

1 **Title:** Identification of an ANCA-Associated Vasculitis Cohort Using Deep Learning and
2 Electronic Health Records

3 **Running head:** Deep Learning in ANCA-Associated Vasculitis Identification

4 **Authors:** Liqin Wang, PhD¹; John Novoa-Laurentiev, MS¹; Claire Cook, MPH²; Shruthi
5 Srivatsan, BA²; Yining Hua, MS³; Jie Yang, PhD⁴; Eli Miloslavsky, MD⁵; Hyon K. Choi,
6 MD, DrPH²; Li Zhou, MD, PhD¹, Zachary S. Wallace, MD, MSc²

7 **Affiliations:**

8 ¹Division of General Internal Medicine and Primary Care, Department of Medicine,
9 Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts,
10 USA

11 (*lwang@bwh.harvard.edu, jlaurentiev@bwh.harvard.edu, lzhou@bwh.harvard.edu*)

12 ² Rheumatology and Allergy Clinical Epidemiology Research Center and Division of
13 Rheumatology, Allergy, and Immunology, and Mongan Institute, Department of
14 Medicine, Massachusetts General Hospital, Boston, MA, USA

15 (*ccook13@mgh.harvard.edu, ssrivatsan1@mgh.harvard.edu, hchoi@mgh.harvard.edu,*
16 *zswallace@mgh.harvard.edu*)

17 ³Department of Epidemiology, Harvard T.H. Chan School of Public Health

18 (*yining_hua@hms.harvard.edu*)

19 ⁴Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine,
20 Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts,
21 USA (*jyang66@bwh.harvard.edu*)

22 ⁵ Division of Rheumatology, Allergy, and Immunology, Massachusetts General Hospital
23 and Harvard Medical School, Boston, MA, USA (*emiloslavsky@mgh.harvard.edu*)

24 **Funding support:** Authors of this work are supported by the National Institutes of
25 Health [grant numbers K23AR073334, R03AR078938, L30AR070520, R03AR078938,
26 K99AG075190]. The funders had no role in the design and conduct of the study;
27 collection, management, analysis, and interpretation of the data; preparation, review, or
28 approval of the manuscript; and decision to submit the manuscript for publication.

29 **Conflict of Interest:** ZSW reports research support from Amgen, Bristol-Myers Squibb,
30 and Principia/Sanofi, consulting fees from Amgen, Viela Bio, Horizon, Zenas
31 Biopharma, PPD, and MedPace, and has served on advisory boards for Amgen,
32 Horizon, Novartis, Sanofi, Shinogi, and Visterra/Otsuka. All other authors report no
33 disclosures.

34 **Corresponding Author:**

35 Liqin Wang, Ph.D.

36 Division of General Internal Medicine and Primary Care

37 Brigham and Women's Hospital

38 399 Revolution Drive, AR01-13-E58

39 Somerville, MA 02145, USA

40 Tel: 857-282-4089 | Fax: 857-282-5754 | E-mail: lwang@bwh.harvard.edu

41

42 **Content:**

43 **Word count:** Abstract 273; Manuscript 3899

44 **Figures:** 2

45 **Tables:** 4

46 **References:** 14

47 **ABSTRACT**

48 **Background:** ANCA-associated vasculitis (AAV) is a rare but serious disease.

49 Traditional case-identification methods using claims data can be time-intensive and may
50 miss important subgroups. We hypothesized that a deep learning model analyzing
51 electronic health records (EHR) can more accurately identify AAV cases.

52 **Methods:** We examined the Mass General Brigham (MGB) repository of clinical
53 documentation from 12/1/1979 to 5/11/2021, using expert-curated keywords and ICD
54 codes to identify a large cohort of potential AAV cases. Three labeled datasets (I, II, III)
55 were created, each containing note sections. We trained and evaluated a range of
56 machine learning and deep learning algorithms for note-level classification, using
57 metrics like positive predictive value (PPV), sensitivity, F-score, area under the receiver
58 operating characteristic curve (AUROC), and area under the precision and recall curve
59 (AUPRC). The deep learning model was further evaluated for its ability to classify AAV
60 cases at the patient-level, compared with rule-based algorithms in 2,000 randomly
61 chosen samples.

62 **Results:** Datasets I, II, and III comprised 6,000, 3,008, and 7,500 note sections,
63 respectively. Deep learning achieved the highest AUROC in all three datasets, with
64 scores of 0.983, 0.991, and 0.991. The deep learning approach also had among the
65 highest PPVs across the three datasets (0.941, 0.954, and 0.800, respectively). In a test
66 cohort of 2,000 cases, the deep learning model achieved a PPV of 0.262 and an
67 estimated sensitivity of 0.975. Compared to the best rule-based algorithm, the deep
68 learning model identified six additional AAV cases, representing 13% of the total.

69 **Conclusion:** The deep learning model effectively classifies clinical note sections for
70 AAV diagnosis. Its application to EHR notes can potentially uncover additional cases
71 missed by traditional rule-based methods.

