

1 **Title:** Algorithmic identification of treatment-emergent adverse events from clinical notes using large
2 language models: a pilot study in inflammatory bowel disease

3
4 **Authors:** Anna L Silverman, MD^{1,2,*}; Madhumita Sushil, PhD^{3,*}; Balu Bhasuran, PhD^{3,*}; Dana Ludwig,
5 MD³; James Buchanan, PharmD³; Rebecca Racz, PharmD⁴; Mahalakshmi Parakala⁵; Samer El-Kamary,
6 MD, MPH⁴⁺; Ohenewaa Ahima, MD⁴; Artur Belov, PhD⁴; Lauren Choi, PharmD⁴; Monisha Billings, DDS,
7 MPH, PhD⁴; Yan Li, B.Pharm, Ph.D.⁴; Nadia Habal, MD⁴; Qi Liu, PhD⁴; Jawahar Tiwari, PhD⁴; Atul J Butte,
8 PhD^{3,6}; Vivek A Rudrapatna, MD, PhD^{3,7}

9
10 ¹Division of Gastroenterology and Hepatology, Department of Medicine, Mayo Clinic, Phoenix, AZ, USA

11 ²Department of Medicine, University of California, San Diego, La Jolla, CA, USA

12 ³Bakar Computational Health Sciences Institute, San Francisco, CA, USA

13 ⁴United States Food and Drug Administration, Silver Spring, MD, USA

14 ⁵Department of Public Health, University of California Berkeley, Berkeley, USA

15 ⁶Center for Data-Driven Insights and Innovation, University of California Health, Oakland, CA, USA

16 ⁷Division of Gastroenterology and Hepatology, Department of Medicine, University of California, San
17 Francisco, San Francisco, CA, USA

18 *These authors contributed equally to this work and share first authorship.

19 + This work was done while at the FDA. Current affiliation University of Maryland School of Medicine,
20 Baltimore, MD, USA

21
22 **Word Count:** 3,599

23 24 **Corresponding Author**

25 Vivek A. Rudrapatna, MD, PhD
26 Assistant Professor of Medicine
27 University of California, San Francisco Bakar Institute, Box 2993
28 490 Illinois Street, Floor 2
29 San Francisco, CA 94143
30 Email: vivek.rudrapatna@ucsf.edu

31 32 **Funding Source:**

33 This publication was supported by the Food and Drug Administration (FDA) of the U.S. Department of
34 Health and Human Services (HHS) as part of a financial assistance award Center of Excellence in
35 Regulatory Science and Innovation grant to University of California, San Francisco, U01FD005978,
36 totaling \$79,250 with 33% percentage funded by FDA/HHS and \$158,500, 66% percentage funded by the
37 UCSF Division of Gastroenterology and UCSF Bakar Computational Health Sciences Institute, and 1%
38 funded by the National Library of Medicine of the National Institutes of Health under Award Number
39 K99LM014099. Additional support for clinical data resources were provided by National Center for
40 Advancing Translational Sciences, National Institutes of Health, through UCSF-CTSI Grant Number
41 UL1TR001872. The contents are those of the authors and do not necessarily represent the official views
42 of, nor an endorsement, by HHS or the U.S. Government.

43 44 45 **Disclosures:**

46 ALS: nothing to disclose

47 MS: nothing to disclose

48 BB: nothing to disclose

49 DL: nothing to disclose

50 JB: nothing to disclose

51 RR: nothing to disclose

52 MP: nothing to disclose

53 SE: nothing to disclose

54 OA: nothing to disclose

55 AB: nothing to disclose

56 LC: nothing to disclose

57 MB: nothing to disclose

58 YL: nothing to disclose

59 NH: nothing to disclose

60 QL: nothing to disclose

61 JT: nothing to disclose

62 AJB: AJB is a co-founder and consultant to Personalis and NuMedii; consultant to Mango Tree
63 Corporation, and in the recent past, Samsung, 10x Genomics, Helix, Pathway Genomics, and Verinata
64 (Illumina); has served on paid advisory panels or boards for Geisinger Health, Regenstrief Institute,
65 Gerson Lehman Group, AlphaSights, Covance, Novartis, Genentech, and Merck, and Roche; is a
66 shareholder in Personalis and NuMedii; is a minor shareholder in Apple, Meta (Facebook), Alphabet
67 (Google), Microsoft, Amazon, Snap, 10x Genomics, Illumina, Regeneron, Sanofi, Pfizer, Royalty Pharma,
68 Moderna, Sutro, Doximity, BioNtech, Invitae, Pacific Biosciences, Editas Medicine, Nuna Health, Assay
69 Depot, and Vet24seven, and several other non-health related companies and mutual funds; and has

70 received honoraria and travel reimbursement for invited talks from Johnson and Johnson, Roche,
71 Genentech, Pfizer, Merck, Lilly, Takeda, Varian, Mars, Siemens, Optum, Abbott, Celgene, AstraZeneca,
72 AbbVie, Westat, and many academic institutions, medical or disease specific foundations and
73 associations, and health systems. Atul Butte receives royalty payments through Stanford University, for
74 several patents and other disclosures licensed to NuMedii and Personalis. Atul Butte’s research has
75 been funded by NIH, Peraton (as the prime on an NIH contract), Genentech, Johnson and Johnson, FDA,
76 Robert Wood Johnson Foundation, Leon Lowenstein Foundation, Intervalien Foundation, Priscilla Chan
77 and Mark Zuckerberg, the Barbara and Gerson Bakar Foundation, and in the recent past, the March of
78 Dimes, Juvenile Diabetes Research Foundation, California Governor’s Office of Planning and Research,
79 California Institute for Regenerative Medicine, L’Oreal, and Progenity.
80 VAR: Receives grant support from Merck, Alnylam, Genentech, Stryker, Blueprint Medicines, Takeda,
81 and Janssen.

