

Original Paper

# Patient Phenotyping for Atopic Dermatitis With Transformers and Machine Learning: Algorithm Development and Validation Study

Andrew Wang<sup>1</sup>, BSc; Rachel Fulton<sup>2</sup>, MD; Sy Hwang<sup>1</sup>, MS; David J Margolis<sup>1\*</sup>, MD, PhD; Danielle Mowery<sup>1\*</sup>, MS, PhD

<sup>1</sup>University of Pennsylvania, Philadelphia, PA, United States

<sup>2</sup>Lankenau Medical Center, Wynnewood, PA, United States

\*these authors contributed equally

**Corresponding Author:**

Danielle Mowery, MS, PhD

University of Pennsylvania

A206 Richards Building

3700 Hamilton Walk

Philadelphia, PA, 19104

United States

Phone: 1 2157466677

Email: [dlmowery@pennmedicine.upenn.edu](mailto:dlmowery@pennmedicine.upenn.edu)

## Abstract

**Background:** Atopic dermatitis (AD) is a chronic skin condition that millions of people around the world live with each day. Performing research into identifying the causes and treatment for this disease has great potential to provide benefits for these individuals. However, AD clinical trial recruitment is not a trivial task due to the variance in diagnostic precision and phenotypic definitions leveraged by different clinicians, as well as the time spent finding, recruiting, and enrolling patients by clinicians to become study participants. Thus, there is a need for automatic and effective patient phenotyping for cohort recruitment.

**Objective:** This study aims to present an approach for identifying patients whose electronic health records suggest that they may have AD.

**Methods:** We created a vectorized representation of each patient and trained various supervised machine learning methods to classify when a patient has AD. Each patient is represented by a vector of either probabilities or binary values, where each value indicates whether they meet a different criteria for AD diagnosis.

**Results:** The most accurate AD classifier performed with a class-balanced accuracy of 0.8036, a precision of 0.8400, and a recall of 0.7500 when using XGBoost (Extreme Gradient Boosting).

**Conclusions:** Creating an automated approach for identifying patient cohorts has the potential to accelerate, standardize, and automate the process of patient recruitment for AD studies; therefore, reducing clinician burden and informing the discovery of better treatment options for AD.

(*JMIR Form Res* 2024;8:e52200) doi: [10.2196/52200](https://doi.org/10.2196/52200)

## KEYWORDS

atopic dermatitis; classification; classifier; dermatitis; dermatology; EHR; electronic health record; health records; health; informatics; machine learning; natural language processing; NLP; patient phenotyping; phenotype; skin; transformer; transformers

## Introduction

### Background

Atopic dermatitis (AD) is a common skin disease with a population prevalence of approximately 30% [1]. It is often diagnosed in early childhood, but onset can occur at any age [2-5]. Symptoms of AD include inflamed, red, irritated, and

itchy skin and can cause significant physical and emotional distress. AD is often associated with other allergic illnesses, including asthma, seasonal allergies, and food allergies [2,3,5-7].

AD is thought to be associated with skin barrier dysfunction and immune dysregulation [5]. AD has also been associated with genetic variation as well as environmental factors [5]. Classic treatment for AD has included the use of moisturizers, topical steroids, and other topical anti-inflammatory agents [8].

However, in the past few years, there have been significant treatment advances, which include systemic agents that alter immune function, such as dupilumab. Therefore, due to the widespread nature of AD, the need for improved knowledge of the natural history of AD, the need to understand the efficacy of new treatments, and the need to develop new treatments, there is an urgent need to understand the clinical course of individuals with AD. However, identifying appropriate cohorts of patients for medical studies can be difficult and time-consuming. Because AD is so common as well as being diagnosed and managed by many different clinicians in varying health care settings, a potential source population would be patients from a health system's electronic health records (EHRs) [9]. Investigators often ascertain a patient's illness using International Classification of Disease (ICD) hospital billing codes as recorded during routine office visits. However, it has been previously demonstrated that reliance on ICD codes is not an accurate method for the ascertainment of study cohorts with AD [9,10]. Furthermore, epidemiologic studies have used different methods and algorithms, including the UK Working Party (UKWP) diagnostic criteria and the Hanifin and Rajka (HR) criteria [11,12]. Investigators attempting to conduct clinical trials and observational studies have also relied on manual, large-scale chart review, a process that is inefficient, slow, and tedious [9]. This motivates the need for a standard method to accurately, automatically, and efficiently identify potential patient cohorts from their text medical records by using natural language processing (NLP) and machine learning (ML) techniques.

### Previous Work

Previously, researchers aimed to phenotype patients with AD using EHR data. In particular, Gustafson et al [10] trained a logistic regression model with lasso regularization to identify cases of AD from the Northwestern Medical Enterprise Data Warehouse, which contained both structured data (ICD Ninth and Tenth Revision codes, medication prescriptions, and laboratory results) as well as unstructured data (clinician notes from patient encounters). A gold standard diagnosis was assigned to each patient in their data set by 2 rheumatologists following a chart review when using the UKWP criteria and (alternatively) when using the HR criteria.

Although similar, this study differs in the following ways: (1) we survey a wide range of supervised ML algorithms as opposed to only using lasso regularized logistic regression, (2) we use transformer embeddings of sentences to represent information in each patient's records and aggregate these embeddings with multilayer perceptron (MLP) networks to create a patient vector representation for patient phenotyping, and (3) we performed an ablation study of processing methods to compare the impact on performance in using a probability-based versus binary label of whether each patient meets various AD diagnostic criteria when creating a vector to represent each patient for input to our final patient phenotyping algorithms.

### Contributions

The primary contributions of this study are as follows:

- We introduce and validate a rules-based approach for aggregating information from patient EHR data to generate binary-valued patient vectors that are used with standard ML algorithms for patient phenotyping.
- We introduce and validate a transformer-based approach for aggregating information and patient phenotyping by using "Bidirectional Encoder Representations from Transformers" (BERT) models (ie, BERT Base Uncased and BioClinical BERT) to generate patient vectors of probabilities, which are used with standard ML algorithms for patient phenotyping.
- We compare the aforementioned approaches to (1) discern whether a transformer model pretrained on clinical text can provide performance benefits over a transformer model not pretrained on clinical text, and (2) discern whether a transformer-based approach for aggregating information could outperform a rules-based approach for aggregating information.
- We demonstrate that MLP networks can be used with BERT sentence embeddings to identify which sentences in patient records are relevant to the diagnosis of AD. These MLP networks can then be used during clinician chart review to highlight sentences that are relevant to diagnosis and therefore accelerate the process of chart review during clinical trial recruitment.

## Methods

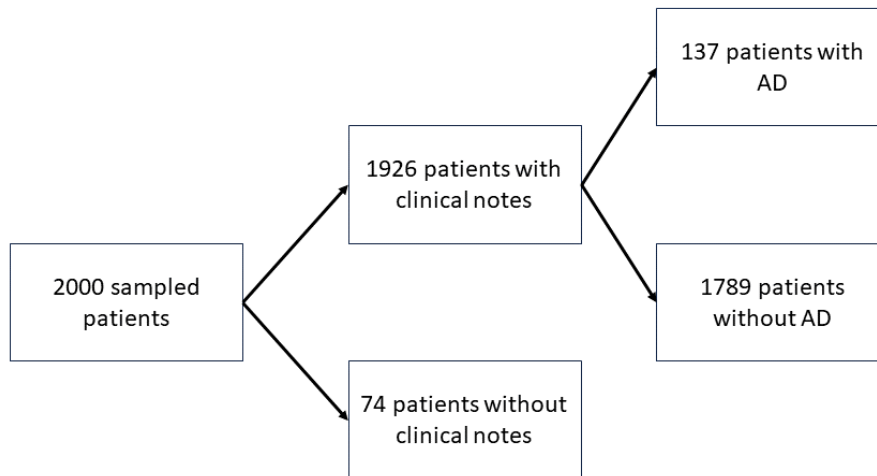
### Overview

To predict whether a patient may qualify as a participant for an AD study based on their EHR, we first assigned patients in our data set to either the training or testing sets. Then, for each patient, we aggregated the text from their EHR and constructed a vector representation of clinical features indicative of AD according to the UKWP criteria. Lastly, we leveraged our vectorized patient representations to train several ML classifiers to predict whether each patient has AD. In the following sections, we detail this process.

