: transformer for electronic health records. *Scientific reports* **10**, 1-12 (2020).
- 20 Pascual, D., Luck, S. & Wattenhofer, R. Towards BERT-based automatic ICD coding: Limitations and opportunities. *arXiv preprint arXiv:2104.06709* (2021).
- 21 Zhang, Z., Liu, J. & Razavian, N. BERT-XML: Large scale automated ICD coding using BERT pretraining. *arXiv preprint arXiv:2006.03685* (2020).
- 22 Levy, J., Vattikonda, N., Haudenschild, C., Christensen, B. & Vaickus, L. Comparison of machine-learning algorithms for the prediction of current procedural terminology (CPT) codes from pathology reports. *Journal of Pathology Informatics* **13**, 100165 (2022).
- 23 Redd, T. K. *et al.* Variability in Electronic Health Record Usage and Perceptions among Specialty vs. Primary Care Physicians. *AMIA ... Annual Symposium proceedings* **2015**, 2053-2062 (2015).
- 24 Cohen, G. R., Friedman, C. P., Ryan, A. M., Richardson, C. R. & Adler-Milstein, J. Variation in Physicians' Electronic Health Record Documentation and Potential Patient Harm from That Variation. *Journal of general internal medicine : JGIM* **34**, 2355-2367 (2019). [https://doi.org:10.1007/s11606-019-05025-3](https://doi.org/10.1007/s11606-019-05025-3)
- 25 Bastani, H., Goh, J. & Bayati, M. Evidence of Upcoding in Pay-for-Performance Programs. *Management Science* **65**, 1042-1060 (2019). [https://doi.org:10.1287/mnsc.2017.2996](https://doi.org/10.1287/mnsc.2017.2996)
- 26 Brunt, C. S. CPT fee differentials and visit upcoding under Medicare Part B. *Health Econ* **20**, 831-841 (2011). [https://doi.org:10.1002/hec.1649](https://doi.org/10.1002/hec.1649)
- 27 Centers for, M. & Medicaid, S. Physician Code Creep: Evidence in Medicaid and State Employee Health Insurance Billing : Health Care Financing Review  
2007 ASI 4652-1.915. *Physician Code Creep: Evidence in Medicaid and State Employee Health Insurance Billing : Health Care Financing Review* (2007).
- 28 Chan, B., Anderson, G. M. & Theriault, M. E. Fee code creep among general practitioners and family physicians in Ontario: Why does the ratio of intermediate to minor assessments keep climbing? *Canadian Medical Association journal* **158**, 749-754 (1998).
- 29 Bowblis, J. R. & Brunt, C. S. Medicare skilled nursing facility reimbursement and upcoding. *Health Econ* **23**, 821-840 (2014). [https://doi.org:10.1002/hec.2959](https://doi.org/10.1002/hec.2959)

- 30 Shin, H., Lee, J., An, Y. & Cho, S. A scoring model to detect abusive medical institutions based on patient classification system: Diagnosis-related group and ambulatory patient group. *J Biomed Inform* **117**, 103752 (2021). <https://doi.org/10.1016/j.jbi.2021.103752>
- 31 Kushida, C. A. *et al.* Strategies for de-identification and anonymization of electronic health record data for use in multicenter research studies. *Med Care* **50 Suppl**, S82-101 (2012). <https://doi.org/10.1097/MLR.0b013e3182585355>
- 32 Meystre, S. M. *et al.* Text de-identification for privacy protection: A study of its impact on clinical text information content. *Journal of biomedical informatics* **50**, 142-150 (2014). <https://doi.org/10.1016/j.jbi.2014.01.011>
- 33 Bellamy, D., Hernán, M. A. & Beam, A. A structural characterization of shortcut features for prediction. *European Journal of Epidemiology* **37**, 563-568 (2022).
- 34 Obeid, J. S. *et al.* Impact of De-Identification on Clinical Text Classification Using Traditional and Deep Learning Classifiers. *Studies in health technology and informatics* **264**, 283-287 (2019). <https://doi.org/10.3233/SHTI190228>
- 35 Ferrandez, O. *et al.* BoB, a best-of-breed automated text de-identification system for VHA clinical documents. *J Am Med Inform Assoc* **20**, 77-83 (2013). <https://doi.org/10.1136/amiajnl-2012-001020>
- 36 Sepas, A., Bangash, A. H., Alraoui, O., El Emam, K. & El-Hussuna, A. Algorithms to anonymize structured medical and healthcare data: A systematic review. *Front Bioinform* **2**, 984807 (2022). <https://doi.org/10.3389/fbinf.2022.984807>
- 37 Luo, W. *et al.* Guidelines for Developing and Reporting Machine Learning Predictive Models in Biomedical Research: A Multidisciplinary View. *J Med Internet Res* **18**, e323 (2016). <https://doi.org/10.2196/jmir.5870>
- 38 CMS. *List of CPT/HCPCS Codes*, <[https://www.cms.gov/medicare/fraud-and-abuse/physicianselfreferral/list\\_of\\_codes](https://www.cms.gov/medicare/fraud-and-abuse/physicianselfreferral/list_of_codes)> (2023).
- 39 Alsentzer, E. *et al.* Publicly available clinical BERT embeddings. *arXiv preprint arXiv:1904.03323* (2019).
- 40 Sanh, V., Debut, L., Chaumond, J. & Wolf, T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- 41 Li, Y., Wehbe, R. M., Ahmad, F. S., Wang, H. & Luo, Y. A comparative study of pretrained language models for long clinical text. *Journal of the American Medical Informatics Association* **30**, 340-347 (2023).