82 **Involvement with the Manuscript:**

83 Anna L. Silverman: study concept and design; acquisition of data; analysis and interpretation of data;
84 lead annotation protocol; lead drafting of the manuscript; lead critical revision of the manuscript for
85 important intellectual content

86 Madhumita Sushil: co-lead model architect, study concept and design; acquisition of data; analysis and
87 interpretation of data; drafting of the manuscript; critical revision of the manuscript for important
88 intellectual content; technical support; implementation and model design

89 Balu Bhasuran: co-lead model architect, study concept and design; acquisition of data; analysis and
90 interpretation of data; drafting of the manuscript; critical revision of the manuscript for important
91 intellectual content; technical support; implementation and model design

92 Dana Ludwig: study concept and design; acquisition of data; drafting of the manuscript; critical revision
93 of the manuscript for important intellectual content; technical support; implementation and model
94 design
95 James Buchanan: study concept and design; acquisition of data; analysis and interpretation of data;
96 critical revision of the manuscript for important intellectual content
97 Rebecca Racz: study concept and design; acquisition of data; analysis and interpretation of data; critical
98 revision of the manuscript for important intellectual content
99 Mahalakshmi Parakala: acquisition of data
100 Samer El-Kamary: lead investigator at the FDA for this project, administrative responsibility and study
101 team supervision at the FDA, intellectual contribution during the conduct of the study, and critical
102 revision of the manuscript for important intellectual content.
103 Ohenewaa Ahima: critical revision of the manuscript for important intellectual content; study
104 supervision
105 Artur Belov: critical revision of the manuscript for important intellectual content; study supervision
106 Lauren Choi: critical revision of the manuscript for important intellectual content; study supervision
107 Monisha Billings: critical revision of the manuscript for important intellectual content; study supervision
108 Yan Li: critical revision of the manuscript for important intellectual content; study supervision
109 Nadia Habal: critical revision of the manuscript for important intellectual content; study supervision
110 Qi Liu: critical revision of the manuscript for important intellectual content; study supervision
111 Jawahar Tiwari: critical revision of the manuscript for important intellectual content; study supervision
112 Atul Butte: study supervision
113 Vivek A Rudrapatna: study concept and design; acquisition of data; analysis and interpretation of data;
114 technical support; critical revision of the manuscript for important intellectual content; study
115 supervision

116 **Disclaimer:**

117 **The contents of this article reflect the views of the authors and should not be construed to represent**
118 **the FDA's views or policies. No official support or endorsement by the FDA is intended or should be**
119 **inferred.**

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138 **Abstract**

139 **Background and Aims:**

140 Outpatient clinical notes are a rich source of information regarding drug safety. However, data
141 in these notes are currently underutilized for pharmacovigilance due to methodological limitations in
142 text mining. Large language models (LLM) like BERT have shown progress in a range of natural language
143 processing tasks but have not yet been evaluated on adverse event detection.

144 **Methods:**

145 We adapted a new clinical LLM, UCSF BERT, to identify serious adverse events (SAEs) occurring
146 after treatment with a non-steroid immunosuppressant for inflammatory bowel disease (IBD). We
147 compared this model to other language models that have previously been applied to AE detection.

148 **Results:**

149 We annotated 928 outpatient IBD notes corresponding to 928 individual IBD patients for all SAE-
150 associated hospitalizations occurring after treatment with a non-steroid immunosuppressant. These
151 notes contained 703 SAEs in total, the most common of which was failure of intended efficacy. Out of 8
152 candidate models, UCSF BERT achieved the highest numerical performance on identifying drug-SAE pairs
153 from this corpus (accuracy 88-92%, macro F1 61-68%), with 5-10% greater accuracy than previously
154 published models. UCSF BERT was significantly superior at identifying hospitalization events emergent to
155 medication use ($p < 0.01$).

156 **Conclusions:**

157 LLMs like UCSF BERT achieve numerically superior accuracy on the challenging task of SAE detection
158 from clinical notes compared to prior methods. Future work is needed to adapt this methodology to
159 improve model performance and evaluation using multi-center data and newer architectures like GPT.
160 Our findings support the potential value of using large language models to enhance pharmacovigilance.

161 **Keywords:** pharmacovigilance, artificial intelligence, natural language processing, large language

162 models, BERT, adverse event detection, inflammatory bowel disease

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182 Introduction

183 The accurate detection of treatment-emergent adverse events (AEs) is critical to ensure that
184 clinicians and patients can make well-informed treatment decisions that balance risks with benefits. This
185 is particularly true of non-steroid immunosuppressants which are commonly needed long-term for the
186 treatment of inflammatory bowel diseases (IBD).

187 Existing approaches for AE surveillance may involve prospective registry studies, spontaneous
188 postmarketing reporting (e.g., the Food and Drug Administration's AE Reporting System [FAERS])¹,
189 literature searches, and/or analyses of the structured data from claims and electronic health records
190 databases^{2,3}. These approaches have provided important data on the postmarket safety of medications
191 but are limited by expense, small numbers, under/over-reporting^{4,5}, missing data, limitations in inferring
192 causality, and suboptimal sensitivity and specificity. Clinical notes are a rich source of AE data because
193 treating clinicians often document actions in response to potential AEs, including treatment
194 discontinuation and hospitalization. However, these notes have been underutilized for surveillance due
195 to methodological limitations in effective text mining.

196 Recent years have seen impressive advances in natural language processing following the release of
197 the large language model known as BERT (Bidirectional Encoder Representations from Transformers)⁶.
198 However, its adaptation to domain-specific arenas like medical language has been limited, in part due to
199 the unavailability of safe platforms for processing this protected health information until recently. In
200 prior work, our group of academic researchers has developed a new BERT model specifically designed to
201 interpret clinical text as typically documented in electronic health records (EHR) systems⁷. This model,
202 UCSF BERT, was trained on 75 million clinical notes documented across a range of specialties over the
203 last 10 years at the University of California, San Francisco (UCSF). Evaluations of UCSF BERT on several
204 general benchmarks show that it performs as well as or better than other comparable BERT models not
205 specifically trained from scratch using a diverse corpus of notes derived from EHRs⁷. However, these

206 prior evaluations were general tasks and are limited by the quality of currently available, publicly
207 benchmarked tasks.

208 An open question motivating this study was whether the BERT model could help automate specific
209 tasks of established clinical importance, particularly one as challenging as AE detection. Many aspects
210 make the task of AE detection from clinical notes particularly difficult. These include the length of typical
211 clinical notes, the need to infer relationships between medications and documented AEs, to encode AEs
212 in a standardized way, and to overcome inherent vagueness in the documentation of clinical notes.
213 Some current examples of automated AE detection come from the National Natural Language
214 Processing Clinical Challenges (N2C2)⁸ adverse event detection challenge, a nationwide clinical data
215 science competition that was held in 2018. Models from this competition were evaluated on highly
216 simplified benchmark tasks that do not reflect the typical patterns of clinical documentation such as
217 short snippets of notes rather than full length notes. Notably, none of the candidate models from the
218 competition were large language models as it was held prior to the wide-spread adoption of large
219 language models.