### Data Set Creation

We initially sampled 2000 patients and their clinical records from Epic Clarity, Penn Medicine's EHR database. We selected Penn Medicine patients who were diagnosed with a subset of AD-related ICD codes [9]. As shown in Figure 1, of the 2000 sampled patients, we identified 1926 patients who had clinical notes for processing. We then deidentified these patient records according to the Safe Harbor method using the "Protected Health Information filter" (Philter) [13]. Each patient in the data set was also manually reviewed and labeled according to the UKWP diagnostic criteria for AD. According to the UKWP criteria, in order to qualify as having AD, a patient must have an itchy skin condition along with 3 or more of the following: a history of flexural involvement, a history of asthma or hay fever, a history of dry skin, an onset of rash when aged 2 years or younger, or a visible flexural dermatitis. Our data set was validated by 2 clinicians (a board-certified dermatologist [DJM] and a medical fellow [RF]), resulting in 137 patients with AD and 1789 patients without AD.

**Figure 1.** Waterfall diagram of cohort. AD: atopic dermatitis.



**Training and Testing Split**

We first created our training set. Due to the heavy class imbalance in our data set, we decided to create a balanced training set to prevent biasing the model toward either patients with AD or patients without AD. In particular, we created the training set by assigning 80% (109/137) of the 137 patients with AD to our training set and undersampling the patients without AD to match the number of patients with AD. The remaining 20% (28/137) of the 137 patients were assigned to both of our testing sets. This resulted in a training set that had 109 patients with AD and 109 patients without AD.

Next, we created 2 testing sets. The first testing set was class-balanced and was intended to show how our patient classification model can generalize to unseen samples if the class distribution is kept the same. The second testing set was class-imbalanced (28/91, 30% of patients with AD and 63/91, 70% of patients without AD) and was intended to show how our patient classification model can perform when the class-distribution of the data set matches the prevalence of AD in the United States.

We created the first (balanced) testing set by including the 20% (28/137; previously reserved for testing) of the 137 patients with AD and combining them with an equal number of patients without AD who have not been used during training. This resulted in a balanced testing set that had 28 patients with AD and 28 patients without AD.

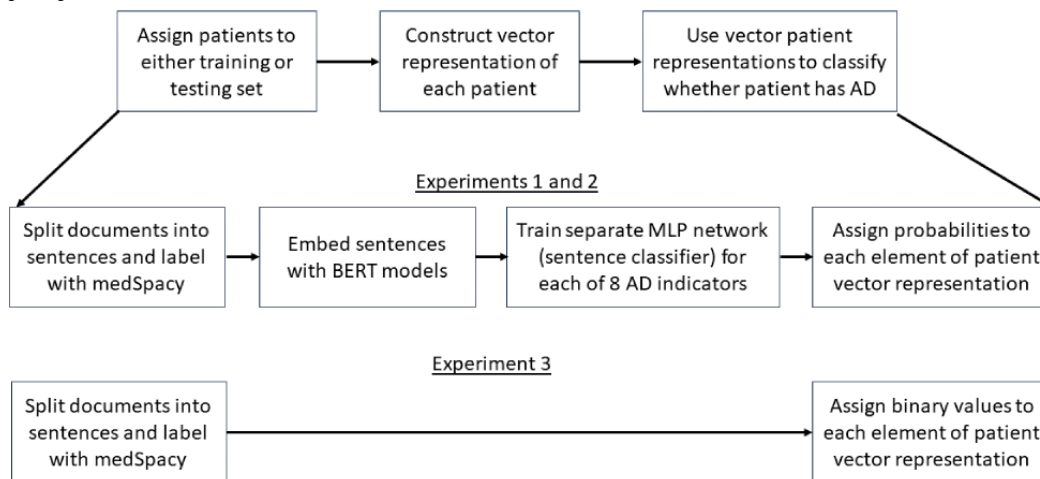
Furthermore, we created the second (unbalanced) testing set by including the same 20% (28/137) patients with AD but instead combining them with a greater number of patients without AD to match the 30% prevalence rate of AD found in the United States [1]. This resulted in an unbalanced testing set with 28 patients who have AD and 63 patients without AD.

We chose not to create a separate hyperparameter tuning set and instead applied cross-validation for hyperparameter tuning on the training set due to the data-scarce setting of our experiments.

**Vector Representation for AD Classification**

Next, we created a vector representation for each patient. We performed 3 experiments to compare different methods of creating each patient’s vector representation (Figure 2).

**Figure 2.** Atopic dermatitis (AD) phenotyping pipeline across all 3 experiments. BERT: Bidirectional Encoder Representations from Transformers; MLP: multilayer perceptron.



### Description of Patient Vector Representation

Each patient's vector representation is 8 elements long, where each element of the vector is representative of whether the patient fulfills a different AD diagnosis criteria based on the UKWP criteria as well as clinician feedback (Table 1). Across all 3 experiments, each element in the patient vector corresponds to a distinct classification task; however, in experiments 1 and 2, each element is a probability, and in experiment 3, each element is a binary value.

In experiments 1 and 2, elements 1-8 of each patient's vector represent the highest probability that any sentence in the patient's EHR mentions (1) AD or synonyms of AD, (2) keywords that suggest hay fever allergies, (3) keywords that suggest atopic allergies, (4) keywords that suggest eczema or rashes, (5) keywords that indicate dry or itchy skin, (6) keywords denoting nonasthma medications, (7) keywords suggesting the presence of asthma, and (8) keywords indicating the use of asthma medications.

In experiment 3, instead of each element representing a probability, each element represents a binary value of whether there was at least 1 sentence in the corresponding patient record suggesting the presence of the corresponding AD indicator.

In the first 2 experiments, each patient's vector elements represent probabilities (ranging from 0 to 1). Each probability value is derived from a distinct MLP classifier. Experiments 1

and 2 were performed to compare the use of 2 BERT models (BERT Base Uncased [14,15] in experiment 1 and BioClinical BERT [16,17] in experiment 2) for creating sentence embeddings used to train MLP networks (or alternatively, sentence classifiers). A separate MLP network is trained for each element of the patient vector. Each MLP network is trained to distinguish sentences in 1 of the 8 AD indicator categories from sentences in all other categories. Furthermore, *medSpacy* (Eyre et al [18]) was used to split documents into sentences and label each sentence with different categories. After each sentence classifier is trained, embeddings of all sentences in each patient's full EHR are passed through each sentence classifier, and an aggregation function (max operator) is used to assign a value to each element of each patient's vector. Our goal in experiments 1 and 2 was to test the hypothesis that a BERT model pretrained on clinical text (BioClinical BERT) could outperform a BERT model trained on nonclinical text (BERT Base Uncased).

In experiment 3, each patient's vector elements are binary (either 0 or 1). Each element corresponds to a diagnostic criterion and represents whether *medSpacy* was able to identify at least 1 sentence in the patient's record with a keyword and affirming context that suggests the patient meets the corresponding diagnostic criteria. Our goal was to conduct an ablation study to test the hypothesis that an AD phenotyping classifier leveraging BERT embeddings to create the patient vector representation will better discern whether a patient has AD than an AD Phenotyping Classifier without BERT embeddings.

**Table 1.** Meaning of each patient vector element.

Element	AD <sup>a</sup> indicator (diagnostic criteria)
1	EHR <sup>b</sup> directly mentions patient has AD
2	Patient has hay fever allergies
3	Patient has atopic allergies
4	Patient has eczema or rashes
5	Patient has dry or itchy skin
6	Patient uses nonasthma medications related to treating AD
7	Patient has asthma
8	Patient uses asthma medications

<sup>a</sup>AD: atopic dermatitis.

<sup>b</sup>EHR: electronic health record.

