

Towards Reducing Diagnostic Errors with Interpretable Risk Prediction

Anonymous ACL submission

Abstract

Many diagnostic errors occur because clinicians cannot easily access relevant information in patient Electronic Health Records (EHRs). In this work we propose a method to use LLMs to identify pieces of evidence in patient EHR data that indicate increased or decreased risk of specific diagnoses; our ultimate aim is to increase access to evidence and reduce diagnostic errors. In particular, we propose a Neural Additive Model to make predictions backed by evidence with individualized risk estimates at time-points where clinicians are still uncertain, aiming to specifically mitigate delays in diagnosis and errors stemming from an incomplete differential. To train such a model, it is necessary to infer temporally fine-grained retrospective labels of eventual “true” diagnoses. We do so with LLMs, to ensure that the input text is from *before* a confident diagnosis can be made. We use an LLM to retrieve an initial pool of evidence, but then refine this set of evidence according to correlations learned by the model. We conduct an in-depth evaluation of the usefulness of our approach by simulating how it might be used by a clinician to decide between a pre-defined list of differential diagnoses.

1 Introduction

A major source of poor patient outcomes and unnecessary costs in healthcare are missed or delayed diagnoses. A recent report estimated that diagnostic errors result in around 795,000 serious harms annually (Newman-Toker et al., 2023). Furthermore, many diagnostic errors result from information transfer problems (Zwaan et al., 2010). This is unsurprising given “note bloat”, i.e., the widespread problem of information overload in EHR notes, often due to copied or irrelevant information which obfuscates relevant information. All of this motivates the potential of providing more efficient mechanisms to access relevant information in EHRs as a means to reduce these errors.

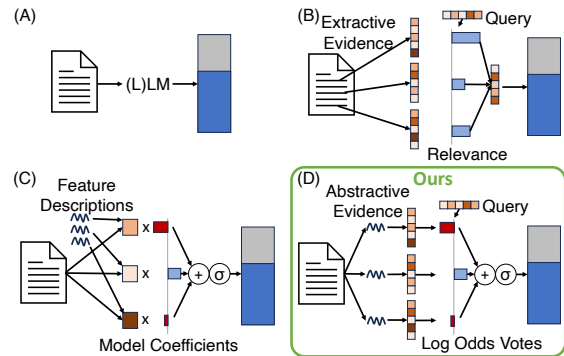


Figure 1: **Inherently “interpretable” approaches to prediction.** Typically, ‘interpretable’ models trade off between the expressiveness of intermediate representations and the faithfulness of the resulting interpretability to the models’ true mechanisms. Our approach (D) manages to use very expressive intermediate representations in the form of abstractive natural language evidence while still maintaining true transparency during aggregation of this evidence. See Table 1 for more details.

One approach to helping practitioners make use of EHR is to train NLP models on free-text notes to provide predictions about patient risk for various illnesses (Rasmy et al., 2021; Li et al., 2021; Yang et al., 2023), but these systems are often lack transparency. Even when systems have high accuracy, clinicians may still prefer simple linear models as clinical decision support tools (Goldstein et al., 2016). Prior work has focused on developing inherently interpretable¹ models with minimal tradeoff in predictive performance, e.g., in the general domain with Neural Additive Models (Agarwal et al., 2020) and in healthcare with GA²Ms (Caruana et al., 2015). Recently, zero-shot instruction-tuned LLMs have been shown capable of extracting information from clinical text (Agrawal et al., 2022), which in turn facilitates interpretable predictions (McInerney et al., 2023; Alsentzer et al., 2023).

¹Interpretability is a famously ambiguous term; we are focused on having explicit measure of the contribution of individual pieces of evidence to an output.

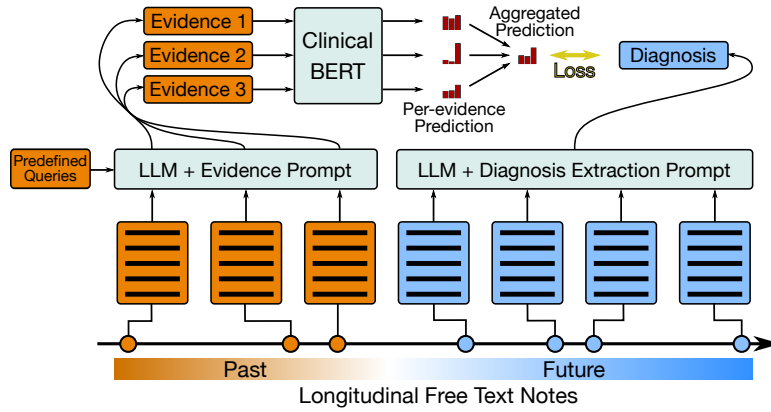


Figure 2: **Explainable Risk Prediction and Training.** An overview of our approach. Left: We retrieve evidence snippets from past notes with an LLM for predefined queries posed by a clinician. Then we use our risk prediction model to estimate risk of various diagnoses given each piece of evidence individually, and aggregate these scores. Right: We automatically extract diagnosis ‘labels’ from future reports with an LLM to use to train the risk predictor.

In this work, we combine the power and flexibility of zero-shot instruction-tuned LLMs with the transparency and modeling ability of Neural Additive Models (NAMs) to train a risk-prediction model that can also surface evidence to support predictions. We use an LLM (FLAN-T5-XXL; Chung et al. 2022) to generate abstractive “evidence” from EHR, which is then processed by a simpler model (Clinical BERT; Alsentzer et al. 2019) to produce features for a Neural Additive Model (Figure 2). This provides flexibility—the model can make inferences and condense information into fluent text snippets—but brings risk of “hallucinations”.

We view this approach as “interpretable” in that it produces “evidence” in the form of human-understandable intermediate variables: abstractive text with associated risks, providing insights into factors important to a given prediction. Relative to other approaches to inherent interpretability (Figure 1), like those that use a relevance weights to weight and combine information from different sentences (B) and those that use large language model prompts to infer feature values (C), our approach permits greater flexibility in comparison to (C), while maintaining a more faithful interpretability in comparison to (B); see Table 1 for contrasts.

One complication is that we would like fine-grained, accurate labels to train our predictor (see section 4.1); ICD codes do not meet these criteria (Searle et al., 2020). Instead of ICD codes, which are noisy and temporally coarse (observed at the end of an encounter with discharge summaries), we propose to synthetically extract diagnosis labels from each report using an LLM. In some cases, this has been shown to be more aligned with true

diagnoses (Alsentzer et al., 2023).

We focus our evaluation on how this system impacts clinical decision-making. Specifically, we examine settings where risk of misdiagnosis is high and the consequences severe. Our methods work within the confines of data present in electronic health record, which allows the model to be trained on any EHR. LLMs can be run locally and are only used for inference, so privacy and compute resources are not an issue.