220 We hypothesized that adaptations of the UCSF BERT, a large language model, would outperform
221 previously published methods on multiple tasks related to AE detection, due in large part to its prior
222 training on a large volume of EHR notes. In this pilot study, we trained UCSF BERT to identify
223 hospitalization-associated serious adverse events (SAEs) from notes written in the outpatient IBD clinic,
224 and we compared its performance to a range of baselines including previously published models.

225

226 **Methods**

227 *Ethics*

228 This single-center study of natural language processing algorithms for adverse event detection
229 was approved by the UCSF Institutional Review Board (#18-24588).

230 *Target of prediction*

231 The target of prediction was treatment-emergent, outpatient SAEs requiring hospitalization with
232 exposure to non-steroid immunosuppressant drugs for IBD, as documented in the outpatient
233 gastroenterology clinic notes at UCSF. The candidate list of drugs included all biologics and small-
234 molecule medications, except steroids, that were approved by the FDA for the treatment of ulcerative
235 colitis or Crohn's disease as of 2020, as well as off-label medications that are occasionally used to treat
236 these conditions. A complete list of included medications can be found in the supplemental materials.
237 Although the FDA's definition of SAEs includes multiple categories⁹, we limited our scope to only SAEs
238 associated with a hospitalization event, as these are more likely to be well-documented in clinical notes
239 due to their clinical importance. We defined treatment-emergent as an SAE that occurred while the
240 patient was actively receiving scheduled doses of a given medication, having been absent pre-
241 treatment. For example, if a patient was hospitalized for pneumonia 6 weeks after receiving an infusion
242 that was prescribed to be given every 8 weeks, the hospitalization event would be considered a
243 treatment-emergent SAE. Once the clinical decision to discontinue a given treatment plan was
244 documented, subsequent hospitalization events were no longer considered treatment-emergent SAEs.
245 Worsening of previously existing conditions that prompted hospitalization were included in line with
246 internationally used guidelines on AE reporting¹⁰. A definitive assessment of potentially causal
247 relationships between treatments and SAEs was beyond the planned scope of this analysis.

248 *Document identification strategy*

249 To identify the target notes for this study, we used a deidentified research database consisting
250 of structured EHR data at UCSF as well as clinical notes that had been subjected to automated redaction
251 of protected health information¹¹. We queried the database to identify all notes associated with the
252 gastroenterology department and an IBD diagnosis code (ICD-9 555/556; ICD-10 K50/K51). We selected
253 notes written between 1/1/2018 and 12/31/2020 and utilized the most recent note for each patient

254 who was at least 18 years old during this period. We selected this timeframe to maximize the capture of
255 a wide range of FDA-approved treatments. We used the most recent note per patient to avoid double
256 counting SAEs mentioned in multiple notes, and to take advantage of the fact the documented histories
257 tend to be inclusive of prior events. All included notes were written by a gastroenterology physician or
258 advanced practice provider in the IBD outpatient clinic.

259 *Document Preprocessing*

260 The history of present illness (HPI) section of the notes was extracted using rule-based approaches
261 developed specifically for this project (supplemental methods). The HPI section of the note was the only
262 portion of the note utilized for downstream analysis as this section of the note often contains a
263 cumulative source of information on treatment exposures and outcomes, particularly out-of-system
264 events (i.e., hospitalizations and SAEs that occurred outside of UCSF but were relayed to the
265 gastroenterology provider at the time of routine follow-up). The HPI was pre-labeled with medications
266 of interest, hospitalization and signs and symptoms using, named entity recognition functions, from the
267 clinical natural language processing software *cTAKES*¹², as well as regular expressions (i.e., the ability to
268 locate pre-defined key-words). To minimize downstream algorithmic confusion in learning medication
269 names, the medication brand names were replaced with the generic name using the RxNorm
270 Application Program Interfaces (APIs) in Unified Medical Language System (UMLS)¹³.

271 *Note Annotation*

272 All notes that met the above inclusion criteria were annotated to fine-tune UCSF BERT on a
273 variety of AE detection-related tasks and to evaluate its performance against comparator models. A
274 team of five annotators, consisting of gastroenterologists, pharmacists, pharmacovigilance experts, and
275 patients carried out all annotation related tasks. These included the development and finalization of an
276 annotation protocol, participation in interrater reliability assessments, and annotation of all target
277 notes. The annotation protocol was collectively developed and refined over the course of weekly team

278 meetings utilizing an initial subset of notes. Using LabelStudio¹⁴, an open-source annotation platform,
279 annotators marked up the prelabelled HPI section of candidate notes according to triplets of medication
280 mentions, hospitalization mentions, and SAE mentions, if they corresponded to a hospitalization as per
281 the protocol (supplemental methods) (Figure 1). These annotations became the basis of the subsequent
282 efforts to train UCSF BERT and other models to automate this process.

283 All five annotators participated in an interrater reliability assessment on a sample of 19 notes. The
284 results of these assessments were reviewed in weekly meetings to improve the protocol as well as
285 annotator compliance to it. We computed a Fleiss' kappa statistic to characterize the interrater
286 reliability on the final round of assessments. Following this training and assessment phase, the protocol
287 was locked, and the remainder of the corpus was annotated.

288 289 *Modeling*

290 We defined several prediction tasks, asking the model to classify whole HPIs according to the
291 occurrence of: (task 1) all candidate medication mentions given prior to a hospitalization (task 2)
292 adverse event (AE) as reason for hospitalization and (task 3) the combination of task 1 and task 2 the
293 medication-hospitalization-AE triple (Figure 1). We trained models of different architectures to
294 determine which were best suited for the task of AE detection. We used scikit-learn¹⁵ to train several
295 baseline Bag of Words (BoW) models such as Logistic Regression, K-Nearest Neighbors, Decision Trees,
296 Random Forest, and XGBoost (supplemental methods). We used AutoGluon¹⁶ to train the automated
297 machine learning models. The annotated notes were split into 80% training, 10% validation and 10%
298 testing. These served as a baseline to compare the performance of our UCSF BERT model. We adapted
299 deep learning models architectures such as Convolutional Neural Network (CNN^{17,18}), Bidirectional Long
300 Short Term Memory Network (Bi-LSTM¹⁹) and Bi-LSTM with attention. These are deep learning models
301 adapted from the top performing entries in the N2C2⁸ adverse event detection challenge. All BERT
302 results are from the median performance of Macro F1 score over 5 runs of the model with different

303 seeds. Comparative model performance significance was evaluated using Fisher's exact test and chi
304 square with Yates's correlation when values were large enough to require it.