### Preprocessing for Experiments 1-3

Before each experiment, we applied the same preprocessing steps to assign 1 or more labels to each sentence in our corpus of documents in both our training and testing sets. Each sentence can be labeled as applying to 1, multiple, or none of the 8 AD indicators previously defined.

For each of the 8 diagnostic criteria, we first created a list of keywords and phrases (for each vector element) that suggested the presence of the corresponding diagnostic criteria. Next, we used *medSpacy* with the ConText (Harkema et al [20]) algorithm to split each document into sentences and categorize each sentence [18]. Using *medSpacy* allows us to obtain sentences

that suggest the presence of each of the 8 diagnostic criteria due to *medSpacy*'s use of regex and rules-based keyword matching. Furthermore, *medSpacy*'s implementation of the ConText algorithm allows us to discern between sentences that affirm from negated assertions. We define negated sentences for each AD indicator as sentences where the indicator is ruled out, sentences where the indicator is experienced by someone other than the patient, and sentences where the existence of the indicator is hypothetical [19-22].

After assigning 1 or more categorical labels to each sentence with *medSpacy*, we then performed 3 different experiments to create a vectorized representation of each patient.

In Tables 2 and 3, we include some statistics on the data set obtained after preprocessing.

As shown in Table 2, patients with AD have approximately twice as many sentences as patients without AD. The average number of documents and sentences is the same (within patients with AD and similarly within patients without AD) between BERT Base Uncased and BioClinical BERT experiments because these values are only dependent on *medSpacy*'s preprocessing of documents. Furthermore, using BioClinical BERT to tokenize sentences tends to yield more tokens (on average) per patient and per document. We hypothesize this is because the BioClinical BERT tokenizer is able to recognize more clinical terms and therefore yields more tokens for the same sentence than using the tokenizer from BERT Base Uncased.

As shown in Table 3, sentences in category 5 (relating to dry or itchy skin) tend to have the most tokens, whereas sentences in category 6 (relating to the use of nonasthma medications related to treating AD) tend to have the least number of tokens. We hypothesize that this is because categories where the average number of tokens per sentence is greater tend to correspond to more general categories where many terms and sentences could apply, whereas categories where the average number of tokens per sentence is lower tend to correspond to more specific categories, thus yielding a lower average number of tokens per sentence. Additionally, similarly to before, we can see that using BioClinical BERT tends to result in a greater number of tokens per sentence than using BERT Base Uncased for the same sentence.

**Table 2.** Differences in the number of documents, sentences, and tokens between patients with atopic dermatitis (AD) and those without AD.

	Patients with AD		Patients without AD	
	BERT <sup>a</sup> Uncased	BioClinical BERT	BERT Uncased	BioClinical BERT
Average number of documents (per patient)	23.44	23.44	7.99	7.99
Average number of sentences (per patient)	392.99	392.99	193.69	193.69
Average number of tokens (per patient)	16035.39	17054.11	7241.02	7674.35
Average number of sentences (per document)	16.77	16.77	24.25	24.25
Average number of tokens (per document)	684.16	727.63	906.45	960.69
Average number of tokens (per sentence)	40.80	43.40	37.38	39.62

<sup>a</sup>BERT: Bidirectional Encoder Representations from Transformers.

**Table 3.** Mean number of tokens for sentences identified in each category.

	BERT <sup>a</sup> Uncased (tokens per sentence), mean	BioClinical BERT (tokens per sentence), mean
Category 1	99.49	106.16
Category 2	81.18	92.41
Category 3	79.20	82.07
Category 4	83.74	92.55
Category 5	106.64	112.58
Category 6	74.93	80.17
Category 7	92.85	109.40
Category 8	76.13	83.57

<sup>a</sup>BERT: Bidirectional Encoder Representations from Transformers.

### **Experiments 1 and 2: Patient Vector Construction With BERT Embeddings**

In experiments 1 and 2, we first used the sentences *medSpacy* identified in each category to create class-balanced training and testing sets for each MLP network classifier, as shown in Table 4. The same training and testing set was used for both experiment 1 (BioClinical BERT) and experiment 2 (BERT Base Uncased).

Next, we used pretrained BERT models to generate embeddings of the sentences in each classifier's training and testing set. We

incorporated pretrained BERT models because these models have been trained on a much larger corpus than our existing data set, and BERT provides a context-sensitive embedding of text that other techniques, such as bag of words, do not provide. Furthermore, we used BERT Base Uncased in experiment 1 and Alsentzer et al's [16] BioClinical BERT in experiment 2 because we wanted to quantify how much of a difference in performance using a model pretrained on clinical text can provide over a model that has not been pretrained on clinical text.



Using these embeddings, we trained a MLP network to distinguish sentence embeddings in each category from sentence embeddings that are not in the corresponding category. Each of our MLPs was trained with the following architecture: a fully connected input layer of shape  $768 \times 100$ , followed by a Rectified Linear Unit (ReLU) activation, further followed by a fully connected output layer of shape  $100 \times 2$ . We trained each of our MLPs for 10 epochs with the cross-entropy loss function, the stochastic gradient descent (SGD) optimizer, a learning rate of 0.001, and a momentum value of 0.9. The final layer of each MLP can then be used to obtain the probability that any given sentence embedding comes from the category for which the MLP is being trained by passing the logits of the final layer to the softmax function.

We used the ReLU activation function as defined below, where  $x$  is the input to the ReLU function:

$$ReLU(x) = \max(0, x)$$

We also used the softmax function as defined below, where  $e$  is the standard exponential function and  $\vec{x}_i$  is element at index  $i$  of the  $K$  element long input vector  $\vec{x}$ .

$$softmax(\vec{x})_i = \frac{e^{x_i}}{\sum_{i=1}^K e^{x_i}}$$

We chose to embed our sentences once with pretrained BERT models and then feed these saved embeddings to our MLP networks as opposed to adding a classification head (a linear

layer) to the end of our pretrained BERT models. Although doing so only allows us to fine-tune the weights in our MLP network (as opposed to also fine-tuning the weights BERT uses to embed the sentences), doing so allows us to iterate over different experiments more quickly and with less computational power. In particular, we are able to (1) avoid the large computational expense of gradient calculations during backpropagation for all 12 layers of transformers used by BERT when fine-tuning the model, (2) avoid the computational expense of repeatedly generating the same embeddings from BERT multiple times (if we choose to freeze the weights of BERT and only fine-tune an added classification head or linear layer), and (3) iterate more efficiently over different hyperparameter combinations across different experiments with our MLP networks.

After training a separate MLP network for each of the 8 categories, we generated a vector representation for each patient, where each of the 8 vector elements represents the highest probability that any given sentence in the patient record affirms the presence of the corresponding AD indicator (Figure 3). We accomplished this by iterating through all sentences in each patient's full EHR and passing the sentence embedding through each of our 8 trained MLP networks to obtain 8 probabilities for each sentence corresponding to the probability that the sentence affirms each of the 8 AD indicators we previously selected. Then, for each patient and for each AD indicator, we kept the highest probability that any given sentence in the patient's record affirms the presence of the AD indicator.

**Table 4.** Training and testing data set size for each classifier.

Classifier	Number of training samples, n	Number of testing samples, n
1	2766	862
2	1302	392
3	532	168
4	9822	2454
5	1466	354
6	9114	2316
7	1596	520
8	4764	1070

**Figure 3.** Patient vector representations of atopic dermatitis indicators in experiments 1 and 2. BERT: Bidirectional Encoder Representations from Transformers; MLP: multilayer perceptron.