72 **Keywords:** ANCA-associated vasculitis; Case Identification; Deep Learning; Machine
73 learning; Electronic Health Records; Clinical Notes

74
75 **SIGNIFICANCE AND INNOVATIONS:**

- 76 - Traditional approaches to identifying AAV cases for research have relied on
77 registries assembled through clinical care and/or on billing codes which may miss
78 important subgroups.
- 79 - Unstructured data entered as free text by clinicians document a patient's
80 diagnosis, symptoms, manifestations, and other features of their condition which
81 may be useful for identifying AAV cases
- 82 - We found that a deep learning approach can classify notes as being indicative of
83 AAV and, when applied at the case level, identifies more cases with AAV than
84 rule-based algorithms.

85 **INTRODUCTION**

86 Anti-neutrophil cytoplasmic autoantibody (ANCA)-associated vasculitis (AAV) is a rare
87 immune-mediated inflammatory disease associated with substantial morbidity, mortality,
88 and resource utilization.^{1, 2} The disease presents in patients with heterogeneous
89 manifestations (e.g., glomerulonephritis, sinusitis, skin rash, pulmonary nodules) to
90 clinicians from a variety of specialties (e.g., rheumatology, nephrology, otolaryngology,
91 intensive care) across the care spectrum within healthcare systems (e.g., community
92 hospital, tertiary care hospital, outpatient clinic, emergency departments). AAV case
93 identification in electronic health records (EHR), an increasingly important source for
94 epidemiologic research, is limited by a lack of well-performing methods to identify
95 cases.

96
97 To enable outcomes and comparative effectiveness studies using large, phenotypically
98 diverse cohorts from big data, a novel AAV case-finding algorithm is needed. Previous
99 studies have demonstrated that rule-based algorithms relying on ICD-9 codes for AAV
100 case identification have poor performance, partly because there is no specific ICD-9
101 code for microscopic polyangiitis (MPA), a subtype of AAV, and many MPA patients
102 may be miscoded using less specific ICD-9 codes. Additionally, the previously
103 developed algorithms, which require a positive ANCA test result or the use of multiple
104 ICD codes, may exclude important and informative subgroups of patients, including
105 ANCA-negative granulomatosis with polyangiitis (GPA) and those who die soon after
106 diagnosis from severe disease or complications. The performance of algorithms that
107 incorporate ICD-10 codes has not been previously assessed.

108

109 In addition to billing code data and test results, EHR data include unstructured data
110 entered as free text by clinicians documenting a patient's diagnosis, symptoms,
111 manifestations, and other features of their condition. We have previously demonstrated
112 that these notes can be leveraged to characterize the temporal course of AAV.³ Other
113 studies have suggested that unstructured data can enhance the performance of case-
114 finding algorithms for other conditions, but this remains underexplored for prototypic
115 rare conditions like AAV.⁴⁻⁷ Here, we hypothesized that machine learning methods could
116 be utilized to develop case-finding algorithms that accurately identify AAV patients and
117 that these algorithms would outperform or produce more phenotypically diverse cohorts
118 than rule-based algorithms.

119

120 **MATERIALS AND METHODS**

121 **Overview**

122 This study was conducted at Mass General Brigham (MGB), a large, integrated
123 healthcare delivery system in the Greater Boston area, Massachusetts. We used data
124 from MGB's research patient data registry (RPDR). The study was approved by the
125 MGB's Institutional Review Board (IRB number: 2016P000633). **Figure 1** illustrates the
126 overall process for the development of the AAV case-finding algorithms. We first
127 created a screening cohort of potential AAV cases. We then created three labeled
128 datasets from three cohorts for the development and evaluation of multiple machine
129 learning models for AAV case identification from unstructured clinical notes. The deep
130 learning model was further deployed to identify AAV patients from a random sub-cohort

131 of patients. The performance of this model was compared against rule-based
132 approaches.

133

134 **Study Cohorts**

135 The screening cohort was constructed using all data from RPDR, spanning its inception
136 on December 1, 1979, through May 11, 2021. We identified potential AAV patients
137 based on the presence of at least one AAV-related ICD code in a diagnosis field or the
138 use of a keyword in clinical notes. **Supplemental Table 1** lists keywords selected by the
139 subject-matter experts (EM and ZSW). To develop and evaluate the performance of
140 machine learning algorithms for AAV case identification, we used three distinct cohorts.
141 Cohort A comprised 700 patients with confirmed AAV, as they were previously identified
142 as part of the MGB AAV Cohort between January 01, 2002, and December 31, 2019.⁸
143 Cohorts B and C were random samples of 1000 patients each from the screening
144 cohort.

145

146 **Processing of Clinical Notes**

147 For each study cohort, we extracted all available clinical notes included in RPDR
148 database any time before July 23, 2021. These notes encompassed visit notes,
149 progress notes, ambulatory notes, history and physical exam notes, and discharge
150 summaries. Clinical notes contain rich information, such as clinical manifestations,
151 physical exams, and differential diagnoses, which are useful to determine whether
152 patient does or does not have AAV. However, due to the voluminous number of notes
153 that each patient accumulates over years of interaction with a healthcare system, the

154 model could miss true signals when handling large datasets. Thus, we processed each
155 note into smaller sections and then applied the models being tested to evaluate their
156 ability to predict whether or not the text refers to a diagnosis of AAV. To split the notes
157 into sections, we used Medical Text Extraction, Reasoning, and Mapping System
158 (MTERMS), an in-house developed natural language processing (NLP) system.⁹

159

160 **Definition of AAV Classification Task**

161 We approached the identification of AAV as a classification task, training models for
162 binary classification of note sections as pertaining to AAV. Throughout model
163 development, we assessed and compared the effectiveness of these models at the note
164 section level. In practical applications for patient identification, a patient was classified
165 as having AAV if any of their note sections were predicted positive, suggesting an AAV
166 diagnosis. Conversely, a patient was deemed not to have AAV if all their note sections
167 were predicted negative.