305 *Note Length Handling*

306 To include the entire HPI section, which was often longer than the typical maximum input length
307 used by other BERT models, we developed a hierarchical version of the UCSF-BERT²⁰ model (H-UCSF-
308 BERT). This model learns to process text using input sequences of 512 tokens (roughly equivalent to
309 words), in the same manner as a typical BERT model would. It then combines them into a longer-
310 sequence representation by integrating an additional transformer layer on top of these chunk
311 representations. We encoded sequences up to 2560 tokens, which is 5 times the usual processing limit
312 of a BERT model. We used the Mann-Whitney test to evaluate the possible association between note
313 length and the presence of SAEs.

314 *Handling of Class Imbalance*

315 SAEs were seen in 44% of notes in our corpus, however SAEs were uncommon once the notes
316 were subdivided into chunks that were ingestible by H-UCSF-BERT, creating a potential problem for
317 training models to learn to positively identify these SAEs when they do occur. To optimize learning in the
318 face of this imbalance in the dataset, the *training data* examples without AEs were randomly
319 undersampled. We explored a range of sampling ratios and identified the ratio of 1:4 positive to
320 negative examples as being best for model performance. This was applied to the training dataset for all
321 downstream tasks. We also explored additional strategies such as weighting the optimization loss based
322 on class distributions, as well as learning these weights dynamically²¹. However, we obtained the most
323 promising result by undersampling the majority dataset.

324 *MedDRA*

325 All SAEs were manually coded using the Medical Dictionary for Regulatory Activities (MedDRA)
326 version 23.0²².

327 **Results**

328 *Source Corpus, Patient Population, and SAE Dataset*

329 From a deidentified dataset of 110 million machine redacted clinical notes at UCSF, we
330 identified a total of 928 notes corresponding to 928 adults with IBD who were seen during the 2018-
331 2020 period. The patients in our study were 53% female with an average age of 45 years old (Table 1).
332 The most common race of patients was white. We annotated all 928 notes and performed interrater
333 reliability testing on a set of 19 notes to characterize the quality of the annotated dataset. The mean
334 observed agreement among the five annotators was 93-99% across all annotation categories
335 (Supplemental Table 1).

336 We identified a total of 703 SAEs in the 928 annotated notes from 928 patients with IBD. All
337 SAEs were associated with hospitalization as defined by the annotation protocol. Out of the 928
338 annotated notes, 411 documented at least one SAE (Table 2 and Supplemental Table 3). Importantly,
339 some notes included more than one SAE due to multiple distinct hospitalizations in the note. The notes
340 documenting an SAE tended to be longer than those without an SAE ($p < 0.001$). Over 60% of SAEs in our
341 corpus were associated with anti-tumor necrosis factor agents (anti-TNF). Infliximab was associated with
342 179 SAEs, the most of any drug, followed closely by adalimumab (136 SAEs; Table 3). This finding was
343 expected given that infliximab was the first biologic to be approved for IBD, and more patients have
344 been exposed to this medication than any other due to its longer availability. Additionally, given the
345 relative absence of alternative treatments in the early 2000s, it is likely that patients remained on
346 infliximab and other anti-TNFs for a longer period (even after experiencing SAEs), compared to the
347 current era with multiple approved medications. The most common SAE was failure of intended efficacy
348 ($N=299$), followed by infections ($N=94$) (Figure 2 and Table 4). However, SAEs were found for every
349 organ system and every non-steroid immunosuppressant. Our corpus contained only one episode of
350 cancer, sarcoma, which occurred in a patient receiving an anti-TNF drug. The complete list of SAEs

351 mapped by clinical note terms and MedDRA terms can be found in the supplemental materials
352 (Supplemental Table 3). However, to enable a more user-friendly exploration of trends in the data, we
353 have developed an interactive web application (see <https://ibd-ade.streamlit.app/>).

354 *Performance of UCSF-BERT on the Task of SAE Detection*

355 We established three targets of prediction for all downstream models: (task 1) identify all
356 candidate medication mentions given prior to a hospitalization, (task 2) identify adverse event as reason
357 for hospitalization and (task 3) the combination of task 1 and task 2 the medication-hospitalization-AE
358 triple (Figure 1). The annotated data was transformed and then split into training, validation, and testing
359 datasets for each of these binary classification tasks (Table 2 and Supplemental Table 2). We developed
360 and trained several variations of the UCSF-BERT model to address each of these targets. We then
361 evaluated its performance against several other comparator models, including several of the top entries
362 from the 2018 N2C2 adverse event detection challenge (supplemental methods).

363 On the task of medication prior to hospitalization, H-UCSF-BERT was the most performant model
364 with a Macro F1 of 62% (Table 5). It was significantly more accurate than the next-best model by a
365 margin of 11% ($p < 0.01$). Similarly, H-UCSF-BERT was the best model at the task of identifying
366 hospitalization relations to AEs with an accuracy of 96% and Macro F1 of 62%. We hypothesized that
367 long distances between mentions of a hospitalization and the associated SAEs could be reducing model
368 accuracy. Indeed, we found that restricting the input to SAEs mentioned within a two-sentence span of
369 the hospitalization, Macro F1 increased to 68% from 62% ($p < 0.01$). However, when compared to the
370 next performant model, BiLSTM, UCSF BERT was not significantly superior ($p = 0.40$). The ultimate goal
371 was to have our model accurately detect triples which include the mention of a non-steroid
372 immunosuppressant prior to a hospitalization plus the hospitalization plus the associated SAEs. For the
373 triples task, H-UCSF-BERT was again the best performer with a Macro F1 of 61%, however again this was
374 not significantly different than BiLSTM ($p = 0.40$).