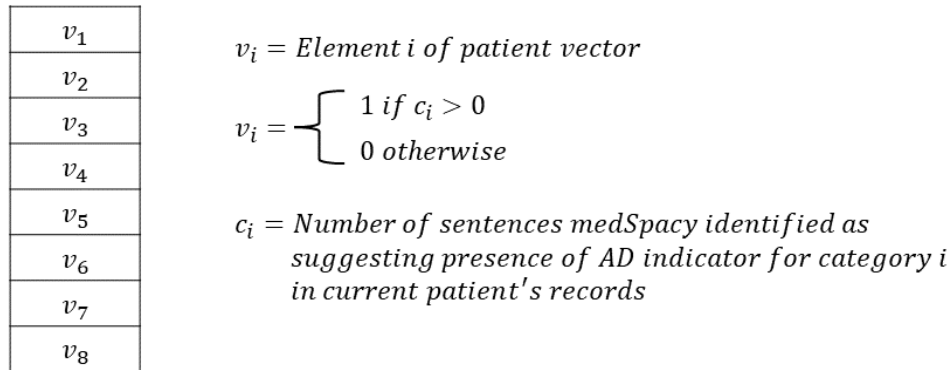
$v_1 = \max(p_{1,1}, p_{1,2}, \dots, p_{1,n})$	$v_i = \text{Element } i \text{ of patient vector}$
$v_2 = \max(p_{2,1}, p_{2,2}, \dots, p_{2,n})$	$v_i = \max(p_{i,1}, p_{i,2}, \dots, p_{i,n})$
$v_3 = \max(p_{3,1}, p_{3,2}, \dots, p_{3,n})$	$n = \text{Total number of sentences in current patient's record}$
$v_4 = \max(p_{4,1}, p_{4,2}, \dots, p_{4,n})$	$p_{i,j} = \text{Probability that sentence } j \text{ in patient record}$
$v_5 = \max(p_{5,1}, p_{5,2}, \dots, p_{5,n})$	$\text{indicates presence of category } i$
$v_6 = \max(p_{6,1}, p_{6,2}, \dots, p_{6,n})$	$p_{i,j} = \text{softmax}(MLP_i(s_j))$
$v_7 = \max(p_{7,1}, p_{7,2}, \dots, p_{7,n})$	$MLP_i = \text{Multilayer perceptron network (or alternatively, Sentence Classifier)}$
$v_8 = \max(p_{8,1}, p_{8,2}, \dots, p_{8,n})$	$\text{trained to distinguish sentences in category } i \text{ from other sentences}$
	$s_j = \text{BERT embedding of sentence } j \text{ from patient's record}$

### Experiment 3: Patient Vector Construction Without BERT Embeddings

In experiment 3, we generated each patient's vector representation by assigning a 1 to each element of the patient vector if *medSpacy* with the ConText algorithm identified at

least 1 sentence in the patient's record that affirms or suggests the presence of the AD indicator for which the vector element corresponds (Figure 4). Experiment 3 was conducted as an ablation study to quantify the performance benefit (if at all) of using contextual BERT text embeddings to generate probability scores that the patient meets various AD indicators.

**Figure 4.** Patient vector representations of atopic dermatitis (AD) indicators in experiment 3.



### AD Phenotyping With Vector Representations

In all 3 experiments, after generating a vector representation for each patient, we collated each patient's vector representation with the corresponding label our clinicians assigned the patient when validating the data set. Then, we fed the vector patient representation and corresponding patient label through a variety of classification algorithms. These include logistic regression, support vector machines (SVM), decision trees, random forests, k-nearest neighbor (KNN), Extreme Gradient Boosting (XGBoost), and Adaptive Boosting (AdaBoost). During training for each of the previously mentioned classifiers, we used 5-fold cross validation to determine the best set of hyperparameters to use (as opposed to creating a separate validation set) due to the data-scarce setting of our experiments. We then used the selected hyperparameters to train each algorithm on the entire training set and evaluated performance on the unbalanced and balanced testing sets. In addition to using the previously mentioned classifiers, we also used the stacking algorithm provided by scikit-learn to obtain an ensemble prediction from the different classifiers [23]. To quantify performance, we calculated the accuracy, precision, recall,  $F_1$ -score, negative predictive value (NPV), and specificity of each algorithm on both testing sets.

We define accuracy, precision, and recall as follows, where TP is the number of true positives, TN is the number of true negatives, FP is the number of false positives, and FN is the number of false negatives:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

Additionally, we define the  $F_1$ -score, NPV, and specificity as follows:

$$F_1 = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$NPV = \frac{TN}{TN + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

### Ethical Considerations

This research protocol was reviewed and approved by the University of Pennsylvania Institute Review Board and determined to be exempt (IRB#843922).

## Results

### Performance of MLP Networks

In this section, we compare the performance of several MLP classifiers in distinguishing sentences relevant to the diagnosis of AD. This corresponds to the "Train separate MLP network (sentence classifier) for each of 8 AD indicators" box in Figure 2.

As part of our AD phenotyping pipeline, we trained various MLP networks to classify when a given sentence embedding indicates the presence of an AD indicator, and we compared the performance of BioClinical BERT embeddings to BERT Base Uncased embeddings when training these MLP networks. In both cases, the classifier with the highest accuracy was the classifier for category 1 (sentences with direct mentions of AD). The classifiers with the 2 lowest accuracies were either the classifier for category 5 (sentences with mentions of dry or itchy skin) or the classifier for category 7 (sentences with mentions of asthma) for both the use of BioClinical BERT embeddings and the use of BERT Base Uncased embeddings. However, the accuracy in classifier 7 was lower when using BERT Base Uncased embeddings than when using BioClinical BERT embeddings.

In experiment 1, the accuracies across AD indicator classifiers ranged from 0.7373 (classifier 5) to 0.9002 (classifier 1), as shown in [Table 5](#).

In experiment 2, the accuracies across AD indicator classifiers ranged from 0.7269 (classifier 7) to 0.9153 (classifier 1), as shown in [Table 6](#).

**Table 5.** Accuracy of different multilayer perceptron networks in discerning sentences by atopic dermatitis (AD) indicator categories using “BioClinical Bidirectional Encoder Representations from Transformers” sentence embeddings.

Classifier	AD indicator	Accuracy
1	Direct mention of AD	0.9002
2	Mention of hay fever allergies	0.8954
3	Mention of atopic allergies	0.8214
4	Mention of eczema or rash	0.8284
5	Mention of dry or itchy skin	0.7373
6	Mention of nonasthma medications	0.8204
7	Mention of asthma	0.7712
8	Mention of asthma medications	0.8299

**Table 6.** Accuracy of different multilayer perceptron networks in discerning sentences by atopic dermatitis (AD) indicator categories using “Bidirectional Encoder Representations from Transformers Base Uncased” sentence embeddings.

Classifier	AD indicator	Accuracy
1	Direct mention of AD	0.9153
2	Mention of hay fever allergies	0.7730
3	Mention of atopic allergies	0.7976
4	Mention of eczema or rash	0.8439
5	Mention of dry or itchy skin	0.7288
6	Mention of nonasthma medications	0.8096
7	Mention of asthma	0.7269
8	Mention of asthma medications	0.8738

## AD Phenotyping With Patient Vector Representations

In this section, we compare performance in patient classification when using different methods for creating patient vector representations. This encompasses all 3 experiments and corresponds to the “Use vector patient representations to classify whether patient has AD” box in [Figure 2](#).

In experiment 1, we leveraged BioClinical BERT sentence embeddings to train various MLP networks to discern sentence embeddings in different AD indicator categories. Then, we applied these trained MLP networks (sentence classifiers) along with an aggregation function (max operator) to assign values to each element of each patient’s vector representation. Lastly, we used each patient’s vector representation with their validated label to train various ML algorithms. We evaluated these on both a balanced and unbalanced testing set.

As shown in [Table 7](#), the accuracy on the balanced testing set ranges from 0.5893 (decision tree) to 0.7321 (logistic regression and SVM).

As shown in [Table 8](#), the range of accuracies on the unbalanced testing set is slightly lower, ranging from 0.5824 (decision tree) to 0.7253 (stacking classifier).

In experiment 2, we followed the same process as in experiment 1; however, we used BERT Base Uncased instead of BioClinical BERT. As shown in [Table 9](#), the accuracy of our AD classifiers on the balanced testing set ranges from 0.5179 (AdaBoost) to 0.6250 (random forest).

As shown in [Table 10](#), the range of accuracies of our AD classifiers on the unbalanced testing set is slightly higher, ranging from 0.5714 (logistic regression and SVM) to 0.6703 (random forest).