168

169 **Development of Labeled Datasets**

170 We developed three labeled datasets to train, test, and compare multiple machine
171 learning algorithms. Datasets I and II contained note sections that were derived from the
172 same population, i.e., patients with validated AAV (Cohort A) as well as patients with
173 possible AAV (Cohort B). To increase the positive case density in Dataset I, we applied
174 a list of expert-curated keywords to filter for sections that likely contain references to a
175 diagnosis of AAV or the presence of AAV manifestations (**Supplemental Table 2**). To
176 assess the generalizability of the trained models to note sections, regardless of keyword

177 presence, we established Dataset II by randomly selecting note sections from those not
178 included in Dataset I. Specifically, Dataset I included 5,000 sections from cohort A and
179 1,000 sections from cohort B, all mentioning specific keywords. Dataset II contains
180 2,000 sections from cohort A and 1,000 sections from cohort B, selected randomly
181 without considering keyword references. Dataset III consists of 7,500 note sections
182 randomly selected from Cohort C, a subset of the screening cohort, to evaluate model
183 performance in identifying potential AAV patients.

184

185 **Dataset Annotation**

186 Two subject matter experts (ZSW and CC) labeled each selected section from clinical
187 notes for whether it indicated that the diagnosis of AAV was present. Cases were
188 classified as AAV based on a prior algorithm for identifying AAV in epidemiologic
189 studies.¹⁰ In cases where there was limited data available to apply this algorithm, we
190 classified a case as AAV if the treating provider and both chart reviewers agreed on the
191 classification of the case as AAV. Patients with eosinophilic granulomatosis with
192 polyangiitis (EGPA) were classified as negative. Although EGPA is a specific type of
193 AAV, it presents a different etiology, pathology, and clinical features compared to other
194 AAV types, such as GPA and MPA.¹¹ The annotators first individually labeled 100
195 sections and any conflicts were resolved by consensus. Then in the second dataset of
196 100 sections, two annotators achieved near-perfect agreement with a Cohen's kappa of
197 0.897. The remaining note sections were each annotated by one of the annotators. Any
198 cases for which labeling was uncertain were resolved through consensus by ZSW and
199 CC.

200

201 **Model Development**

202 We first implemented four basic statistical machine learning algorithms, including
203 logistic regression, random forest, support vector machine (SVM), and XGBoost.¹² The
204 note sections were processed into n-grams (where n=1). Each section was converted
205 into term frequency-inverse document frequency vectors based on n-grams. The
206 algorithms were trained and tested with 5-fold cross validation using the Dataset I.

207

208 We also implemented a hierarchical attention-based deep neural network, which
209 includes a convolutional neural network for handling word variations (e.g., plural,
210 misspelling), a recurrent neural network for handling context (e.g., negation), and
211 attention layers for interpreting predictions. We chose this algorithm as it was previously
212 proved to be effective in allergic reaction detection from hospital safety reports⁵ and
213 cognitive decline detection from clinical notes.⁴ When implementing this deep learning
214 algorithm, each note section was treated as a sequence of tokens, and individual words
215 were represented by word embedding. We used pre-trained word embedding, named
216 BioWordVec, which is an open set of biomedical word vectors that integrated
217 biomedical text with Medical Subject Headings (MeSH) using the fastText model.¹³

218

219 Additionally, we implemented BioClinicalBERT, a domain-specific Bidirectional Encoder
220 Representation from Transformers (BERT) model,¹⁴ on the labeled dataset. BERT is
221 one of the most widely-used deep contextualized language models, achieving state-of-
222 the-art performance on various NLP tasks, including named entity recognition,

223 sentiment analysis, and question answering. We previously leveraged this algorithm to
224 identify patient gender identity in the EHR.⁷

225

226 **Evaluation for Model Generalizability**

227 To evaluate the generalizability of the models to the dataset regardless of keywords, we
228 applied the models trained in Dataset I to Dataset II. To evaluate the generalizability of
229 the models to the screening cohort, we applied the models trained in Datasets I and II to
230 Dataset III and reported the models' performance.

231

232 **Comparison of machine learning-based models with rule-based approaches**

233 To assess the feasibility of applying the deep learning model to identify AAV cases, we
234 compared its efficacy with that of two rule-based algorithms derived from administrative
235 claims data. Specifically, Rule 1 identifies patients who have any ICD-9 or ICD-10 codes
236 documented on at least two separated occasions. Rule 2 identifies patients with any
237 ICD-9 or ICD-10 codes recorded on at least two separated dates, and who also
238 received an AAV medication within a 6-month window (± 6 months) of the first ICD code
239 recorded. For the deep learning model, we used the optimal model to analyze the
240 clinical notes of Cohort D. The highest section-level prediction probability was
241 considered as the patient-level model prediction probability. Patients with a predicted
242 probability of 1 were classified as positive for AAV.