375 **Discussion**

376 We adapted an EHR-specific clinical language model, UCSF BERT, to multiple tasks pertaining to
377 the detection of treatment-emergent serious adverse events. We have evaluated its performance in the
378 context of a specific use case: the use of non-steroid immunosuppressants for the treatment of IBD. We
379 have generated a gold standard corpus of 928 clinical notes as the basis of training and evaluating this
380 model against several baselines. Inter-rater reliability testing indicates good to excellent concordance
381 across annotators. UCSF BERT performed well in a range of tasks pertaining to SAE detection from
382 clinical notes. It achieves macro F1 scores ranging from 61-68% and accuracies from 88-92%. This model
383 numerically outperforms existing models for SAE detection associated with the N2C2 Challenge⁸ as well
384 as a range of strong baseline models, including several trained using automated machine learning. On
385 the task of accurately determining a medication of interest mentioned prior to a hospitalization, UCSF
386 BERT was significantly superior to all other models.

387 We found that the most common errors made by the models involved chains of reasoning
388 across many events. For example, instances where the reason of hospitalization is not explicitly
389 mentioned but merely implied from the clinical context. In addition, the model struggled in the setting
390 of both long-distance dependency where there were many sentences between entities of interest and
391 long chronology of events where several medication changes occurred over many sentences. Lastly,
392 when there were both non-specific adverse events such as pain or vomiting as well as more specific
393 terms such as small bowel obstruction or ulcerative colitis flare the combination was challenging for the
394 model to handle.

395 The last few decades have seen a significant expansion in FDA-approved therapies for IBD. In the
396 current era of IBD treatment with numerous agents available, continued monitoring for new safety
397 information on these agents is helpful to inform optimal treatment selection. The most frequent SAE
398 found in our corpus of outpatient IBD clinical notes at a tertiary referral center was failure of intended

399 efficacy followed by infections. This is in line with previously published data, especially in the setting of
400 more than 60% of the SAEs in our corpus associated with anti-tumor necrosis factor agents²³. We did not
401 account for concurrent use of steroids which are known to increase the risk of infection. However, our
402 corpus includes SAEs from every organ system. Of note, the non-steroid medications of interest are not
403 being prescribed with equal frequency; thus, prescribing practices are likely to influence the frequency
404 of events as well as frequencies of possible AEs associated with the medication. The strength of
405 association with SAEs and classes of non-steroid immunosuppressants can be explored using our
406 interactive web application (see <https://ibd-ade.streamlit.app/>). The goal of developing text-based
407 automation tools like this is to enable more precise characterizations of adverse events in the context of
408 routine clinical care to help validate known safety profiles of these drugs as well as identify previously
409 unrecognized SAEs which can point to areas of inquiry. For instance, our dataset included a patient
410 receiving an anti-TNF who was hospitalized for a new diagnosis of sarcoma. Sarcomas have previously
411 been reported in the context of children with IBD using anti-TNFs²⁴, although multiple long-term
412 observational studies have not consistently found a link between anti-TNF use and an increased risk of
413 cancer in adults^{25,26}. Future directions of this work include external validation using data from additional
414 centers, expansion to additional disease states outside of IBD, and downstream studies designed to
415 identify new drug-SAEs associations more rigorously using aggregated data.

416 Our work has many notable strengths. We have used transparent methods for developing the
417 training corpus and assessing its quality, including interrater reliability. Because our models have been
418 trained on de-identified clinical data, we intend to make them publicly available for others to reproduce
419 and enhance multiple aspects of this work. Of note, the N2C2 national challenge which produced
420 models for detection of treatment-emergent adverse events from clinical notes prior to our work was
421 before the release of BERT. We suspect that the underlying architecture of BERT in addition to our pre-

422 training from scratch on a sizeable clinical corpus are driving our improved performance compared to
423 prior models.

424 Some limitations of our work, outside of those common to retrospective research, include
425 imperfect accuracy of the model, which on certain tasks did not perform statistically significantly
426 superior to other models. We suspect this largely the result of long-distance dependence and long
427 chains of reasoning across many events in a clinical note. In addition, we have not yet assessed the
428 generalizability of our model across other diseases, treatments, or health systems. As well, our
429 interrater agreement is potentially optimistic as it was calculated iteratively on the same 19 notes.
430 However, there are no universally accepted standards of the Fleiss' kappa statistic²⁷ for good
431 agreement. Future work aimed at improving upon our current model includes annotating a larger corpus
432 at an outside health system to evaluate generalizability and over-sampling for SAEs to have more
433 positive examples for the model to learn from. Overall, our approach, utilizing novel methods from the
434 field of artificial intelligence, has the potential to address unmet needs in drug safety surveillance, an
435 area of central importance to regulatory agencies across the globe and to public health in general.

436 **Conclusion**

437 We have successfully adapted a new clinical language model, UCSF BERT, to the task of mining
438 outpatient clinic notes for SAEs occurring in patients with IBD administered non-steroid
439 immunosuppressants. This model performs well on the tasks of SAE detection, especially identifying
440 target medications prior to hospitalizations. The success of this model appears to stem from its
441 pretraining on a large and diverse corpus of notes derived from real-world clinical care and use of
442 hierarchical modeling which allows for long sequence document classification tasks. These results
443 suggest the feasibility of adapting artificial intelligence methods to address important unmet needs in
444 the field of pharmacovigilance, with the potential to substantially reduce the manual efforts needed to
445 review notes and identify events of concern. Our work is a step closer to a future of automated drug

446 surveillance algorithms embedded within EHR systems which can facilitate pharmacovigilance activities
447 ranging from health system reporting of SAEs to large-scale safety evaluations across multiple EHR
448 systems without the limitations of using billing codes as surrogates for actual AEs.

449 **Study Highlights**

- 450 • What is the current knowledge on the topic?
 - 451 ○ Prior work in automated adverse event (AE) detection from routine clinic notes utilized note
 - 452 fragments and simplified AE detection tasks. In addition, the newest model architectures
 - 453 were not widely available when automated AE detection was evaluated in a national natural
 - 454 language processing challenge in 2018.
- 455 • What question did this study address?
 - 456 ○ Are the newest model architectures trained on clinical notes capable of detecting serious
 - 457 AEs in routine clinical notes as written by clinicians seeing patients with inflammatory bowel
 - 458 diseases at a tertiary medical center.
- 459 • What does this study add to our knowledge?
 - 460 ○ Our model, UCSF BERT, trained on a large corpus of real-world clinical notes performs better
 - 461 than prior models previously designed for this task. Notably, our hierarchical model
 - 462 architecture is able to digest information five times the usual processing limit of BERT.
- 463 • How might this change clinical pharmacology or translational science?
 - 464 ○ Our work is a step closer to a future of automated drug surveillance algorithms embedded
 - 465 within EHR systems which can facilitate pharmacovigilance activities.