In experiment 3, we performed an ablation study and assigned binary labels to the elements of each patient’s vector based on whether *medSpacy* was able to identify at least 1 sentence in each of the AD indicator categories that each vector element corresponds to. As shown in [Table 11](#), the accuracy across our AD classifiers on the balanced testing set ranges from 0.6964 (KNN) to 0.8036 (XGBoost).

As shown in [Table 12](#), the lower bound of the range of accuracies across our AD classifiers on the unbalanced testing set is higher, and the upper bound of the accuracies is lower. The accuracies on the unbalanced testing set range from 0.7143 (Stacking Classifier) to 0.7582 (Random Forest and Stacking Classifier).



**Table 7.** Atopic dermatitis phenotyping performance on balanced testing set in experiment 1 (BioClinical Bidirectional Encoder Representations from Transformers).

Model	Accuracy	Precision	Recall	$F_1$ -score	NPV <sup>a</sup>	Specificity
Logistic regression	0.7321	0.7241	0.7500	0.7368	0.7407	0.7500
SVM <sup>b</sup>	0.7321	0.7826	0.6429	0.7059	0.6970	0.7857
Decision tree	0.5893	0.6316	0.4286	0.5106	0.5676	0.7500
Random forest	0.6964	0.7037	0.6786	0.6909	0.6897	0.8214
KNN <sup>c</sup>	0.6786	0.7273	0.5714	0.6400	0.6471	0.7857
XGBoost <sup>d</sup>	0.6071	0.6154	0.5714	0.5926	0.6000	0.8571
AdaBoost <sup>e</sup>	0.6429	0.6538	0.6071	0.6296	0.6333	0.7857
Stacking classifier	0.6964	0.7391	0.6071	0.6667	0.6667	0.7500

<sup>a</sup>NPV: negative predictive value.

<sup>b</sup>SVM: support vector machines.

<sup>c</sup>KNN: k-nearest neighbor.

<sup>d</sup>XGBoost: Extreme Gradient Boosting.

<sup>e</sup>AdaBoost: Adaptive Boosting.

**Table 8.** Atopic dermatitis phenotyping performance on unbalanced testing set in experiment 1 (BioClinical Bidirectional Encoder Representations from Transformers).

Model	Accuracy	Precision	Recall	$F_1$ -score	NPV <sup>a</sup>	Specificity
Logistic regression	0.6813	0.4884	0.7500	0.5915	0.8542	0.6984
SVM <sup>b</sup>	0.6923	0.5000	0.6429	0.5625	0.8181	0.7302
Decision tree	0.5824	0.3438	0.3929	0.3667	0.7119	0.7143
Random forest	0.7143	0.5313	0.6071	0.5667	0.6845	0.7619
KNN <sup>c</sup>	0.6593	0.4571	0.5714	0.5079	0.7857	0.7937
XGBoost <sup>d</sup>	0.6264	0.4211	0.5714	0.4848	0.7736	0.7619
AdaBoost <sup>e</sup>	0.6044	0.4048	0.6071	0.4857	0.7755	0.7302
Stacking classifier	0.7253	0.5429	0.6786	0.6032	0.8393	0.6984

<sup>a</sup>NPV: negative predictive value.

<sup>b</sup>SVM: support vector machines.

<sup>c</sup>KNN: k-nearest neighbor.

<sup>d</sup>XGBoost: Extreme Gradient Boosting.

<sup>e</sup>AdaBoost: Adaptive Boosting.

**Table 9.** Atopic dermatitis phenotyping performance on balanced testing set in experiment 2 (Bidirectional Encoder Representations from Transformers Base Uncased).

Model	Accuracy	Precision	Recall	$F_1$ -score	NPV <sup>a</sup>	Specificity
Logistic regression	0.5893	0.5758	0.6786	0.6230	0.6087	0.5000
SVM <sup>b</sup>	0.6071	0.5938	0.6786	0.6333	0.6250	0.5357
Decision tree	0.6071	0.6071	0.6071	0.6071	0.6071	0.6071
Random forest	0.6250	0.6522	0.5357	0.5882	0.6061	0.7143
KNN <sup>c</sup>	0.5536	0.5714	0.4286	0.4898	0.5429	0.6786
XGBoost <sup>d</sup>	0.5536	0.5556	0.5357	0.5455	0.5517	0.5714
AdaBoost <sup>e</sup>	0.5179	0.5185	0.5000	0.5091	0.5172	0.5357
Stacking classifier	0.6071	0.6071	0.6071	0.6071	0.6071	0.6071

<sup>a</sup>NPV: negative predictive value.

<sup>b</sup>SVM: support vector machines.

<sup>c</sup>KNN: k-nearest neighbor.

<sup>d</sup>XGBoost: Extreme Gradient Boosting.

<sup>e</sup>AdaBoost: Adaptive Boosting.

**Table 10.** Atopic dermatitis phenotyping performance on unbalanced testing set in experiment 2 (Bidirectional Encoder Representations from Transformers Base Uncased).

Model	Accuracy	Precision	Recall	$F_1$ -score	NPV <sup>a</sup>	Specificity
Logistic regression	0.5714	0.3878	0.6786	0.4935	0.7857	0.5238
SVM <sup>b</sup>	0.5714	0.3878	0.6786	0.4935	0.7857	0.5238
Decision tree	0.6484	0.4474	0.6071	0.5152	0.7925	0.6667
Random forest	0.6703	0.4737	0.6429	0.5455	0.8113	0.6825
KNN <sup>c</sup>	0.6264	0.4000	0.4286	0.4138	0.7377	0.7143
XGBoost <sup>d</sup>	0.6374	0.4286	0.5357	0.4762	0.7679	0.6825
AdaBoost <sup>e</sup>	0.5934	0.3784	0.5000	0.4308	0.7407	0.6349
Stacking classifier	0.6484	0.4474	0.6071	0.5152	0.7925	0.6667

<sup>a</sup>NPV: negative predictive value.

<sup>b</sup>SVM: support vector machines.

<sup>c</sup>KNN: k-nearest neighbor.

<sup>d</sup>XGBoost: Extreme Gradient Boosting.

<sup>e</sup>AdaBoost: Adaptive Boosting.

**Table 11.** Atopic dermatitis phenotyping performance on balanced testing set in experiment 3 (binary vector encoding).

Model	Accuracy	Precision	Recall	$F_1$ -score	NPV <sup>a</sup>	Specificity
Logistic regression	0.7679	0.7586	0.7857	0.7719	0.7778	0.7500
SVM <sup>b</sup>	0.7857	0.7857	0.7857	0.7857	0.7857	0.7857
Decision tree	0.7857	0.7667	0.8214	0.7931	0.8077	0.7500
Random forest	0.7857	0.8077	0.7500	0.7778	0.7667	0.8214
KNN <sup>c</sup>	0.6964	0.7391	0.6071	0.6667	0.6667	0.7857
XGBoost <sup>d</sup>	0.8036	0.8400	0.7500	0.7925	0.7742	0.8571
AdaBoost <sup>e</sup>	0.7857	0.7857	0.7857	0.7857	0.7857	0.7857
Stacking classifier	0.7500	0.7500	0.7500	0.7500	0.7500	0.7500

<sup>a</sup>NPV: negative predictive value.

<sup>b</sup>SVM: support vector machines.

<sup>c</sup>KNN: k-nearest neighbor.

<sup>d</sup>XGBoost: Extreme Gradient Boosting.

<sup>e</sup>AdaBoost: Adaptive Boosting.

**Table 12.** Atopic dermatitis phenotyping performance on unbalanced testing set in experiment 3 (binary vector encoding).

Model	Accuracy	Precision	Recall	$F_1$ -score	NPV <sup>a</sup>	Specificity
Logistic regression	0.7253	0.5366	0.7857	0.6377	0.8800	0.6984
SVM <sup>b</sup>	0.7473	0.5641	0.7857	0.6567	0.8846	0.7302
Decision tree	0.7473	0.5610	0.8214	0.6667	0.9000	0.7143
Random forest	0.7582	0.5833	0.7500	0.6563	0.8727	0.7619
KNN <sup>c</sup>	0.7363	0.5667	0.6071	0.5862	0.8197	0.7937
XGBoost <sup>d</sup>	0.7582	0.5833	0.7500	0.6563	0.8727	0.7619
AdaBoost <sup>e</sup>	0.7473	0.5641	0.7857	0.6567	0.8846	0.7302
Stacking classifier	0.7143	0.5250	0.7500	0.6176	0.8627	0.6984

<sup>a</sup>NPV: negative predictive value.