Our contributions are summarized as follows:

Interpretable Risk Prediction with LLMs. We propose an approach to risk prediction that offers a particular form of interpretability in that it can expose faithful relationships between specific pieces of retrieved evidence and an output prediction.

Extracting Future Targets with LLMs. We present a method to extract target diagnoses for use in training from the unstructured text in the future of a patient’s medical record that are more granular than ICD codes in the time dimension, and we validate with clinician annotations that the extracted labels are accurate.

In-depth Annotation of Usefulness. We validate how much evidence-wise interpretability can positively impact a clinician’s expert judgement in high-impact settings which feature the greatest risk of misdiagnosis.

2 Dataset

We use MIMIC-III (Johnson et al., 2016a,b), an open-source dataset of EHRs from ICU patients. The ICU is one of the hospital settings (along with, e.g., the ER and Radiology) where misdiagnosis or

Modeling Approach	Intermediate representation(s)	Aggregation	Interpretability
(A) Direct Black-box Prediction (e.g., zero/few-shot, fine-tuned LLM)	None	CLS or last token embedding + classification or LM head	No inherent interpretability
(B) Aggregating chunked input with relevance weights	<i>Extractive</i> text snippets	Weighted avg. of CLS embeddings + class. head	Positive, real-valued relevance scores per query
(C) Logistic regression with LLM-inferred features	<i>Inferred</i> , real-valued numbers relating to predefined natural language queries	Logistic regression	Negative and positive real-valued static model coefficients
(D) Log odds voting with LLM-inferred text snippets (ours)	<i>Inferred/abstractive</i> text snippets relating to predefined natural language queries	Neural Additive Model (conditioned on the query/condition vector)	Negative and positive real-valued dynamic impact scores

Table 1: Types of interpretability afforded by the different modeling approaches for EHR data visualized in Figure 1. Red and green denote negative and positive aspects of each model.

128 delayed diagnosis are often caused by incomplete
129 information, since clinicians typically do not have
130 enough time to fully examine a patient’s EHR.

131 In healthcare, cancer, inflection, and vascular
132 dysfunction (termed the “big three”) account
133 for about 75% of all mis-diagnosis-related harms
134 (Newman-Toker et al., 2023). Within the ICU, the
135 latter two categories mostly manifest as pneumonia,
136 and pulmonary edema (which in this paper we treat
137 as interchangeable with congestive heart failure).
138 For this reason, we will focus on predicting the
139 risk of ICU patients for cancer, pneumonia, and
140 pulmonary edema. These are also conditions for
141 which clinical correlation with notes from the past
142 EHR is important for diagnosis. We use all patients
143 in the MIMIC dataset so that we have both negative
144 and positive examples of the conditions.

145 3 An Interpretable Risk Prediction Model

146 We propose a multi-stage approach to risk predic-
147 tion, capitalizing on a modern LLM, FLAN-T5-
148 XXL (Chung et al., 2022; Wei et al., 2022) in this
149 case, to implement each of the following steps.

150 **Retrieval (Section 3.1).** We generate abstractive
151 evidence from free text notes by prompting an LLM
152 with appropriate queries. The evidence snippets
153 provide a form of interpretability, in that they can
154 be inspected directly to verify predictions.

155 **Risk Prediction (Section 3.2).** We input the ev-
156 idence into the risk predictor, which models rela-
157 tionships between the evidence and each of the
158 potential diagnoses and outputs multi-label classi-
159 fication probabilities, i.e. the predicted risk that the
160 patient will be diagnosed with each condition.

161 **Evidence Re-ranking (Section 3.3).** The retrieved
162 evidence may still be too large a pool to review
163 given the time constraints of the clinician. There-

164 fore, we re-rank the evidence so as to only show
165 that which promotes risk predictions that most de-
166 viate from the baseline risks of each condition.

167 To train risk prediction models we use use syn-
168 thetic labels extracted from *future* notes in a pa-
169 tient’s record (Section 4). Figure 2 provides an
170 overview of our model and training approach.

171 3.1 Evidence Retrieval

172 Following prior work (Ahsan et al., 2023), we use a
173 sequential prompting strategy to retrieve evidence
174 that is relevant to a queried diagnosis or a risk fac-
175 tor. Specifically, we first ask the LLM for a binary
176 response as to whether evidence for a condition
177 exists; if the answer is affirmative, we then issue a
178 second prompt tasking the LLM to generate sup-
179 porting evidence. Formally, we define the evidence
180 retrieved for report n and query q_i as follows:

$$181 e_{n,q_i} = \begin{cases} \text{GetEvidence}(r_n, q_i) & \\ \text{if EvidenceExists}(r_n, q_i) = \text{“yes”} & \\ \text{null} & \text{otherwise} \end{cases} \quad (1)$$

182 where “GetEvidence” and “EvidenceExists” rep-
183 resent the corresponding prompt functions.

184 This approach does have limitations. For ex-
185 ample, it cannot produce more than one snippet
186 of evidence per report/query pair. Retrieved evi-
187 dence may also be abstractive rather than extrac-
188 tive, which introduces the risk of model “halluci-
189 nations”, but permits flexibility and interpretability
190 (Ahsan et al., 2023). It also significantly reduces
191 the amount of text (therefore requiring a relatively
192 small context window) by going from all reports
193 to sentence-length snippets for some reports. The
194 resulting “summarization” in the form of evidence
195 snippets is also controllable through the querying
196 process and works zero-shot, i.e., it requires no

specialized or in-domain training. Queries, written by a clinician co-author, are in appendix Figure 9.

3.2 Risk Prediction

Because a patient can have more than one diagnosis, we treat risk prediction as a multi-label classification problem where each label corresponds to a diagnosis. To realize interpretability, we use a Neural Additive Model (Agarwal et al., 2020). Specifically, we do not model *interactions* between evidence snippets. Instead, we predict scores individually for each piece of evidence, and average these² to obtain a logit for risk prediction:

$$p(\hat{y}_i = 1|e_{1:E}) = \sigma(b_i + w_i \cdot (\frac{1}{E} \sum_{j=1}^E f_{\theta}^{\text{BERT}}(e_j))) \quad (2)$$

where $w_i \in \mathbb{R}^d$ is the embedding of diagnosis i , $e_{1:E}$ is the flattened list of evidence snippets with null evidence omitted, f_{θ}^{BERT} is the ClinicalBERT (Alsentzer et al., 2019) [CLS] embedding function (which yields a d -dimensional vector), and $b_i \in \mathbb{R}$ is the bias for diagnosis i . The prior over conditions can be defined as the same equation excluding the evidence term: $p(\hat{y}_i) = \sigma(b_i)$, and the **relative risk** follows as $p(\hat{y}_i|e_{1:E})/p(\hat{y}_i)$.

While the bias could be learned, we instead simply set it to the inverse sigmoid of the observed prevalence of the disease in the training sample distribution: $b_i = \sigma^{-1}(\text{prevalence}_i^{\text{train}})$. This means that if we wanted to transfer the model to a new population, where the prevalence differed but the contributions of different evidence were assumed to remain, we could simply update the b_i term.