243

244 After identifying potential AAV cases using either the rule-based or the deep learning
245 model, we conducted a manual review of the EHR for these cases to pinpoint true

246 positive cases. From the cases not deemed to be AAV (i.e., those not identified by the
247 rules or the models), we randomly selected a subset (n=100) for manual chart review to
248 assess the false negative rate.

249

250 **Statistical Analysis**

251 We evaluated four statistical machine learning models, one deep learning model, and a
252 large language model for AAV case detection from clinical notes. Performance was
253 assessed based on the area under the receiver operating characteristic curve (AUROC),
254 the area under the precision-recall curve (AUPRC), positive predictive value (PPV),
255 sensitivity, and F-1 score which accounts for both precision and recall by taking the
256 harmonic mean. Both the AUROC and AUPRC were computed using the scikit-learn
257 Python library (scikit-learn Developers). We estimated the 95% confidence intervals (CI)
258 using 2,000 bootstrap iterations (Python, version 3.7; Python Software Foundation). To
259 compare the rule-based approaches with the top-performing AAV case identification
260 model, we computed the PPVs and sensitivities of all methods in Cohort D. Here the
261 total number of true positive patients was determined by adding those identified by both
262 the rule-based approaches and the top-performing AAV case identification model.

263

264 **RESULTS**

265 Cohort A, termed the 2002-2019 MGB AAV Cohort, included 700 PR3- or MPO-ANCA+
266 AAV patients. From these patients, 134,506 notes were extracted from the RPDR,
267 which included progress notes, ambulatory notes, and discharge summaries. These
268 notes were further split into 1,927,286 sections, of which 320,038 contained keywords.

269 Dataset I comprised 6,000 note sections, representing 5,765 notes from 1,638 patients
270 (968 [59.1%] female). Dataset II contained 3,008 sections from 2,970 notes,
271 representing 1,501 patients (885 [60.0%] female). Dataset III included 7,500 sections,
272 representing 5,429 notes from 1,000 patients (568 [56.8%] female) (**Table 1**). In Dataset
273 I, evidence of AAV was present in 2,669 sections (44.5%). Dataset II had 206 (6.8%)
274 sections positive for AAV. Out of the 3,008 sections in Dataset II, 457 contained at least
275 one keyword, with 203 (44.4%) of these containing evidence of AAV. Dataset III had 50
276 (0.67%) AAV-positive sections, and, of the 7,500 sections, 219 (2.92%) had one or
277 more keywords. Of those with keywords, 45 (20.5%) contained evidence of AAV.

278

279 The performance of the five models in each dataset is outlined in **Table 2**. The
280 hierarchical attention-based deep learning model demonstrated the best performance
281 on Dataset I during cross-validation, significantly outperforming other models, with an
282 AUROC of 0.983 (95% CI, 0.980-0.986) and an AUPRC of 0.977 (95% CI, 0.972-
283 0.982). Compared to the deep learning model, Bio_ClinicalBERT had slightly worse
284 performance in AUROC and AUPRC in Dataset I; however, it achieved better results in
285 precision, recall, and F-1 score.

286

287 Overall, all the models generalized well to Dataset II. The deep learning model exhibited
288 an AUROC of 0.991 (95% CI, 0.981-0.997) and an AUPRC of 0.962 (95% CI, 0.941-
289 0.980), with a 0.015 drop in AUPRC compared to its performance in Dataset I. Notably,
290 in Dataset II and among all the models, the XGBoost model achieved the best

291 performance in AUPRC, though the difference was not statistically significant, and
292 Bio_ClinicalBERT achieved the best sensitivity and F-1 score.

293
294 In Dataset III, the deep learning model outperformed other algorithms in all metrics.
295 Compared to its performance in Datasets I and II, it maintained a high AUROC of 0.991
296 (95% CI, 0.982-0.998); however, the AUPRC decreased to 0.760 (95% CI, 0.620-
297 0.885). Both Bio_ClinicalBERT and XGBoost saw greater decrease in performance from
298 Datasets I and II to Dataset III.

299
300 Among 2,000 patients from the screening cohort, Rule 1 identified 218 with two or more
301 AAV-related ICD codes (**Table 3**). After excluding 12 patients due to insufficient
302 information to ascertain AAV status, 40 (19.4%) were confirmed to have AAV. Rule 2
303 identified 52 patients meeting Rule 1 criteria and receiving a medication prescription
304 within 6 months of the first ICD code; 11 (21.2%) had confirmed AAV. Among the 2,000
305 patients, 1,977 had clinical notes reviewed using the deep learning model, which
306 predicted AAV in 177 patients with a probability of 1. After excluding 5 patients due to
307 insufficient information, 45 (26.2%) patients were confirmed to have AAV.

308
309 A review of 100 randomly selected cases, which were not predicted as AAV by either
310 method, confirmed the absence of AAV cases. If both the rules and the deep learning
311 algorithms identified all positive AAV cases among the 2000 cases reviewed, the total
312 number of positive cases amounted to 46. The estimated sensitivities for the deep
313 learning model, Rule 1, and Rule 2 were 97.5%, 87.0%, and 23.9%, respectively.