<sup>b</sup>SVM: support vector machines.

<sup>c</sup>KNN: k-nearest neighbor.

<sup>d</sup>XGBoost: Extreme Gradient Boosting.

<sup>e</sup>AdaBoost: Adaptive Boosting.

## Discussion

### Sentence Classification Results

We hypothesized that using BioClinical BERT sentence embeddings to train sentence classifiers would provide better performance than using BERT Base Uncased sentence embeddings due to the clinical setting of our data. Given the results in Tables 5 and 6, we observed that this was most often true in the context of sentence classification because we were able to achieve better performance in the majority (5 out of 8) of the sentence classification tasks when using BioClinical BERT embeddings as opposed to BERT Base Uncased embeddings.

Using BioClinical BERT sentence embeddings yielded stronger performance when distinguishing sentences in 5 of the 8

sentence categories: category 2 (mentions of hay fever allergies), category 3 (mentions of atopic allergies), category 5 (mentions of dry or itchy skin), category 6 (mentions of nonasthma medications), and category 7 (mentions of asthma). More specifically, we observed higher accuracies when using BioClinical BERT sentence embeddings for classifiers 2 (0.8954), 3 (0.8214), 5 (0.7373), 6 (0.8204), and 7 (0.7712) than their corresponding counterparts when using BERT Base Uncased embeddings for classifiers 2 (0.7730), 3 (0.7976), 5 (0.7288), 6 (0.8096), and 7 (0.7269). We observed that the differences in performance between using BioClinical BERT embeddings and BERT Base Uncased embeddings are most pronounced for classifiers 2 and 7, which correspond to mentions of hay fever allergies and asthma mentions, respectively. We hypothesize this is because hay fever allergies and asthma (and their synonyms) may be very common terms

in clinical notes; therefore, models trained on clinical data (BioClinical BERT) may be able to provide stronger performance than models trained on nonclinical text (BERT Base Uncased), which may not have as many mentions of hay fever allergies or asthma.

Conversely, using BERT Base Uncased embeddings yielded stronger performance when distinguishing sentences in the other 3 of 8 sentence categories: category 1 (direct mentions of AD), category 4 (mentions of eczema or rashes), and category 8 (mentions of asthma medications). More specifically, we observed higher accuracies when using BERT Base Uncased sentence embeddings for classifiers 1 (0.9153), 4 (0.8439), and 8 (0.8738) than their corresponding counterparts when using BioClinical BERT embeddings for classifiers 1 (0.9002), 4 (0.8284), and 8 (0.8299). We observed differences in performance between using BERT Base Uncased embeddings and BioClinical BERT embeddings, which are most evident for classifier 8, which corresponds to mentions of asthma medications. Although this is counterintuitive at first (we would expect a classifier using embeddings generated from BioClinical BERT to be able to better recognize allergy medicines), we believe that the performance benefit from using BERT Base Uncased can be attributed to the list of terms we gave to *medSpacy* when asking it to identify sentences in category 8. Many of the asthma medications in category 8 sentences are either monoclonal antibody medications ending in -mab (benralizumab, mepolizumab, omalizumab, etc) or hydrofluoroalkanes (hfa; atrovent hfa, flovent hfa, xopenex hfa, etc). Because monoclonal antibodies are very specialized types of medication, they may not occur as frequently as other terms in the corpus used to train BioClinical BERT, so a more general model such as BERT Base Uncased may provide more robust performance. Additionally, because the hydrofluoroalkane allergy medications in category 8 sentences are often abbreviated with “hfa,” which can have alternate medical meanings such as high-functioning autism or health facility administrator, the BioClinical BERT embeddings might not be representative of the presence of allergy medications in the sentence, so a more general model such as BERT Base Uncased may be able to provide better performance.

More broadly, looking at the results in [Tables 5 and 6](#), we can see that the least accurate classifier has an accuracy of 0.7288, while the most accurate classifier is able to achieve an accuracy of 0.9153. Furthermore, when aggregating the most accurate classifiers from both tables we can see that we are able to achieve accuracies of 0.9153 (classifier 1) for identifying sentences that directly suggest the patient has AD, 0.8954 (classifier 2) for identifying sentences that mention hay fever allergies, 0.8214 (classifier 3) for identifying sentences that mention atopic allergies, 0.8439 (classifier 4) for identifying sentences that mention eczema or skin rashes, 0.7373 (classifier 5) for identifying sentences that mention dry or itchy skin, 0.8204 (classifier 6) for identifying sentences that mention nonasthma medications related to diagnosis of AD, 0.7712 (classifier 7) for identifying sentences that mention asthma, and 0.8738 (classifier 8) for identifying sentences that mention asthma medications. Because our training and testing sets were both class-balanced and the majority (6 of the 8) of the most

accurate classifiers previously mentioned achieved accuracies between 0.8204 and 0.9153, we believe these results are promising and indicate that our sentence classifiers could potentially be used to save time in a clinical setting during chart review by identifying (and highlighting for review) sentences relevant to the diagnosis of AD when recruiting for clinical trials.

## AD Phenotyping Results

As per [Tables 7-10](#), our earlier hypothesis holds: using clinical embeddings (BioClinical BERT) to generate the patient vector representation does provide better performance in patient phenotyping than using nonclinical embeddings (BERT Base Uncased). Comparing evaluations on the balanced testing set in [Tables 7 and 9](#), we observe that using BioClinical BERT embeddings provides higher accuracy in almost all models, with the exception of Decision Trees where BERT Base Uncased provides better performance (accuracy of 0.6071) as compared with BioClinical BERT (accuracy of 0.5893). Comparing evaluations on the unbalanced testing set in [Tables 8 and 10](#), we observed that the same trend follows: using BioClinical BERT embeddings provides higher accuracy in almost all models, with the exception of Decision Trees and XGBoost, where using BERT Base Uncased embeddings provides better performance (accuracy of 0.6484 for Decision Trees and 0.6374 for XGBoost) as compared with their counterparts with BioClinical BERT embeddings (accuracy of 0.5824 for Decision Trees and 0.6264 for XGBoost).

As part of our experimental design, we included an ablation study in experiment 3 so we could compare the difference in performance during patient phenotyping when removing the use of BERT models to create each patient’s vector representations. On the class-balanced testing set, we observed that accuracies range from 0.6071 to 0.7321 when using BioClinical BERT embeddings in [Table 7](#), accuracies range from 0.5179 to 0.6250 when using BERT Base Uncased embeddings in [Table 9](#), and accuracies range from 0.6964 to 0.8036 when removing the use of BERT models in [Table 11](#) (experiment 3). On the unbalanced testing set, we observed that accuracies range from 0.5824 to 0.7253 when using BioClinical BERT embeddings in [Table 8](#), accuracies range from 0.5714 to 0.6703 when using BERT Base Uncased embeddings in [Table 10](#), and accuracies range from 0.7143 to 0.7582 when removing the use of BERT models in [Table 12](#) (experiment 3).

In both cases (evaluation on the balanced testing set and evaluation on the unbalanced testing set), we found that models in experiment 3 (ablation study) generally outperform (or are as good as) their corresponding counterparts in experiments 1 and 2 (BERT experiments) across all metrics (accuracy, precision, recall,  $F_1$ -score, NPV, and specificity), with the exception that the stacking classifier in experiment 1 (BioClinical BERT) has marginally stronger accuracy and precision than the stacking classifier in experiment 3. This shows that traditional rules-based approaches (experiment 3) can outperform BERT-based approaches for generating a patient vector representation for downstream patient phenotyping.