Excluding interactions between evidence snippets is a sacrifice in model complexity, but it also allows us to compute an interpretable “vote” for any individual piece of evidence as

$$p(\hat{y}_i|e_j) = \sigma(b_i + w_i \cdot f_{\theta}^{\text{BERT}}(e_j)) \quad (3)$$

and compute an individualized relative risk for each piece of evidence using this value.

Conveniently, forcing the bias term to be the inverse sigmoid of the training prevalence, by definition, also means we can interpret the evidence term in Equations 2 and 3 as the **log odds ratio**, i.e., the difference between the logits when conditioning vs. not conditioning on the evidence. The

²Neural Additive Models typically use a sum instead of an average, but we found that given varying amount of evidence retrieved, it worked better to use an average.

model is effectively estimating this log odds ratio directly. This variable’s expected value does not change if we sample conditions for training with a frequency different from the the natural prevalence of the conditions (Simon, 2001). Because of this, we can estimate the likelihood and the relative risk during inference on a differently sampled population by simply changing the bias term in the prior and in equations 2 and 3 to reflect the estimate of the natural prevalence of the conditions (Zhang and Kai, 1998), which we can get from the training set before sampling: $b'_i = \sigma^{-1}(\text{prevalence}_i^{\text{train}})$.

3.3 Evidence Re-ranking

Because of the simplicity of the risk prediction, we can use the internal variables it exposes to re-rank evidence. The intuition behind the re-ranking is that the most important evidence will be that which most changes our risk assessment from the prior over the diagnoses, and we would like the chosen metric to capture this across all of the potential diagnoses. We use Mean Squared Error (MSE) of the predicted logits with the logits of the prior $p(y)$. This makes the formulation of the MSE metric simple as the mean (over Q conditions) of the squares of the log odds ratio for a piece of evidence:

$$\text{MSE}(\sigma^{-1}p(\hat{y}|e_j), \sigma^{-1}p(\hat{y})) = \frac{1}{Q} \sum_{i=1}^Q (w_i \cdot f_{\theta}^{\text{BERT}}(e_j))^2. \quad (4)$$

It is necessary to use the *log odds* ratio term in this score function because we care not only about increasing but also about decreasing the probability of a condition, so it makes most sense to compare and sum these two different effects in log space. The reason to choose MSE over other scores (e.g. the absolute distance) comes from the intuition that it is more important to see the evidence that is “very opinionated” about one condition rather than to see evidence that is “slightly opinionated” about many. Therefore, it is necessary to square this log odds ratio before averaging across conditions to reflect this idea when sorting evidence.

4 Certain Diagnosis Extraction

We make an assumption about the EHR of patients that eventually receive a diagnosis that there is some period of time in the record where a diagnosis is “uncertain” before it becomes “certain”, and the eventual “certain” diagnosis is correct. Of course

285 just because a diagnosis is definitive as noted by
286 clinician in the record does not necessarily mean
287 that it is correct—sometimes clinicians are wrong.

288 However, it is hard to detect such cases, so here
289 we focus on reducing delayed diagnosis errors
290 where we assume some evidence in the medical
291 record from that “uncertain” period could have in-
292 fluenced a clinician to make a diagnosis or order a
293 certain kind of test sooner than they did, or keep
294 a diagnosis in the running list of differentials for
295 longer. If notes are incorporated into the input
296 where the diagnosis is already certain, the predic-
297 tion problem becomes too easy, which is why a
298 time-wise fine-grained label is necessary—such a
299 label could more accurately weed out all of this ob-
300 vious evidence. To extract these certain diagnoses
301 with an LLM, we use three sequential prompts and
302 a normalization step.

303 4.1 3-Stage Extraction with LLMs

304 In this section we describe the prompts for certain
305 diagnosis extraction, which are shown in full in the
306 appendix (Section B). Following prior work (Ahsan
307 et al., 2023), we first prompt the LLM with a binary
308 question asking if there exists a confident diagnosis
309 for a patient. If the answer is “yes”, we then ask the
310 model for the diagnoses. Unfortunately, creating
311 a list of diagnosis terms from the answer to this
312 prompt is not just a matter of parsing because we
313 found that the model will often return extended
314 phrases that are not easily mapped to diagnoses.
315 Therefore, we issue one more prompt that only
316 takes in the output of the previous prompt to create
317 a structured list of diagnostic terms. We then parse
318 this final output of the LLM into a list of strings.

319 4.2 Normalization

320 To normalize produced diagnostic terms, we take
321 a two-step approach. First we use string matching
322 heuristics to handle easy cases. Then we embed
323 sentences with SentenceTransformers (Wang
324 et al. 2020; Reimers and Gurevych 2019; specif-
325 ically, all-MiniLM-L6-v2) and calculate cosine
326 similarities, matching a term in the parsed list to
327 the most similar term (with similarity $>.85$) in the
328 predefined set (“cancer”, “pneumonia”, and “pul-
329 monary edema”). We ignore terms with no match.

330 5 Evaluation

331 Because our targets are synthetically generated us-
332 ing an LM, we first evaluate how well our labels

333 align with the “ground truth” (Section 5.1). Next,
334 we aim to evaluate how well the model can real-
335 istically help with risk prediction. Though it is
336 straightforward to assess the accuracy of the risk
337 prediction itself—we use the standard metrics of
338 precision, recall, F1 and AUROC scores to com-
339 pare to various uninterpretable baselines—it is not
340 as easy to assess what we really care about: How
341 helpful is the interpretability offered by the pro-
342 posed model to clinicians (section 5.2)? For this
343 we resort to manual evaluation by our clinical co-
344 authors and develop bespoke interfaces to facilitate
345 annotation.

346 5.1 Future Target Extraction

347 To evaluate how well the LLM extracts targets in
348 the form of “confident” diagnoses, we enlist our
349 clinical collaborators to annotate the precision with
350 which the LLM infers “confident” diagnoses. In
351 particular, for every report where one of the three
352 diagnoses—cancer, pneumonia, and pulmonary
353 edema—was automatically extracted, an ICU clini-
354 cian is first tasked with answering the question “Is
355 [diagnosis] a confident diagnosis of the patient ac-
356 cording to the report?”. If the answer is “yes”, they
357 are asked: “Is it likely that this confident diagnosis
358 could be identified in earlier reports?”.

359 5.2 Risk Prediction Interpretability

360 To assess the viability of clinicians using this model
361 in practice, we collect in-depth annotations in-
362 tended to simulate the real-world use of this tech-
363 nology. We evaluate a number of baseline models
364 and model ablations to assess the relative benefits
365 of different model components.