466 **Access to Data:**

467 The analytic code to train and evaluate models will be made publicly available at
468 https://github.com/MadhumitaSushil/ADE_detection. A machine-redacted version of the notes-based

469 data can be made available to requesting researchers by mutual agreement and following the execution
470 of a data use agreement.

471 **Acknowledgements:**

472 We gratefully acknowledge the invaluable administrative support provided by Lily Wong. In
473 addition, we would like to acknowledge the [UCSF Information Commons Computational Research](#)
474 [Platform](#), developed and supported by UCSF Bakar Computational Health Sciences Institute and UCSF
475 Academic Research Services. We would like to express our thanks to the Wynton support team and the
476 UCSF high-performance computing cluster, Wynton.

477 **Citations**

- 478
479
- 480 1. Questions and Answers on FDA's Adverse Event Reporting System (FAERS).
481 <https://www.fda.gov/drugs/surveillance/questions-and-answers-fdas-adverse-event-reporting-system-faers#:~:text=The%20FDA%20Adverse%20Event%20Reporting,that%20were%20submitted%20to%20FDA>. Accessed 05/19/2023.
482
483
 - 484 2. Thein D, Egeberg A, Skov L, Loft N. Absolute and Relative Risk of New-Onset Psoriasis
485 Associated With Tumor Necrosis Factor- α Inhibitor Treatment in Patients With Immune-
486 Mediated Inflammatory Diseases: A Danish Nationwide Cohort Study. *JAMA*
487 *dermatology*. 2022.
488
 - 489 3. Chaparro M, Garre A, Ricart E, et al. Short and long-term effectiveness and safety of
490 vedolizumab in inflammatory bowel disease: results from the ENEIDA registry.
491 *Alimentary pharmacology & therapeutics*. 2018;48(8):839-851.
492
 - 493 4. Hazell L, Shakir SA. Under-reporting of adverse drug reactions. *Drug safety*.
494 2006;29(5):385-396.
495
 - 496 5. Varallo FR, Guimarães SdOP, Abjaude SAR, Mastroianni PdC. Causes for the
497 underreporting of adverse drug events by health professionals: a systematic review.
498 *Revista da Escola de Enfermagem da USP*. 2014;48:739-747.
499
 - 500 6. Devlin J, Chang M-W, Lee K, Toutanova K. Bert: Pre-training of deep bidirectional
501 transformers for language understanding. *arXiv preprint arXiv:181004805*. 2018.
502
 - 503 7. Sushil M, Ludwig D, Butte AJ, Rudrapatna VA. Developing a general-purpose clinical
504 language inference model from a large corpus of clinical notes. *arXiv preprint arXiv:221006566*. 2022.
 - 505 8. Henry S, Buchan K, Filannino M, Stubbs A, Uzuner O. 2018 n2c2 shared task on adverse
506 drug events and medication extraction in electronic health records. *Journal of the American Medical Informatics Association*. 2019;27(1):3-12.

- 505 9. Code of Federal Regulations Title 21. In: Administration USFaD, ed. *Title 21, Volume 5,*
506 *21CFR312.32.*
- 507 10. Agency EM. ICH Topic E9 Statistical Principles for Clinical Trials
508 [https://www.ema.europa.eu/en/documents/scientific-guideline/ich-e-9-statistical-](https://www.ema.europa.eu/en/documents/scientific-guideline/ich-e-9-statistical-principles-clinical-trials-step-5_en.pdf)
509 [principles-clinical-trials-step-5_en.pdf.](https://www.ema.europa.eu/en/documents/scientific-guideline/ich-e-9-statistical-principles-clinical-trials-step-5_en.pdf)
- 510 11. Norgeot B, Muenzen K, Peterson TA, et al. Protected Health Information filter (Philter):
511 accurately and securely de-identifying free-text clinical notes. *NPJ digital medicine.*
512 2020;3(1):57.
- 513 12. Savova GK, Masanz JJ, Ogren PV, et al. Mayo clinical Text Analysis and Knowledge
514 Extraction System (cTAKES): architecture, component evaluation and applications.
515 *Journal of the American Medical Informatics Association.* 2010;17(5):507-513.
- 516 13. Bodenreider O. The unified medical language system (UMLS): integrating biomedical
517 terminology. *Nucleic acids research.* 2004;32(suppl_1):D267-D270.
- 518 14. *Label Studio: Data Labeling Software* [computer program]. 2020-2022.
- 519 15. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine learning in Python.
520 *the Journal of machine Learning research.* 2011;12:2825-2830.
- 521 16. Erickson N, Mueller J, Shirkov A, et al. Autogluon-tabular: Robust and accurate automl
522 for structured data. *arXiv preprint arXiv:200306505.* 2020.
- 523 17. Lecun Y, Bottou L, Bengio Y, Haffner P. Gradient-based learning applied to document
524 recognition. *Proceedings of the IEEE.* 1998;86(11):2278-2324.
- 525 18. Yann LeCun BB, John Denker, Donnie Henderson, R. Howard, Wayne Hubbard, Lawrence
526 Jackel. Handwritten Digit Recognition with a Back-Propagation Network. Paper
527 presented at: Advances in Neural Information Processing Systems 21989.
- 528 19. Hochreiter S, Schmidhuber J. Long Short-Term Memory. *Neural Comput.*
529 1997;9(8):1735–1780.
- 530 20. Ji S, Hölttä M, Marttinen P. Does the magic of BERT apply to medical code assignment?
531 A quantitative study. *Computers in Biology and Medicine.* 2021;139:104998.
- 532 21. He J, Cheng MX. Weighting Methods for Rare Event Identification From Imbalanced
533 Datasets. *Frontiers in big Data.* 2021;4:715320.
- 534 22. *Medical Dictionary for Regulatory Activities* [computer program].
- 535 23. Quezada SM, McLean LP, Cross RK. Adverse events in IBD therapy: the 2018 update.
536 *Expert Review of Gastroenterology & Hepatology.* 2018;12(12):1183-1191.
- 537 24. Anderson K, Moss K, Campbell B, Moote D, Kakazu K, Hyams JS. Follicular Dendritic Cell
538 Sarcoma in a Patient With Adolescent-Onset Crohn’s Disease Exposed to Multiple
539 Immunomodulator and Biologic Therapies. *JPGN Reports.* 2022;3(3):e231.
- 540 25. Smith M. Anti-TNF α therapy did not increase short-or medium-term risk for cancer in
541 patients with rheumatoid arthritis. *Annals of Internal Medicine.* 2010;152(20):JC5-13.
- 542 26. Conti F, Atzeni F, Massaro L, et al. The influence of comorbidities on the efficacy of
543 tumour necrosis factor inhibitors, and the effect of tumour necrosis factor inhibitors on
544 comorbidities in rheumatoid arthritis: report from a National Consensus Conference.
545 *Rheumatology.* 2018;57(Supplement_7):vii11-vii22.
- 546 27. *Statsmodels* [computer program]. 2009-2022.
547
548