We hypothesize that models in experiments 1 and 2 showed lower performance because errors from our sentence classifiers

in earlier stages of the pipeline could have propagated to later stages of the pipeline during patient phenotyping. Because we leveraged the max operator to aggregate probabilities that any given sentence in the patient record applies to each category, more sentences in each patient record would lead to a greater chance that an erroneous prediction with a high probability would lead to a false positive error in the creation of each patient's vector representation in experiments 1 and 2.

Although there is a wide range in performance for our patients with AD phenotyping algorithms, we believe that we have reached our goal of developing a system capable of patient with AD phenotyping for clinical trial recruitment because [Tables 11 and 12](#) show promising results. Furthermore, our system can be used as a first step during AD clinical trial recruitment to filter out most patients who may not qualify for AD trials and therefore save valuable clinician time. We believe our pipeline is important and valuable because, unlike other diseases, such as influenza, COVID-19, and cancer, there is no gold-standard test result that can be used to determine when a patient has AD. Instead, clinicians must spend large amounts of time undergoing chart reviews to individually determine whether each patient has AD.

### Limitations

One limitation of this study was the small size of our data set. Although we had a total of 1926 patients in our data set, only 137 of them were validated as having AD. During training, we leveraged 109 of the 137 patients with AD and sampled another 109 patients without AD to create a class-balanced training set. The small size of the training set could lead to overfitting and therefore result in reduced performance on the testing set. Future work could involve obtaining more data from patients with AD as well as exploring the use of an imbalanced data set but using a class-weighted loss function to counteract the class imbalance.

A second limitation of this study was the input-limit size of the large language models that were used. Both BERT Base Uncased and BioClinical BERT had an input limit of 512 tokens. This meant that any input text that was longer than 512 tokens would be ignored when training BERT. Consequently, we could not simply directly concatenate all documents from each patient's EHR and feed the tokenized documents of each patient into BERT with an added classification head for training as well as direct prediction of whether the patient has AD. Instead, we designed a pipeline around distilling information from all documents in each patient's EHR into a patient vector representation and then using this patient vector representation to train various classical ML algorithms for phenotyping the patient. Future work could involve exploring the use of other large language models that are suited for long inputs, such as Longformer or Doc2Vec, for predicting when a patient should be labeled as having AD.

### Acknowledgments

This study was partially funded by the National Institutes of Health (NIH) and the National Institute of Arthritis and Musculoskeletal and Skin Diseases (NIAMS) P30-AR069589 as part of the Penn Skin Biology and Diseases Resource-Based Center (core: DJM and DM).

A third limitation of this study was the list of AD indicators we selected. We did not consider additional AD indicators, and we also did not consider the use of different combinations (or subsets) of the AD indicators selected. This is particularly relevant in considering that (1) our pipeline is intended to be used for identifying patients with AD, and (2) one of our AD indicators (category 1) directly targets whether there is any given sentence in the patient's record that mentions AD, which could be in the context of a family history of AD, a potential (but not confirmed) diagnosis of AD, as well as a confirmed diagnosis of AD, among other possibilities. If this AD indicator is removed, then 1 interesting research question could be whether our pipeline is still able to maintain performance similarly to what it is currently able to achieve. Future work could involve assessing the performance impact of removing or adding the use of various AD indicators. We could then determine if our pipeline is relying too much on or overfitting 1 or more indicators. Furthermore, we could also redesign our patient vector and separate the feature for category 1 (any sentence that mentions AD) into 3 separate indicators: whether there is (1) a family history of AD, (2) an affirmed diagnosis that the patient has AD, and (3) uncertainty of whether the patient has AD. Doing so could potentially improve precision.

### Potential Applications

Given the aforementioned results, we believe our AD classifier could be operationalized to facilitate reliable and efficient EHR chart review. For example, sentence classifiers could visually indicate AD indicators inline text, therefore reducing information foraging efforts by clinicians. Additionally, AD phenotyping classifiers could indicate the strength of a patient match to UKWP criteria, exact or partial, based on AD indicator sentence classifications. Furthermore, ranking patient cases by match strength could reduce the number of cases reviewed to generate both case and matched controls.

### Conclusions

In conclusion, we present and validate a promising pipeline for phenotyping patients with AD during clinical trial recruitment. To do so, we compare a rules-based and transformer-based approach for creating a vector representation of each patient and compare downstream performance in patient phenotyping with various standard ML algorithms. We find that a traditional rules-based approach outperforms using a transformer-based approach (experiment 3). We hope that our pipeline can be deployed in hospital settings during clinical trial recruitment as an initial step to automatically filter candidates before manual review. Additionally, we show that MLP networks can identify whether sentences are relevant to AD diagnosis. These MLP networks can later be deployed in clinical settings to highlight which sentences are relevant for physicians during manual chart review, therefore reducing physician burden. Future work can involve extending our patient phenotyping pipeline to other data sets and other diseases.



## Authors' Contributions

AW designed the experiments, wrote the code, performed the experiments, wrote the first draft of the manuscript, and revised the manuscript. DJM conceptualized and implemented the chart abstraction study, annotated the data set, interpreted the results, and revised the manuscript. RF annotated the data set and revised the manuscript. SH queried and deidentified the data set as well as revised the manuscript. DM conceptualized the study and experiment design, interpreted results, wrote and revised the manuscript, and provided secure storage and computer resources.

## Conflicts of Interest

DJM is or recently has been a consultant for Pfizer, Leo, and Sanofi with respect to studies of atopic dermatitis and served on an advisory board for the National Eczema Association.