366 **Interface and Annotations** To conduct annota-
367 tions, we develop an interface that simulates as
368 closely as possible the envisioned use case: A clini-
369 cian is seeing an ICU patient’s chart for the first
370 time and trying diagnose the patient or determine
371 what they are at risk of. The clinician may not
372 have much time to spend with the patient’s chart,
373 so we ask clinician annotators to work quickly—
374 specifically, to try and keep annotation time to a
375 few minutes—and we record the amount of time
376 they take to review the patient’s record. When they
377 are done, the annotation process starts, and though
378 they are allowed to access the patient’s notes, they
379 are encouraged not to.

380 We first ask if a diagnosis is noted explicitly in
381 the patient’s record. Given that we are aiming to

evaluated records where the diagnosis is not yet clear, we skip the rest of the annotations on the instance if a diagnosis is explicit. If not, we ask for estimates of the likelihood (“unlikely”, “somewhat likely”, or “very likely”) of each of the possible conditions. Note that we explicitly do not show any model predictions until after this question, to avoid bias. Then, we show the annotator the model predictions and ask if the predicted risk for the conditions aligns with intuition.

Moving onto the evidence (appendix Figure 13), we allow the annotator to look at the sorted evidence one snippet at a time along with the individualized risk prediction only based on that snippet. The annotator notes the usefulness of the evidence with respect to each condition. If the evidence is useful, they are asked whether or not the impact of this evidence on the risk scoring (for the particular condition) aligns with intuition, and whether the annotator remembers seeing this piece of evidence during their initial review of the patient’s notes. After two pieces of evidence, if the annotator feels like more evidence is needed to form a reasonable opinion of the patient’s risk, they can request more evidence snippets (up to a maximum of 10), annotating each as they go. Finally, the annotator is asked if any of the evidence presented impacted their original assessment of likelihood.

Ablations While the task of risk prediction is standard, there is less work on the the task of surfacing relevant evidence (abstracted or extracted) to support such predictions. Consequently, there is not a large set of baselines to serve as natural comparators to our approach. Therefore, in our analysis we focus on showing the importance of each component of our model through ablations. We can decompose our approach into two evidence retrieval components, generating the evidence, which we refer to as “LLM Evidence” and reranking it, which we refer to as “Log Odds Sorting”. The following ablations show the importance of both of these components in identifying useful evidence.

We use prior work (Ahsan et al., 2023) as a starting point for generating the evidence, so it is natural to ask what that component can do by itself without re-ranking using the risk prediction scores for each piece of evidence. A natural comparison is to present the same evidence retrieved but in a *random* or *reverse chronological* order (as recency is probably important). But we can also use the model certainty in evidence, given that this has

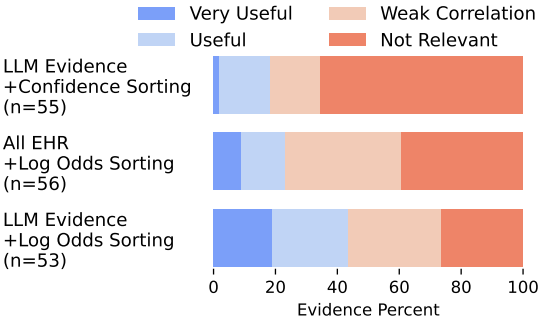


Figure 3: **Evidence Usefulness** (the maximum score across conditions) for our approach and two ablations. “LLM Evidence+Confidence Sorting” uses model evidence, but sorts by (length-normalized) log probability instead of the log odds. “All EHR+Log Odds Sorting” does not use LLM evidence and instead takes the last 1000 sentences in the record as evidence.

been shown to correlate with the utility of snippets (Ahsan et al., 2023). We adopt this approach for comparison and call it “**Confidence Sorting**”.

It is also natural to question the importance of using the language model to abstractively generate evidence at all. We might instead simply use every sentence in the report as evidence and train our prediction model with this retrieved evidence, re-ranking it in the normal way (“Log Odds Sorting”) with the prediction model’s scores. We call this the “**All EHR**” model.

6 Results and Discussion

The majority of our results are based on annotations from 4 annotators on 24 instances and 3 models. Each instance has a maximum of 3 annotators, each annotating different models (assigned randomly). Table 2 reports detailed statistics.

Our main goal is to understand if our approach can retrieve better evidence. To this end, we plot the percentage of evidence annotated in each category of usefulness for each model in Figure 3. Though we record usefulness for each condition individually, here we combine these annotations by taking the maximum score across the conditions for each piece of evidence. The results highlight the necessity of both the “LLM Evidence” retrieval component and the “Log Odds Sorting” method, as both other variants retrieve significantly less “Useful” and “Very Useful” evidence and more “Weakly Correlated” and “Not Relevant” evidence.

How much of the relevant retrieved evidence is

Annotator	LLM Evidence+Confidence Sorting				All EHR+Log Odds Sorting				LLM Evidence+Log Odds Sorting			
	Inst.	Evid.	Rep.	Percent Useful	Inst.	Evid.	Rep.	Percent Useful	Inst.	Evid.	Rep.	Percent Useful
1	8	20	195	5.0	5	14	81	7.1	6	13	154	30.8
2	2	6	26	50.0	2	5	72	40.0	5	14	162	50.0
3	4	13	105	23.1	6	17	224	35.3	5	14	119	50.0
4	5	16	132	18.8	6	20	127	20.0	4	12	85	41.7
Aggregated	19	55	458	24.2	19	56	504	25.6	20	53	520	43.1

Table 2: Annotations. We report the statistics for the number instances annotated, the amount of evidence snippets annotated, the total number of reports in the annotated instances, and the percent of evidence annotated as “Useful” and “Very Useful”. Aggregated statistics are computed by summing over the annotators except in the case of “Percent Useful”, where scores are macro-averaged over annotators. (This is slightly different from Figure 3 where percentages are macro-averaged, i.e., we combine all annotated evidence).

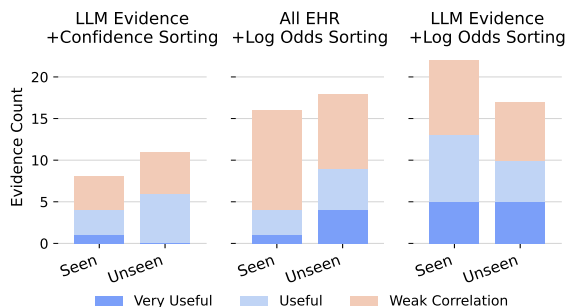


Figure 4: **Seen vs. unseen** evidence counts for all evidence that at least weakly correlates with a condition.

464 redundant with the information already uncovered
465 during the annotator’s initial review of the patient?
466 We plot evidence counts separately for seen vs un-
467 seen evidence in Figure 4 and find that there is a
468 significant amount of unseen evidence that is useful
469 and very useful in all models.