314
315 Significant differences were observed when comparing patients identified by the deep
316 learning model versus rule-based algorithms (**Table 3**). The deep learning model
317 identified a more ethnically diverse group (Hispanic: 24% vs 2.5% and 0%, respectively)
318 and more ANCA-negative AAV cases. The deep learning model found six additional
319 patients not identified by Rule 1, accounting for 13.0% of the positive cases, and Rule 1
320 found one not identified by the deep learning model.

321

322 **Error Analysis**

323 We analyzed the sections of clinical notes which the deep learning model falsely
324 predicted as consistent (false positive) with or not consistent (false negative) with AAV.

325 Many of the false positive errors can be grouped into three categories.

326 1. Ambiguous terms related to AAV. Some terms that appear in clinical notes
327 occasionally correspond to unrelated concepts with identical spelling. For example, the
328 abbreviation “MPA” might denote microscopic polyangiitis, which is pertinent to AAV, or
329 it could represent unrelated concepts like a multipurpose angiographic catheter or a
330 Master of Public Administration degree. Similarly, “GPA” can be used to abbreviate
331 grade point average.

332 2, Hypothetical scenarios in notes. Prediction errors also arose in note sections
333 described conjectural situations, such as guidelines of diagnosing AAV or potential
334 medication side effects (e.g., risks of anti-thyroid medications).

335 3. Notes detailing family history. The model made false positive predictions on notes
336 that mentioned a patient's family member being diagnosed with AAV, even though the
337 patient in question was not diagnosed.

338
339 False negative errors can be attributed to various representation of AAV-related
340 keywords: 1. Dictation errors or misspellings. Some false negative cases contain typos
341 of AAV-related terms or instances of terms transcribed incorrectly (e.g., "Wagner's"
342 instead of "Wegener's").

343 2. Combined terms. Certain terms related to AAV were mentioned as part of large
344 tokens, which might not be recognized by the deep learning model. For instance,
345 "ANCA+MPO+vasculitis" or "GPA/Wegener's" were treated as distinct or unrelated
346 compared to simpler terms like "ANCA+" or "GPA" or "Wegener's".

347 3. Rare variations of AAV-related terms, such as, "WEgeners", "WEGENER'S" or
348 "GRANULOMATOSIS", which might not be well recognized by the deep learning model.

349
350 When assessing model performance in 2,000 random patients from our screening
351 cohort, six patients were identified by the deep learning model but not by the rule-based
352 algorithm. After reviewing their charts, there are two potential reasons for these patients
353 were not captured by the rule-based algorithms. First, three patients were diagnosed
354 with AAV at institutions external to MGB so ICD codes for AAV were not used in
355 encounters in our healthcare system. Second, three patients had only one diagnosis
356 code, which did not meet the criteria of our rule-based algorithm. One patient identified
357 solely by the rule-based algorithm had positive ANCA pathology reports external to

358 MGB, which weren't included in the note screening. Available clinical notes lacked other
359 specific features of AAV in this case.

360

361 **DISCUSSION**

362 We found that a deep learning algorithm that integrated convolutional neural network,
363 recurrent neural network, and an attention mechanism trained using a small set of
364 keyword-identified, manually labeled note sections can be accurate and useful for
365 identifying a rare disease like AAV in a large cohort. The model performance in dataset I
366 showed its great capacity for detecting relevant signals from free-text narratives to make
367 accurate predictions. The model was generalizable to notes, regardless of the presence
368 of keywords. When applied to notes of patients from a large screening cohort for AAV
369 case identification, the deep learning model out-performed the traditional rule-based
370 algorithms which rely on ICD codes with or without medication prescriptions.

371

372 In addition to assessing the performance of the deep learning model, we also evaluated
373 the performance of rule-based algorithms using ICD codes and medication prescriptions
374 in our healthcare system. This is the first study that incorporates ICD-10 codes into an
375 assessment of performance of this rule-based algorithm. We found that these rule-
376 based algorithms had a PPV worse than that of the deep-learning model and that
377 incorporating medication prescriptions into the rule only slightly improved the PPV by
378 1.8%. In contrast, requiring a medication prescription significantly reduced sensitivity.
379 These observations speak to the need for innovative approaches, such as deep
380 learning, for developing new approaches for AAV case identification.

381
382 In comparing the deep learning approach with the ICD/medication-based rules, the
383 former demonstrated higher sensitivity and PPV. Examining clinical notes proved
384 beneficial in identifying additional cases, particularly those with a more remote history of
385 AAV or those diagnosed with AAV outside the MGB system. It also addressed cases
386 missed by the rule-based algorithm due to a limited number of ICD codes in the EHR.
387 This will be helpful for identifying patients, for instance, who have severe disease and
388 die during their initial admission for AAV. This is particularly crucial in rare diseases,
389 where even a small increase in sample size and including patients with the most severe
390 spectrum of disease can significantly impact studies. While clinical notes revealed only
391 6 additional cases in a sample of 2,000, after extrapolating these observations to the
392 entire screening cohort (n=88,902) we suspect that the deep learning model could
393 identify approximately 7,868 patients, with an estimated 2,000 of them having AAV.
394 Compared with a rule-based algorithm approach, the deep learning algorithms could
395 identify an additional 267 patients while reducing the need for extensive chart reviews
396 by more than 1,823 patients.