## References

1. Collins-Williams C. Eczema (atopic dermatitis). In: Paediatric Allergy and Clinical Immunology (as Applied to Atopic Disease): A Manual for Students and Practitioners of Medicine. Toronto. University of Toronto Press; 1973;32-37.
2. Lyons JJ, Milner JD, Stone KD. Atopic dermatitis in children: clinical features, pathophysiology, and treatment. *Immunol Allergy Clin North Am*. 2015;35(1):161-183. [FREE Full text] [doi: [10.1016/j.jac.2014.09.008](https://doi.org/10.1016/j.jac.2014.09.008)] [Medline: [25459583](https://pubmed.ncbi.nlm.nih.gov/25459583/)]
3. Eichenfield LF, Tom WL, Chamlin SL, Feldman SR, Hanifin JM, Simpson EL, et al. Guidelines of care for the management of atopic dermatitis: section 1. Diagnosis and assessment of atopic dermatitis. *J Am Acad Dermatol*. 2014;70(2):338-351. [FREE Full text] [doi: [10.1016/j.jaad.2013.10.010](https://doi.org/10.1016/j.jaad.2013.10.010)] [Medline: [24290431](https://pubmed.ncbi.nlm.nih.gov/24290431/)]
4. Abramovits W. Atopic dermatitis. *J Am Acad Dermatol*. 2005;53(1 Suppl 1):S86-S93. [doi: [10.1016/j.jaad.2005.04.034](https://doi.org/10.1016/j.jaad.2005.04.034)] [Medline: [15968268](https://pubmed.ncbi.nlm.nih.gov/15968268/)]
5. Weidinger S, Beck LA, Bieber T, Kabashima K, Irvine AD. Atopic dermatitis. *Nat Rev Dis Primers*. 2018;4(1):1. [doi: [10.1038/s41572-018-0001-z](https://doi.org/10.1038/s41572-018-0001-z)] [Medline: [29930242](https://pubmed.ncbi.nlm.nih.gov/29930242/)]
6. Schneider L, Hanifin J, Boguniewicz M, Eichenfield LF, Spergel JM, Dakovic R, et al. Study of the atopic march: development of atopic comorbidities. *Pediatr Dermatol*. 2016;33(4):388-398. [FREE Full text] [doi: [10.1111/pde.12867](https://doi.org/10.1111/pde.12867)] [Medline: [27273433](https://pubmed.ncbi.nlm.nih.gov/27273433/)]
7. Del Pozo DV, Zhu Y, Mitra N, Hoffstad OJ, Margolis DJ. The risk of atopic comorbidities and atopic march progression among Black and White children with mild-to-moderate atopic dermatitis: a cross-sectional study. *J Am Acad Dermatol*. 2022;87(5):1145-1147. [doi: [10.1016/j.jaad.2022.02.023](https://doi.org/10.1016/j.jaad.2022.02.023)] [Medline: [35192898](https://pubmed.ncbi.nlm.nih.gov/35192898/)]
8. Eichenfield LF, Tom WL, Berger TG, Krol A, Paller AS, Schwarzenberger K, et al. Guidelines of care for the management of atopic dermatitis: section 2. Management and treatment of atopic dermatitis with topical therapies. *J Am Acad Dermatol*. 2014;71(1):116-132. [FREE Full text] [doi: [10.1016/j.jaad.2014.03.023](https://doi.org/10.1016/j.jaad.2014.03.023)] [Medline: [24813302](https://pubmed.ncbi.nlm.nih.gov/24813302/)]
9. Fulton RL, Mitra N, Chiesa-Fuxench Z, Sockler PG, Margolis DJ. Untapping the potential of utilizing electronic medical records to identify patients with atopic dermatitis: an algorithm using ICD-10 codes. *Arch Dermatol Res*. 2022;314(5):439-444. [doi: [10.1007/s00403-021-02251-w](https://doi.org/10.1007/s00403-021-02251-w)] [Medline: [34081192](https://pubmed.ncbi.nlm.nih.gov/34081192/)]
10. Gustafson E, Pacheco J, Wehbe F, Silverberg J, Thompson W. A machine learning algorithm for identifying atopic dermatitis in adults from electronic health records. *IEEE Int Conf Healthc Inform*. 2017;2017:83-90. [FREE Full text] [doi: [10.1109/ICHI.2017.31](https://doi.org/10.1109/ICHI.2017.31)] [Medline: [29104964](https://pubmed.ncbi.nlm.nih.gov/29104964/)]
11. Hanifin JM, Rajka G. Diagnostic features of atopic dermatitis. *Acta Derm Venereol*. 1980;60:44-47. [FREE Full text] [doi: [10.2340/00015555924447](https://doi.org/10.2340/00015555924447)]
12. Williams HC, Burney PG, Hay RJ, Archer CB, Shipley MJ, Hunter JJ, et al. The U.K. Working Party's diagnostic criteria for atopic dermatitis. I. Derivation of a minimum set of discriminators for atopic dermatitis. *Br J Dermatol*. 1994;131(3):383-396. [doi: [10.1111/j.1365-2133.1994.tb08530.x](https://doi.org/10.1111/j.1365-2133.1994.tb08530.x)] [Medline: [7918015](https://pubmed.ncbi.nlm.nih.gov/7918015/)]
13. Norgeot B, Muenzen K, Peterson TA, Fan X, Glicksberg BS, Schenk G, et al. Protected Health Information filter (Philter): accurately and securely de-identifying free-text clinical notes. *NPJ Digit Med*. 2020;3:57. [FREE Full text] [doi: [10.1038/s41746-020-0258-y](https://doi.org/10.1038/s41746-020-0258-y)] [Medline: [32337372](https://pubmed.ncbi.nlm.nih.gov/32337372/)]
14. Bert-base-uncased. Hugging Face. URL: <https://huggingface.co/bert-base-uncased> [accessed 2023-11-29]
15. Devlin J, Chang MW, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. ArXiv. Preprint posted online on May 24 2019. 2018 [FREE Full text]
16. Alsentzer E, Murphy JR, Boag W, Weng WH, Jin D, Naumann T, et al. Publicly available clinical BERT embeddings. ArXiv. Preprint posted online on June 20 2019. 2019 [FREE Full text] [doi: [10.18653/v1/w19-1909](https://doi.org/10.18653/v1/w19-1909)]
17. emilyalsentzer/Bio\_ClinicalBERT. Hugging Face. URL: [https://huggingface.co/emilyalsentzer/Bio\\_ClinicalBERT](https://huggingface.co/emilyalsentzer/Bio_ClinicalBERT) [accessed 2023-11-29]
18. Eyre H, Chapman AB, Peterson KS, Shi J, Alba PR, Jones MM, et al. Launching into clinical space with medspaCy: a new clinical text processing toolkit in Python. *AMIA Annu Symp Proc*. 2021;2021:438-447. [FREE Full text] [Medline: [35308962](https://pubmed.ncbi.nlm.nih.gov/35308962/)]

19. Chapman BE, Lee S, Kang HP, Chapman WW. Document-level classification of CT pulmonary angiography reports based on an extension of the ConText algorithm. *J Biomed Inform.* 2011;44(5):728-737. [FREE Full text] [doi: [10.1016/j.jbi.2011.03.011](https://doi.org/10.1016/j.jbi.2011.03.011)] [Medline: [21459155](https://pubmed.ncbi.nlm.nih.gov/21459155/)]
20. Harkema H, Dowling JN, Thornblade T, Chapman WW. ConText: an algorithm for determining negation, experienter, and temporal status from clinical reports. *J Biomed Inform.* 2009;42(5):839-851. [FREE Full text] [doi: [10.1016/j.jbi.2009.05.002](https://doi.org/10.1016/j.jbi.2009.05.002)] [Medline: [19435614](https://pubmed.ncbi.nlm.nih.gov/19435614/)]
21. Mowery DL, Kawamoto K, Bradshaw R, Kohlmann W, Schiffman JD, Weir C, et al. Determining onset for familial breast and colorectal cancer from family history comments in the electronic health record. *AMIA Jt Summits Transl Sci Proc.* 2019;2019:173-181. [FREE Full text] [Medline: [31258969](https://pubmed.ncbi.nlm.nih.gov/31258969/)]
22. Mowery DL, Velupillai S, Chapman WW. Medical diagnosis lost in translation-analysis of uncertainty and negation expressions in English and Swedish clinical texts. Presented at: Proceedings of the 2012 Workshop on Biomedical Natural Language Processing (BioNLP 2012); Montreal, Canada, 2012; June 8, 2012. URL: <https://www.aclweb.org/anthology/W12-2407.pdf>
23. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: machine learning in python. *J Mach Learn Res.* 2011;12(85):2825-2830. [FREE Full text]

## Abbreviations

**AD:** atopic dermatitis  
**AdaBoost:** Adaptive Boosting  
**BERT:** Bidirectional Encoder Representations from Transformers  
**EHR:** electronic health record  
**FN:** false negatives  
**FP:** false positives  
**Hfa:** hydrofluoroalkanes  
**HR:** Hanifin and Rajka  
**ICD:** International Classification of Disease  
**KNN:** k-nearest neighbor  
**ML:** machine learning  
**MLP:** multilayer perceptron  
**NLP:** natural language processing  
**NPV:** negative predictive value  
**Philter:** Protected Health Information filter  
**ReLU:** Rectified Linear Unit  
**SGD:** stochastic gradient descent  
**SVM:** support vector machines  
**TN:** true negatives  
**TP:** true positives  
**UKWP:** UK Working Party  
**XGBoost:** Extreme Gradient Boosting

*Edited by G Eysenbach, A Mavragani; submitted 31.08.23; peer-reviewed by J Zaghir, A Chapman; comments to author 06.10.23; revised version received 30.11.23; accepted 04.12.23; published 26.01.24*

*Please cite as:*

Wang A, Fulton R, Hwang S, Margolis DJ, Mowery D

*Patient Phenotyping for Atopic Dermatitis With Transformers and Machine Learning: Algorithm Development and Validation Study*  
*JMIR Form Res* 2024;8:e52200

URL: <https://formative.jmir.org/2024/1/e52200>

doi: [10.2196/52200](https://doi.org/10.2196/52200)

PMID: [38277207](https://pubmed.ncbi.nlm.nih.gov/38277207/)

©Andrew Wang, Rachel Fulton, Sy Hwang, David J Margolis, Danielle Mowery. Originally published in JMIR Formative Research (<https://formative.jmir.org>), 26.01.2024. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Formative Research, is properly cited. The

complete bibliographic information, a link to the original publication on <https://formative.jmir.org>, as well as this copyright and license information must be included.