470 The rated usefulness of evidence does not nec-
471 essarily matter if it does not affect the clinician’s
472 decision. An example of how these models might
473 work in practice is when our LLM Evidence model
474 with Confidence Sorting surfaced the following:
475 “Atrial fibrillation with rapid ventricular response.
476 Compared to the previous tracing atrial fibrillation
477 is seen. Other findings are similar. The patient
478 is at risk of pulmonary edema.” In this case the
479 annotator changed their estimate of the likelihood
480 of pulmonary edema from unlikely to somewhat
481 likely, and it turns out that pulmonary edema did
482 appear in a future report.

483 We show all 7 instances where annotators
484 changed their mind after viewing evidence in Ap-
485 pendix Table 5. Of these we find 2 instances (in-
486 cluding the example above) where annotators’ in-
487 creased their likelihood of conditions that were
488 extracted from future records, and 5 where condi-
489 tion(s) other than the synthetically labeled condi-
490 tion(s) were affected (mostly by increasing the

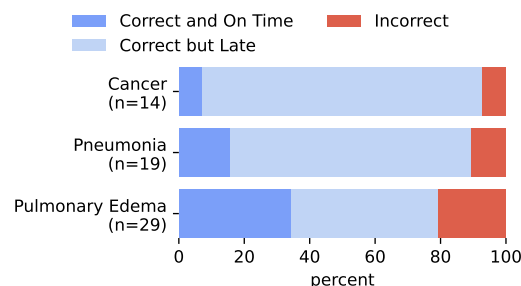


Figure 5: Synthetic label precision. For each confi-
dent diagnosis label extracted by the system, annotators
check whether the diagnosis actually appears in the re-
port (and is definitive), and subsequently if subjectively
they believe that report is likely the *first* time the diag-
nosis was definitive based on the report language.

491 annotators’ risk assessments). Though more data
492 should be collected, this indicates the model might
493 improve annotator recall (though at some cost in
494 precision); recall is arguably more important here.

495 Given that we are using synthetic labels of future
496 diagnoses for both training and evaluation for risk
497 prediction (discussed next), it is important to eval-
498 uate how well our labels align with ground truth.
499 Given that ICD codes are not fine-grained enough
500 and are not always accurate, we turn to manual an-
501 notations of precision for this evaluation. In Figure
502 5, we report the precision of these labels for being
503 correct or for being “correct and on time”. This
504 second category is a stronger correctness in which
505 the annotator also noted that the note where the
506 label was detected subjectively seems to be the first
507 note where that label should have been given as
508 judged using the phrasing in the note.³

509 We see reasonable precision when using auto-
510 matic labeling with the LLM pipeline (about 80
511 percent and above for all conditions). We also

³It would be time-consuming to annotate this directly be-
cause it involves looking at a lot of prior notes.

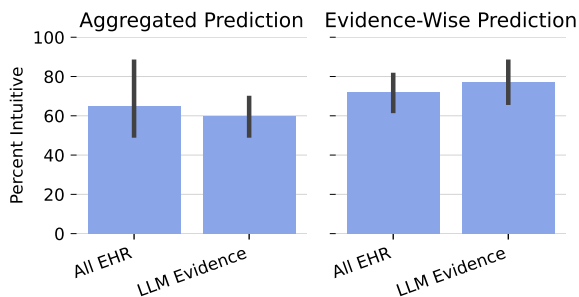


Figure 6: Intuitiveness of predictions macro-averaged across annotators.

512 compute inter-annotator agreement for these an-
 513 notations of precision across the 4 annotators by
 514 enforcing that 8 annotated predictions overlap for
 515 all the annotators. The Fleiss’ Kappa score for
 516 these synthetic label annotations was .68 for the
 517 3-category classification shown in Figure 5 and .86
 518 for the 2-category classification obtained by simpli-
 519 fying the labels into just “Correct” or “Incorrect”.

520 We would also like to assess how well our mod-
 521 els’ risk estimates aligns with the intuitions of clin-
 522 icians with respect to the aggregated and individual
 523 predictions. Though for the aggregated prediction
 524 for an instance, we ask annotators to take the magni-
 525 tude of the risk, not just the direction (i.e. increased
 526 compared to baseline or decreased compared to
 527 baseline) into account, for evidence-level predic-
 528 tions, we ask annotators to take the magnitude with
 529 a grain of salt and mostly judge based on the direc-
 530 tion. This is because the magnitudes appeared to
 531 be somewhat artificially inflated potentially either
 532 due to the strong evidence trying to “componsate”
 533 for the evidence that does not actively contribute
 534 to the log odds (see Figure 12) or because of the
 535 sorting method.⁴ Figure 6 shows that both models
 536 do reasonably well with respect to the aggregated
 537 prediction and the evidence-wise predictions, and
 538 both do slightly better on evidence-wise predictions
 539 than aggregated predictions.

540 Finally, it is important to evaluate the actual pre-
 541 diction performance of our models on our synthetic
 542 labels. Here we also compare against baseline mod-
 543 els that are not interpretable: BERT and Long-
 544 former. These black-box models are trained on
 545 both the All EHR and the concatenated retrieved
 546 LLM evidence. Figure 7 shows that including all
 547 evidence usually helps prediction performance, but

⁴Future work might investigate how to bring make this *magnitude* more interpretable.

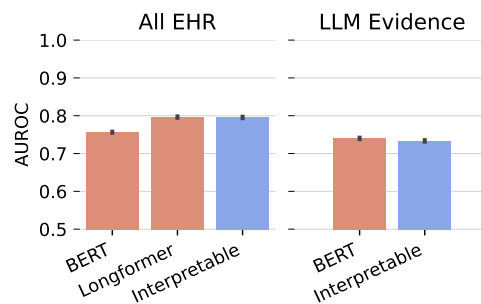


Figure 7: Risk Prediction Performance evaluated on synthetic labels and averaged over 5 random seeds for choosing the which time-point in the EHR to use prior to the diagnosis label. Error bars represent standard deviation of the random seeds. Here, BERT and Longformer refer to Clinical BERT and Clinical Longformer.

548 using the blackbox vs interpretable models on the
 549 same input does not effect performance.

7 Conclusions 550

551 Clinicians should have access to all the pertinent in-
 552 formation to make well-grounded decisions for di-
 553 agnosing a patient, but currently they are inundated
 554 with (unstructured) information from the EHR.
 555 This is exacerbated by the time constraints faced
 556 by practitioners. We have proposed an approach
 557 that aims to facilitate efficient access to potentially
 558 important data within EHR; our method capitalizes
 559 on the capabilities of LLMs to produce digestable,
 560 abstractively generated text evidence, which is then
 561 consumed by a Neural Additive Model (NAM) to
 562 yield a prediction.

563 We find that using NAMs does not sacrifice pre-
 564 dictive quality, but does enable models to surface
 565 useful evidence to clinicians. Using the LLM to
 566 create the starting set of evidence to feed into the
 567 NAM does sacrifice some performance, but it also
 568 significantly increases the usefulness of the evi-
 569 dence in comparison with using the raw sentences
 570 from EHR notes as evidence.