397
398 Compared with rule-based algorithms, we found that the deep learning model more
399 often identified patients of Hispanic background and those with ANCA-negative disease.
400 Why the deep learning model may have yielded a cohort with greater ethnic diversity is
401 unclear. One possibility has to do with differences in the way that ICD codes are used
402 for billing between people of different racial or ethnic backgrounds or because of the
403 way patients of different racial or ethnic backgrounds interact with the healthcare

404 system. The ability of the deep learning model to identify a greater proportion of cases
405 with ANCA-negative granulomatosis with polyangiitis is another strength of this
406 approach. This population is often excluded from observational studies of AAV as well
407 as clinical trials and the ability to identify them easily will facilitate research of this
408 subgroup.

409
410 Our findings suggest that applying a deep learning model may have benefit regarding
411 the efficiency of AAV cases identification. Rule-based approaches to AAV case finding
412 which identify potential AAV cases through billing codes with or without medications
413 typically necessitates a full chart review. In contrast, the deep learning model approach
414 presents a significant advantage because once the model flags sections that are
415 potentially related to AAV, only these specific sections typically require review,
416 potentially reducing the need for comprehensive chart evaluation.

417
418 Our study has several strengths. First, it was conducted in a large healthcare system
419 that includes both quaternary academic medical centers in addition to community
420 hospitals, primary care and specialty clinics, as well as specialty hospitals (e.g., ear,
421 nose throat hospital). Second, we applied four statistical machine learning models and
422 assessed their performance in comparison to commonly used rules-based algorithms.
423 Third, we assessed model performance using training and multiple validation datasets.

424
425 Despite these strengths, this study has certain limitations. First, we used data from a
426 single healthcare system so the model was not evaluated for its generalizability using

427 data from other institutions. This is an important next step in the development of deep
428 learning models to identify AAV cases. Second, the deep learning model was learned
429 from a small dataset, of which the vocabulary size might be relatively small. This might
430 affect the performance of the model when applied to a dataset with a larger vocabulary
431 size. Third, the current approach leveraged only clinical notes to identify potential AAV
432 cases. It is possible that including other data sources for model learning, such as lab
433 results, may improve the performance of algorithms for identifying AAV cases. Fourth,
434 we noted a large decline in the PPV when we applied the deep learning model, which
435 was trained and assessed at the level of sections from notes, to classify at the patient
436 level. This decrease can be attributed to the aggregation of errors from multiple note
437 sections per patient, which, when accumulated at the patient level, magnify the error
438 rate.

439
440 Our findings highlight the potential role of deep learning models for identifying positive
441 AAV cases from large screening cohorts. Moving forward, we intend to leverage our
442 deep learning model to screen the entire cohort, anticipating the identification of over
443 2,000 AAV cases. A cohort of this size would be substantially larger than the current
444 cohort assembled during a similar timeframe which includes fewer than 1,000 cases.
445 Thus, this represents a significant opportunity to expand the current MGB AAV cohort.
446 Furthermore, in the wake of the rise of sophisticated large language models, such as
447 GPT-4, an intriguing avenue of research would be to compare the performance of our
448 deep learning model with these state-of-the-art language models. The evolution of

449 natural language processing tools offers promising opportunities to further enhance the
450 accuracy and efficiency of clinical data mining and disease identification.

451 **CONCLUSION**

452 This study is the first to show that a deep learning algorithm can efficiently and
453 accurately identify cases of AAV, a prototypic rare condition, in part by only using
454 unstructured EHR data. This approach has the potential to identify cases that may be
455 overlooked if only using structured EHR data. This approach to case identification may
456 improve the spectrum of disease captured for observational studies and reduce the time
457 and resources often needed to review electronic health records. Future work will involve
458 deploying the model to screen a broader cohort for potential AAV patients and
459 assessing performance in other healthcare systems.

460 **ACKNOWLEDGEMENTS**

461

462

References

- 463 1. Kitching AR, Anders H-J, Basu N, Brouwer E, Gordon J, Jayne DR, et al. ANCA-
464 associated vasculitis. *Nature reviews Disease primers*. 2020;6(1):71.
- 465 2. Tan JA, Dehghan N, Chen W, Xie H, Esdaile JM, Avina-Zubieta JA. Mortality in
466 ANCA-associated vasculitis: ameta-analysis of observational studies. *Annals of the*
467 *rheumatic diseases*. 2017;76(9):1566-74.
- 468 3. Wang L, Miloslavsky E, Stone JH, Choi HK, Zhou L, Wallace ZS. Topic modeling
469 to characterize the natural history of ANCA-Associated vasculitis from clinical notes: A
470 proof of concept study. *Semin Arthritis Rheum*. 2021;51(1):150-7. PMID: 33383291.
- 471 4. Wang L, Laurentiev J, Yang J, Lo YC, Amariglio RE, Blacker D, et al.
472 Development and Validation of a Deep Learning Model for Earlier Detection of Cognitive
473 Decline From Clinical Notes in Electronic Health Records. *JAMA Netw Open*.
474 2021;4(11):e2135174. PMID: 34792589.
- 475 5. Yang J, Wang L, Phadke NA, Wickner PG, Mancini CM, Blumenthal KG, et al.
476 Development and Validation of a Deep Learning Model for Detection of Allergic
477 Reactions Using Safety Event Reports Across Hospitals. *JAMA Netw Open*.
478 2020;3(11):e2022836. PMID: 33196805.
- 479 6. Shao Y, Zeng QT, Chen KK, Shutes-David A, Thielke SM, Tsuang DW.
480 Detection of probable dementia cases in undiagnosed patients using structured and
481 unstructured electronic health records. *BMC Med Inform Decis Mak*. 2019;19(1):128.
482 PMID: 31288818.