549 **Figure and Table Legends**

550

551 **Figure 1.** Example of the three prediction tasks from a clinical note in the corpus. Medications of interest
552 were pre-annotated in blue, hospitalizations in red and signs and symptoms in yellow. Annotators
553 marked up the HPI section where medications of interest predated a hospitalization (green and blue
554 arrow) and AE causing hospitalization (red arrow). Task 3, also referred to as triple, is the combination of
555 Task 1 and Task 2.

556

557 **Table 1.** Characteristics of patients in our note corpus

558 **Table 2.** Distribution of number of notes containing an SAE

559 **Table 3.** Distribution of hospitalizations and SAEs by medication. Anti-tumor necrosis factor (Anti-TNF),
560 Janus Kinase-inhibitor (JAK-inhibitor), anti-interleukin-12/23 (anti-IL-12/23).

561

562 **Table 4.** Top 7 SAEs in the Study. Anti-tumor necrosis factor (Anti-TNF), Janus Kinase-inhibitor (JAKi),
563 anti-interleukin-12/23 (anti-IL-12/23).

564

565 **Figure 2.** Network Graph of SAEs by medication class. The width of lines indicates the strength of
566 association by frequency. The size of the nodes is relative to the number of exposures in our corpus to
567 each medication. SAE colors are indicative of which medication(s) they were associated with. An
568 interactive version of this figure can be found at <https://ibd-ade.streamlit.app/>

569

570 **Table 5.** Results of UCSF BERT performance on the tasks of SAE detection from real world clinical notes.
571 Results for the three relation tasks to classify whether a pair/triple of specific entities of type
572 medication, hospitalization and adverse event are related. Bolded models correspond to those with the
573 best performance as measured by Macro F1. Only nearby SAEs refer to restricting only SAEs that are
574 mentioned within a two-sentence window of the hospitalization event. Only the best three models are
575 reported. H-UCSF-BERT = Hierarchical University of California San Francisco Bidirectional Encoder
576 Representation from Transformers, TP = true positive, TN = true negative, FP = false positive and FN =
577 false negative.

578

579

Figures and Tables

Task 1: medication before hospitalization	Despite escalated adalimumab dosing (40 mg weekly and later 80mg Q2 weeks) she remains symptomatic. Her most recent colonoscopy on adalimumab 80mg Q2 weeks showed pancolitis. She was hospitalized with diarrhea with 20-30 BMs/day and dehydration.
Task 2: hospitalization for adverse event	Despite escalated adalimumab dosing (40 mg weekly and later 80mg Q2 weeks) she remains symptomatic. Her most recent colonoscopy on adalimumab 80mg Q2 weeks showed pancolitis. She was hospitalized with diarrhea with 20-30 BMs/day and dehydration.
Task 3: medication before hospitalization for adverse event (triple)	Despite escalated adalimumab dosing (40 mg weekly and later 80mg Q2 weeks) she remains symptomatic. Her most recent colonoscopy on adalimumab 80mg Q2 weeks showed pancolitis. She was hospitalized with diarrhea with 20-30 BMs/day and dehydration.

Figure 1. Example of the three prediction tasks from a clinical note in the corpus. Medications of interest were pre-annotated in blue, hospitalizations in red and signs and symptoms in yellow. Annotators marked up the HPI section where medications of interest predated a hospitalization (green and blue arrow) and AE causing hospitalization (red arrow). Task 3, also referred to as triple, is the combination of Task 1 and Task 2.

Patients (N=928)		Count(%)
Sex	Female	489 (52.7)
	Male	438 (47.2)
	Nonbinary	1 (0.1)
Age (years)	18-40	460 (49.6)
	41-60	313 (33.7)
	>60	155(16.7)
Ethnicity	Not Hispanic or Latino	832 (89.7)
	Hispanic or Latino	83 (8.9)
	Unknown/Declined	13 (1.4)

Race	White or Caucasian	670 (72.3)
	Asian	78 (8.4)
	Black or African American	41 (4.4)
	American Indian or Alaska Native	9 (1.0)
	Other Pacific Islander	1 (0.1)
	Other and Unknown/Declined	129 (13.9)
IBD Diagnosis	Crohn's disease	735 (79.2)
	Ulcerative colitis	625 (67.4)

Table 1. Characteristics of patients in our note corpus

	Train (%)	Development (%)	Test (%)
# Annotated notes	742 (80)	93 (10)	93 (10)
SAEs present	335 (82)	37 (9)	39 (9)
No SAEs	406 (79)	56 (11)	54 (10)

Table 2. Distribution of number of notes containing an SAE

Non-steroid Immunosuppressant Class	Non-steroid Immunosuppressant	Number of Hospitalizations	Number of SAEs
Anti-TNF	Adalimumab	136	172
	Certolizumab	19	26
	Etanercept	2	1
	Golimumab	2	2
	Infliximab	179	231
JAK-inhibitor	Tofacitinib	8	8
Anti-IL-12/23	Ustekinumab	100	128
Anti-integrin	Vedolizumab	106	135

Table 3. Distribution of hospitalizations and SAEs by medication. Anti-tumor necrosis factor (Anti-TNF), Janus Kinase-inhibitor (JAK-inhibitor), anti-interleukin-12/23 (anti-IL-12/23).