571 Further, we find that in some cases the surfaced
 572 evidence is able to change a clinician’s mind, in-
 573 creasing the clinician’s recall though decreasing
 574 precision, which warrants future work to improve
 575 on this system. One major concern is that this
 576 type of system could increase clinician’s workload
 577 rather than decrease it. Future work should assess
 578 exactly how and when it might be beneficial to
 579 show snippets to clinicians.

8 Limitations

Though the proposed approach of combining abstractive LLM evidence with Neural Additive Models shows promise, there are still many concerns that need to be addressed in future work. One of the biggest concerns is about the use of abstractive “evidence” produced by LLMs. Hallucinations in this evidence could at best negatively impact trust of clinicians in the system and at worst mislead clinicians and negatively affect patient outcomes. Our work does not directly study this given the substantial extra annotator time needed to check for hallucinations. We also did not experiment much with different prompts or models for producing this evidence given that our main focus was on validating the system-level approach rather than individual components.

Another limitation concerns the lack of a significant number of baseline models. Though not many baselines exist for a task that involves retrieving evidence supporting predictions in EHR, there are still potential baselines that use relevance weights or cosine similarity with clinical BERT that we could have included. However, due to the extensive amount of time needed for just one annotation on one model, we chose to focus on ablating over the LLM evidence retrieval and sorting method components of the model.

Finally, our analysis mostly relies on a relatively small amount of annotations from one dataset. This again stems from the time cost of annotations. Each annotator must first look through a whole patient’s record to get a sense of the patient before even getting to any annotations. On average, this took almost 3 minutes, which is all before annotators even see any of the questions. Then, because the study focuses on just the top evidence presented for each instance, each annotator only annotates 3.2 evidence snippets on average per instance. This time-consuming process did limit the number of annotations we could obtain.

References

Rishabh Agarwal, Nicholas Frosst, Xuezhou Zhang, Rich Caruana, and Geoffrey E. Hinton. 2020. [Neural additive models: Interpretable machine learning with neural nets](#). *ArXiv*, abs/2004.13912.

Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. [Large language models are few-shot clinical information extractors](#).

In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 1998–2022, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Hiba Ahsan, Denis Jered McInerney, Jisoo Kim, Christopher Potter, Geoffrey Young, Silvio Amir, and Byron C. Wallace. 2023. [Retrieving evidence from ehRs with llms: Possibilities and challenges](#).

Emily Alsentzer, John Murphy, William Boag, Weihung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. [Publicly available clinical BERT embeddings](#). In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Emily Alsentzer, Mary-Jette Rasmussen, Raíssa Schmitt Fontoura, Andrew Cull, Brett K. Beaulieu-Jones, Kathryn J. Gray, D. Bates, and Vesela P. Kovacheva. 2023. [Zero-shot interpretable phenotyping of postpartum hemorrhage using large language models](#). *NPJ Digital Medicine*, 6.

Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015. [Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission](#). In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’15*, pages 1721–1730, New York, NY, USA. Association for Computing Machinery.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *arXiv preprint arXiv:2210.11416*.

Benjamin A Goldstein, Ann Marie Navar, Michael J Pencina, and John P A Ioannidis. 2016. [Opportunities and challenges in developing risk prediction models with electronic health records data: a systematic review](#). *Journal of the American Medical Informatics Association*, 24(1):198–208.

Alistair E. W. Johnson, Tom J. Pollard, and Roger G. Mark. 2016a. [MIMIC-III clinical database \(version 1.4\)](#).

Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Liwei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016b. [MIMIC-III, a freely accessible critical care database](#). *Scientific Data*, 3(1):160035.

Yikuan Li, Mohammad Mamouei, Gholamreza Salimi-Khorshidi, Shishir Rao, Abdelaali Hassaine, Dexter Canoy, Thomas Lukasiewicz, and Kazem Rahimi. 2021. [Hi-behrt: Hierarchical transformer-based model for accurate prediction of clinical events using multimodal longitudinal electronic health records](#). *IEEE journal of biomedical and health informatics*, 27:1106 – 1117.

686 Yikuan Li, Ramsey M Wehbe, Faraz S Ahmad, Hanyin
687 Wang, and Yuan Luo. 2023. A comparative study
688 of pretrained language models for long clinical text.
689 *Journal of the American Medical Informatics Associ-*
690 *ation*, 30(2):340–347.

691 Denis McInerney, Geoffrey Young, Jan-Willem van de
692 Meent, and Byron Wallace. 2023. **CHiLL: Zero-shot**
693 **custom interpretable feature extraction from clinical**
694 **notes with large language models**. In *Findings of the*
695 *Association for Computational Linguistics: EMNLP*
696 *2023*, pages 8477–8494, Singapore. Association for
697 Computational Linguistics.

698 David E Newman-Toker, Najilla Nassery, Adam C
699 Schaffer, Chihwen Winnie Yu-Moe, Gwendolyn D
700 Clemens, Zheyu Wang, Yuxin Zhu, Ali S. Saber
701 Tehrani, Mehdi Fanai, Ahmed Hassoon, and Dana
702 Siegal. 2023. **Burden of serious harms from diagnos-**
703 **tic error in the usa**. *BMJ Quality & Safety*.

704 Laila Rasmy, Yang Xiang, Ziqian Xie, Cui Tao, and
705 Degui Zhi. 2021. **Med-BERT: pretrained contextual-**
706 **ized embeddings on large-scale structured electronic**
707 **health records for disease prediction**. *npj Digital*
708 *Medicine*, 4(1):86.

709 Nils Reimers and Iryna Gurevych. 2019. **Sentence-bert:**
710 **Sentence embeddings using siamese bert-networks**.
711 In *Proceedings of the 2019 Conference on Empirical*
712 *Methods in Natural Language Processing*. Associa-
713 tion for Computational Linguistics.

714 Thomas Searle, Zina Ibrahim, and Richard JB Dobson.
715 2020. Experimental evaluation and development of
716 a silver-standard for the mimic-iii clinical coding
717 dataset. *arXiv preprint arXiv:2006.07332*.

718 Stephen D Simon. 2001. Understanding the odds ratio
719 and the relative risk. *Journal of andrology*, 22(4):533–
720 536.

721 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan
722 Yang, and Ming Zhou. 2020. **Minilm: Deep self-**
723 **attention distillation for task-agnostic compression**
724 **of pre-trained transformers**.

725 Jason Wei, Maarten Paul Bosma, Vincent Zhao, Kelvin
726 Guu, Adams Wei Yu, Brian Lester, Nan Du, An-
727 drew Mingbo Dai, and Quoc V. Le. 2022. **Finetuned**
728 **language models are zero-shot learners**.

729 Zhichao Yang, Avijit Mitra, Weisong Liu, Dan
730 Berlowitz, and Hong Yu. 2023. **Transformehr:**
731 **transformer-based encoder-decoder generative model**
732 **to enhance prediction of disease outcomes using elec-**
733 **tronic health records**. *Nature Communications*, 14.