- 483 7. Hua Y, Wang L, Nguyen V, Rieu-Werden M, McDowell A, Bates DW, et al. A
484 deep learning approach for transgender and gender diverse patient identification in
485 electronic health records. *J Biomed Inform.* 2023;147:104507. PMID: 37778672.
- 486 8. Wallace ZS, Fu X, Cook C, Ahola C, Williams Z, Doliner B, et al. Comparative
487 Effectiveness of Rituximab- Versus Cyclophosphamide-Based Remission Induction
488 Strategies in Antineutrophil Cytoplasmic Antibody-Associated Vasculitis for the Risk of
489 Kidney Failure and Mortality. *Arthritis Rheumatol.* 2023;75(9):1599-607. PMID:
490 37011036.
- 491 9. Zhou L, Plasek JM, Mahoney LM, Karipineni N, Chang F, Yan X, et al. Using
492 Medical Text Extraction, Reasoning and Mapping System (MTERMS) to process
493 medication information in outpatient clinical notes. *AMIA Annu Symp Proc.*
494 2011;2011:1639-48. PMID: 22195230.
- 495 10. Watts R, Lane S, Hanslik T, Hauser T, Hellmich B, Koldingsnes W, et al.
496 Development and validation of a consensus methodology for the classification of the
497 ANCA-associated vasculitides and polyarteritis nodosa for epidemiological studies. *Ann*
498 *Rheum Dis.* 2007;66(2):222-7. PMID: 16901958.
- 499 11. Kitching AR, Anders HJ, Basu N, Brouwer E, Gordon J, Jayne DR, et al. ANCA-
500 associated vasculitis. *Nat Rev Dis Primers.* 2020;6(1):71. PMID: 32855422.
- 501 12. Chen T, Guestrin C, editors. Xgboost: A scalable tree boosting system.
502 *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery*
503 *and data mining; 2016.*
- 504 13. Zhang Y, Chen Q, Yang Z, Lin H, Lu Z. BioWordVec, improving biomedical word
505 embeddings with subword information and MeSH. *Scientific data.* 2019;6(1):52.

506 14. Devlin J, Chang M-W, Lee K, Toutanova K. Bert: Pre-training of deep
507 bidirectional transformers for language understanding. arXiv preprint arXiv:181004805.
508 2018.

509

510

511 **Table 1.** Characteristics of the datasets I, II and III for the development and validation of models
512 for identifying evidence of ANCA vasculitis.

Characteristic	Dataset I	Dataset II	Dataset III
Sections, n	6,000	3,008	7,500
Notes, n	5,765	2,970	5,429
Character length per section, mean (SD)	839 (793)	443 (566)	410 (551)
Unique patients, n	1,638	1,501	1,000
Female, n (%)	968 (59.1)	885 (60.0)	568 (56.8)
Keyword present, n (%)	6,000 (100)	457 (15.2)	219 (2.9)
Sections consistent with a diagnosis of AAV, n (%)	2,669 (44.5)	206 (6.8)	50 (0.67)
Sections consistent with AAV diagnosis and containing a keyword of interest, n (%)	2,669 (44.5)	203 (44.4)	45 (20.5)

513 Abbreviations: AAV, ANCA-associated vasculitis; SD, standard deviation

514

Table 2. Performance of four machine learning models for detecting AAV from clinical notes

Model	AUROC (95% CI)	AUPRC (95% CI)	PPV	Sensitivity	F-1 Score
Dataset I (6,000 Note Sections)					
Logistic Regression	0.923 (0.916-0.930)	0.892 (0.879-0.904)	0.857	0.777	0.815
Random Forest	0.929 (0.922-0.935)	0.888 (0.874-0.902)	0.808	0.886	0.845
SVM	0.938 (0.932-0.944)	0.912 (0.901-0.923)	0.848	0.862	0.855
XGBoost	0.957 (0.952-0.962)	0.939 (0.929-0.948)	0.916	0.896	0.906
Bio_ClinicalBERT	0.957 (0.945-0.968)	0.962 (0.950- 0.972)	0.947	0.956	0.952
Deep Learning	0.983 (0.980-0.986)	0.977 (0.972-0.982)	0.941	0.951	0.946
Dataset II (3,008 Note Sections)					
Logistic Regression	0.981 (0.969-0.991)	0.893 (0.858-0.925)	0.709	0.874	0.783
Random Forest	0.983 (0.976-0.989)	0.871 (0.832-0.904)	0.528	0.922	0.671
SVM	0.983 (0.971-0.991)	0.910 (0.880-0.939)	0.645	0.908	0.754
XGBoost	0.990 (0.981-0.997)	0.963 (0.941-0.981)	0.886	0.947	0.915
Bio_ClinicalBERT	0.981 (0.967-0.992)	0.954 (0.933-0.973)	0.939	0.966	0.952
Deep Learning	0.991 (0.981-0.997)	0.962 (0.941-0.980)	0.954	0.898	0.925
Dataset III (7,500 Note Sections)					
Logistic Regression	0.940 (0.906-0.967)	0.390 (0.259-0.534)	0.109	0.720	0.189
Random Forest	0.975 (0.962-0.986)	0.392 (0.258-0.529)	0.064	0.900	0.120
SVM	0.907 (0.861-0.948)	0.402 (0.264-0.538)	0.080	0.700	0.143
XGBoost	0.982 (0.955-0.996)	0.490 (0.359-0.633)	0.320	0.800	0.457
Bio_ClinicalBERT	0.857 (0.792-0.917)	0.570 (0.473-0.659)	0.419	0.720	0.529
Deep Learning	0.991 (0.982-0.998)	0.760 (0.620-0.885)	0.800	0.800	0.800