Serious Adverse Event MedDRA System Organ Class	Steroid Sparing Immunosuppressant Class	Frequency
Infections and infestations	Anti-TNF	66
	Anti-IL 12/23	24
	Anti-integrin	10
	JAKi	2
Failure of intended efficacy	Anti-TNF	216
	Anti-IL 12/23	75
	Anti-integrin	90
	JAKi	6
Gastrointestinal disorders	Anti-TNF	25
	Anti-IL 12/23	12
	Anti-integrin	6
	JAKi	0
Neoplasms	Anti-TNF	1
	Anti-IL 12/23	0
	Anti-integrin	0
	JAKi	0
Cardiac disorders	Anti-TNF	8
	Anti-IL 12/23	1
	Anti-integrin	3
	JAKi	0
General disorders and administration site conditions	Anti-TNF	21
	Anti-IL 12/23	1
	Anti-integrin	9
	JAKi	0
Nervous system disorders	Anti-TNF	7
	Anti-IL 12/23	2
	Anti-integrin	2
	JAKi	0

Table 4. Top 7 SAEs in the Study. Anti-tumor necrosis factor (Anti-TNF), Janus Kinase-inhibitor (JAKi), anti-interleukin-12/23 (anti-IL-12/23).

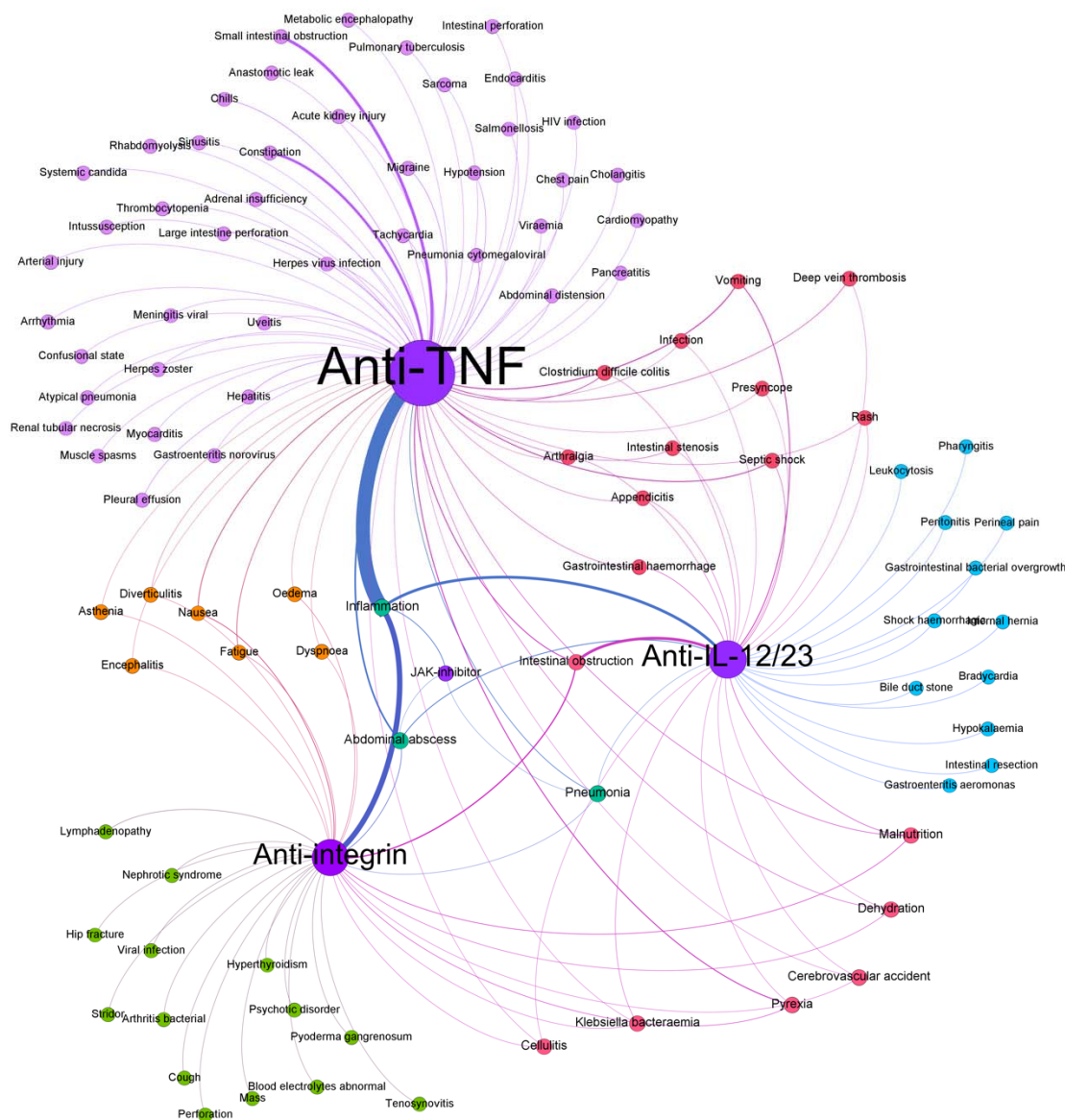


Figure 2. Network Graph of SAEs by medication class. The width of lines indicates the strength of association by frequency. The size of the nodes is relative to the number of exposures in our corpus to each medication. SAE colors are indicative of which medication(s) they were associated with. An interactive version of this figure can be found at <https://ibd-ade.streamlit.app/>

Task	Model	Accuracy (%)	Macro F1 (%)	TP	TN	FP	FN
Medication before hospitalization relations	H-UCSF-BERT	88	62	63	1989	173	105
	CNN	74	49	51	1682	480	117
	XGBoost	73	51	49	1672	490	119

Hospitalization for SAE relations	H-UCSF-BERT	96	62	79	9603	421	16
	H-UCSF-BERT + only nearby SAEs	92	68	34	1078	41	61
	BiLSTM + only nearby SAEs	93	48	7	1091	28	88
Medication before hospitalization for SAE relation (triples)	H-UCSF-BERT + only nearby SAEs	91	61	141	7790	619	178
	CNN + only nearby AEs	94	49	11	8013	396	308
	BiLSTM + only nearby AEs	95	50	11	7953	456	308

Table 5. Results of UCSF BERT performance on the tasks of SAE detection from real world clinical notes. Results for the three relation tasks to classify whether a pair/triple of specific entities of type medication, hospitalization and adverse event are related. Bolded models correspond to those with the best performance as measured by Macro F1. Only nearby SAEs refer to restricting only SAEs that are mentioned within a two-sentence window of the hospitalization event. Only the best three models are reported. H-UCSF-BERT = Hierarchical University of California San Francisco Bidirectional Encoder Representation from Transformers, TP = true positive, TN = true negative, FP = false positive and FN = false negative.