734 Jun Zhang and F Yu Kai. 1998. What’s the relative
735 risk?: A method of correcting the odds ratio in cohort
736 studies of common outcomes. *Jama*, 280(19):1690–
737 1691.

738 Laura Zwaan, Martine de Bruijne, Cordula Wagner,
739 Abel Thijs, Marleen Smits, Gerrit van der Wal, and
740 Daniëlle R. M. Timmermans. 2010. **Patient Record**
741 **Review of the Incidence, Consequences, and Causes**
742 **of Diagnostic Adverse Events**. *Archives of Internal*
743 *Medicine*, 170(12):1015–1021.

A Description of terms.

Table 3 shows all of the terms used to describe different models and settings.

B Certain Diagnosis Extraction Prompts

Prompt 1:

Read the following report:

<input>

Question: Is there a confident diagnosis of the patient’s condition? Choice: -Yes

-No

Answer:

Prompt 2:

Read the following report:

<input>

Answer step by step: What is the correct diagnosis of the patient’s condition?

Answer:

We use Chain of Thought (CoT) prompting here because—similar to the evidence retrieval step—we want the model first to extract the parts of the report that refer to a diagnosis, as this seems to work better than going straight to the list of diagnoses. In initial experiments, using the CoT prompt appeared to more easily elicit these verbose extractions.

Prompt 3:

Here is a diagnosis of a patient:

<confident diagnosis>

Question: Provide a list of diagnostic terms or write none.

Answer:

C Prompting Problems

In our 3-stage prompting process, we initially had some problems with false positives in scenarios where pneumonia was negated (Figure 8). We discovered that this was because our 3rd prompt was originally:

Here is a diagnosis of a patient:

<confident diagnosis>

Question: Based on this diagnosis, provide a list of diagnostic terms.

Answer:

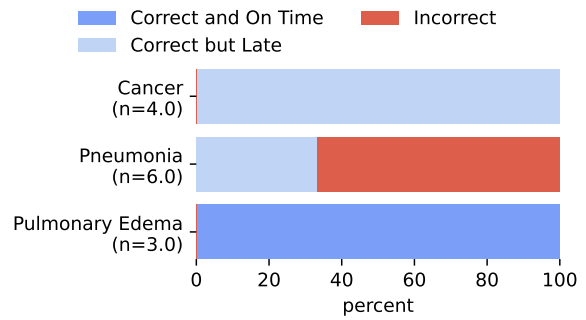


Figure 8: Synthetic labels on validation examples before correcting the prompting problem.

This particular prompt sometimes produced positive synthetic labels for pneumonia when pneumonia was actually negated in the confident diagnosis generated by the previous prompt. We realized this when starting to annotate validation examples, so we changed our prompt (see section 4.1). All of the test annotations reported in the main paper do not include or overlap patients with these annotated validation examples.

D Experiments

We use Clinical BERT for the NAM prediction model. For all models, we train for up to 10 epochs on one Quadro RTX 8000 GPU and pick the best checkpoint (where checkpoints occur every 5 percent of an epoch). For the LLM for both evidence retrieval and synthetic label extraction we use FLAN-T5-XXL (Chung et al., 2022; Wei et al., 2022). In the case of All EHR used as input to the NAM, we split sentences with NLTK. We will make code open-source on acceptance.

E Usefulness of Queries

Unlike (Ahsan et al., 2023), we do not directly evaluate how relevant the retrieved evidence is to the query used to retrieve it; we instead focus on how relevant the evidence is to the risk predictions. However, we would like to examine which queries produce useful evidence. Figure 9 shows counts of evidence in each category separated across which query was used to retrieve that evidence. It seems as though the most useful evidence came from the three queries that directly ask about the condition for which we are predicting risk (the three left-most queries), but a few additional queries sometimes did prove useful.

LLM Evidence	Models that use the evidence retrieved with an LLM.
All EHR	Models that use the all of the text in the EHR. For Interpretable Neural Additive Model, this text is split at the sentence level.
BERT or Longformer	Blackbox models that take either All EHR or LLM Evidence (concatenated) as input. BERT refers to Clinical BERT (Alsentzer et al., 2019) and Longformer refers to Clinical-Longformer (Li et al., 2023).
Interpretable	The proposed Interpretable Neural Additive Model, which can operate either on LLM Evidence or All EHR inputs.
Confidence Sorting	Sorting LLM Evidence by the length-normalized log-likelihood of the evidence under the LLM.
Log Odds Sorting	Sorting either LLM Evidence or All EHR inputs by the mean squared error of the predicted log odds (equation 4).

Table 3: Description of terms.

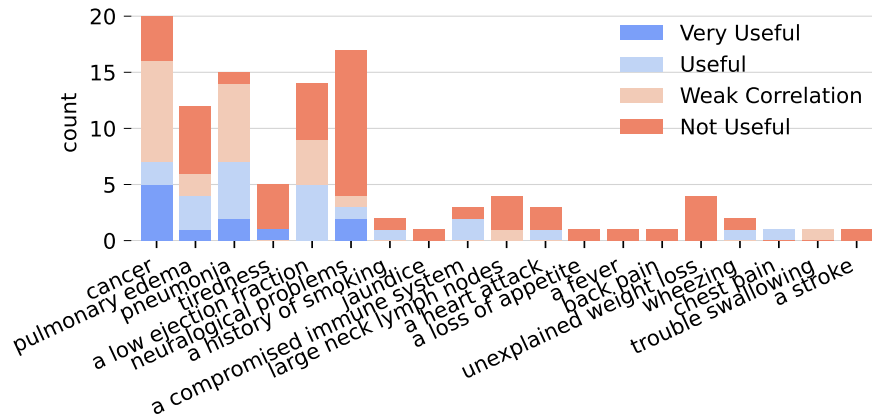


Figure 9: Usefulness per Query

F Full Prediction Performance

We report the full prediction performance in Table 4.

G Annotators Changing Their Minds

Table 5 presents all the occurrences of annotators changing their mind.

H Ablation over amount of evidence used

Figure 10 shows performance if we limit to a set number of evidence that can be used in the Neural Additive Model’s final aggregated score.

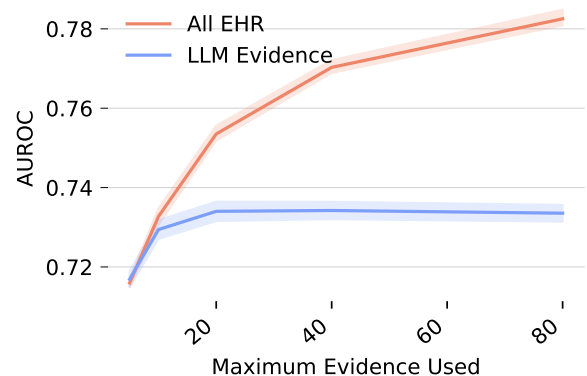


Figure 10: Ablation over amount of evidence used to make a risk prediction.