515 Abbreviations: SVM, support vector machine; AUROC, area under the receiver operating
516 characteristic curve; AUPRC, the area under the precision-recall curve; PPV, positive predictive
517 value. Bold highlights the best performing model according to each measure.

518

519

520

521 **Table 3.** Characteristics of patients identified identified by the deep learning model and two rule-
 522 based algorithms among 2,000 random sample of the screening cohort.

	Deep Learning	Rule 1	Rule 2
N	177	218	52
AAV	45	40	11
Age at diagnosis, years	N = 43 (2 unknown)	N = 38 (2 unknown)	N = 11
Mean (SD)	54.79 (20.61)	55.97 (20.27)	56.64 (19.74)
Median (IQR)	60 (37 – 72)	60 (37.25 – 73.75)	60 (44.5 – 70.5)
Sex, female, n (%)	35 (77.78)	31 (77.50)	9 (81.82)
Race, n (%)			
White	38 (84.4)	33 (82.50)	10 (90.91)
Black	1 (2.2)	1 (2.50)	0
Asian	1 (2.2)	1 (2.50)	0
Other	1 (2.22%)	1 (2.50)	0
Unavailable	3 (6.67%)	3 (7.50)	0
Declined	1 (2.22%)	1 (2.50)	1 (9.09)
Ethnicity, n (%)			
Not Hispanic	33 (73.33%)	29 (72.5)	8 (72.73%)
Hispanic	11 (24.44%)	1 (2.5)	0
Unavailable	1 (2.22%)	10 (25.00%)	3 (27.27%)
ANCA Positive, n (%)	40 (88.89%) (1 unknown)	37 (92.50%)	11 (100%)
ANCA Type, n (%)			
MPO	24 (53.33%)	23 (57.50%)	7 (63.64%)
PR3	14 (31.11%)	13 (32.50%)	3 (27.27%)
Neither	4 (8.89%)	3 (7.50%)	0
Unknown	3 (6.67%)	1 (2.50%)	1 (9.09%)

523 Abbreviation: AAV, ANCA-associated vasculitis; SD, standard deviation; IQR, Interquartile
 524 range; MPO, Myeloperoxidase; PR3, proteinase 3.

525
 526

527 **Figure Legends**

528 **Figure 1.** Two-Phase Process for Identifying ANCA-Associated Vasculitis Cases. Phase
529 1 entails dataset creation along with the training and evaluation of models. Phase 2
530 compared the performance of the deep learning model with two rule-based algorithms.

531

532

533 **Figure 2.** Performance of the Machine Learning and Deep Learning Algorithms on
534 Datasets I, II, and III. A. Precision-recall curves for Dataset I. B. Receiver operating
535 characteristic (ROC) curves for Dataset I. C. Precision-recall curves for Dataset II. D.
536 ROC curves for Dataset II. E. Precision-recall curves for Dataset III. F. ROC curves for
537 Dataset III.

538

539

540

Screening Cohort: Patients identified from the entire RPDR database by ICD codes and/or keywords
(n = 88,902)

Cohort A
Validated patients with AAV
(n = 700)

Cohort B
Randomly selected patients
(n = 1000)

Cohort C
Randomly selected patients
(n = 1000)

Cohort D
Randomly selected patients
(n = 2000)

Sectionization of clinical notes

Sectionization of clinical notes

Sectionization of clinical notes

Identification of sections with keywords

Randomization of sections

5000 sections from **Cohort A** and 1000 sections from **Cohort B**

2000 sections from **Cohort A** and 1008 sections from **Cohort B**

7500 sections from **Cohort C**

Annotation of individual sections for AAV

Dataset I

Dataset II

Dataset III

Feature engineering

Model development and evaluation

Model generalizability evaluation (regardless of keywords)

Model generalizability evaluation (on a subset of testing cohort)

Apply deep learning model to screen clinical notes

Apply rule-based algorithms to screen diagnosis and medication

Statistical analysis

Phase I: Development and Validation of Deep Learning Models for Section Classification

Phase II: AAV Case Identification

medRxiv preprint doi: <https://doi.org/10.1101/2024.06.09.24308603>; this version posted June 10, 2024. 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-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