I Evidence Histograms

Figure 11 shows a histogram of the amount of evidence per each instance, and Figure 12 shows what the distribution over the log odds votes looks like.

J Annotation Interface

Figure 13 shows a screenshot of what the part of the interface dedicated to annotating evidence looks like.

	AUROC	Precision	Recall	F1
BERT (All EHR)	75.6 ± .19	65.6 ± 1.38	16.8 ± .38	26.8 ± .43
Longformer (All EHR)	79.6 ± .22	55.5 ± .32	28.8 ± .43	37.9 ± .38
Interpretable (All EHR)	79.5 ± .23	56.5 ± .57	20.5 ± .58	30.1 ± .60
BERT (LLM Evidence)	74.0 ± .27	51.6 ± 1.32	22.7 ± .27	31.5 ± .42
Interpretable (LLM Evidence)	73.3 ± .27	53.6 ± 1.09	15.0 ± .36	23.4 ± .48

Table 4: **Risk Prediction Performance** on the synthetic labels averaged over 5 different random seeds used for choosing the time-point in each patient that separates the past from the future.

Annotator	Model	Sorting	Changes	Best Evidence	Usefulness	Synthetic Label
2	LLM Evidence	Confidence Sorting	Pneumonia: Unlikely → Somewhat likely	There is a small right pneumothorax. There is extensive consolidation of the right upper lobe. Consolidation in the right lower lobe is mostly located in the superior segment. The left lung is grossly clear. There. Signs: There is extensive consolidation of the right upper lobe. Consolidation in the right lower lobe is mostly located in the superior segment. The left lung is grossly clear. There is no left pleural effusion. There is	Useful for Pneumonia	Pneumonia
4	LLM Evidence	Confidence Sorting	Pulmonary Edema: Unlikely → Somewhat likely	Atrial fibrillation with rapid ventricular response. Compared to the previous tracing atrial fibrillation is seen. Other findings are similar. The patient is at risk of pulmonary edema.	Useful for Pulmonary Edema	Pulmonary Edema
3	All EHR	Log Odds Sorting	Cancer: Unlikely → Very likely	Basal cell skin ca. [**27**].	Useful for Cancer	Pulmonary Edema
4	All EHR	Log Odds Sorting	Cancer: Unlikely → Somewhat likely	o.b.resident to see pt., pt.waiting for a "biopsy".	Useful for Cancer	Pulmonary Edema
4	All EHR	Log Odds Sorting	Pulmonary Edema: Somewhat likely → Unlikely, Pneumonia: Somewhat likely → Very likely	There is increased opacity in the. retrocardiac left lower lobe, as well as the right lower lobe, which could be. due to atelectasis, aspiration, or possibly pneumonia.	Very Useful for Pneumonia	
1	LLM Evidence	Log Odds Sorting	Pneumonia: Somewhat likely → Very likely	CXR showed L middle/lower lobe PNA, prob asp PNA.	Very Useful for Pneumonia	
4	LLM Evidence	Log Odds Sorting	Cancer: Unlikely → Very likely	CLL. Signs: id: pmh of CLL	Very Useful for Cancer	

Table 5: Examples of the 5 instances where annotators changed their mind based on evidence shown.

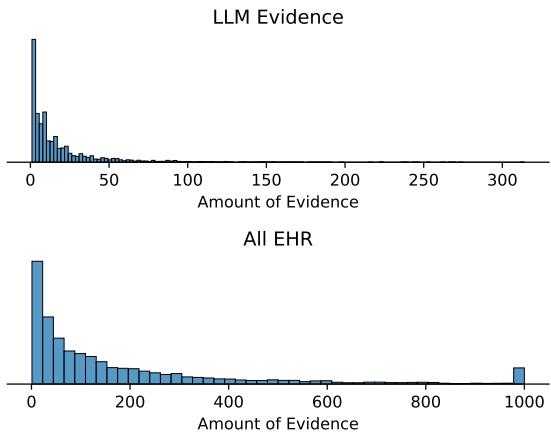


Figure 11: Histogram of the number of text snippets for each instance.

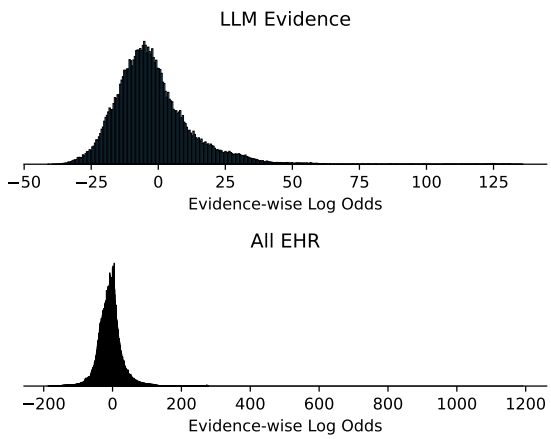


Figure 12: Histogram of the log odds of each individual piece of evidence.

1.

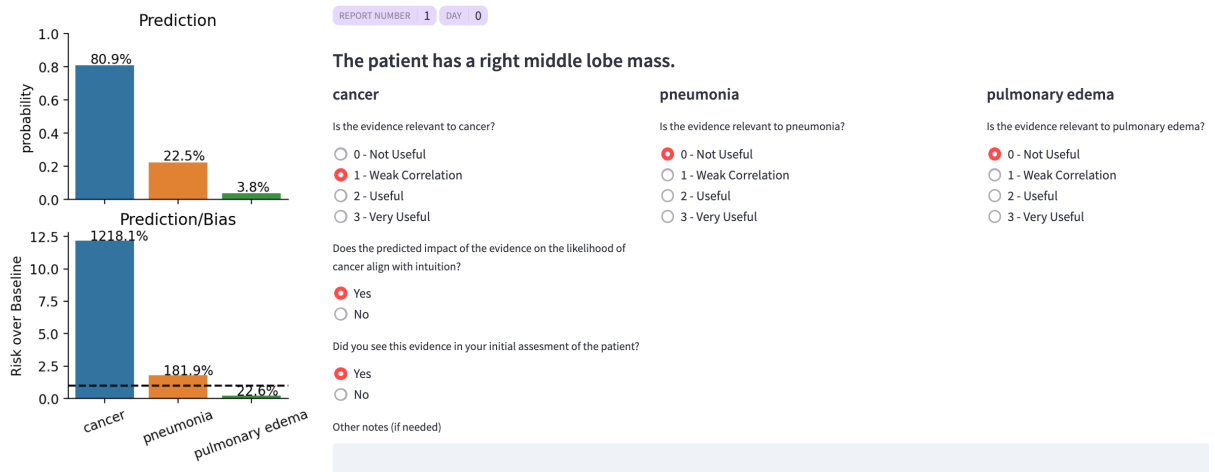


Figure 13: An example part of the **evidence** annotation interface. The plots on the left indicate the predicted likelihood (top) and the odds ratio (bottom).