# Highlights

**Informing Clinical Assessment by Contextualizing Post-Hoc Explanations of Risk Prediction Models in Type-2 Diabetes**

Shruthi Chari, Prasant Acharya, Daniel M. Gruen, Olivia Zhang, Elif K. Eyigoz, Mohamed Ghalwash, Oshani Seneviratne, Fernando Suarez Saiz, Pablo Meyer, Prithwish Chakraborty, Deborah L. McGuinness

- Generate contextual explanations for a comorbidity risk prediction setting.

- Explanations tie predictions and posthoc features to evidence from the context of use.

- Explanations extracted from guidelines using LLMs and knowledge-augmentations (KAs).

- Evaluate 5 clinical LLMs and KAs on an entire dataset and divide them by diseases.

- An expert panel found value in providing such contextual explanations.

# Informing Clinical Assessment by Contextualizing Post-Hoc Explanations of Risk Prediction Models in Type-2 Diabetes

Shruthi Chari[a], Prasant Acharya[a], Daniel M. Gruen[a], Olivia Zhang[b], Elif K. Eyigoz[b], Mohamed Ghalwash[b], Oshani Seneviratne[a], Fernando Suarez Saiz[c], Pablo Meyer[b], Prithwish Chakraborty[b], Deborah L. McGuinness[a]

[a]*Rensselaer Polytechnic Institute, 110 8th St, Troy, 12180, NY, US*
[b]*Center for Computational Health, IBM Research, 1101 Kitchawan Rd
, Yorktown Heights, 10598, NY, US*
[c]*IBM Watson Health, 75 Binney St, Cambridge, 02142, MA, US*

## Abstract

Medical experts may use Artificial Intelligence (AI) systems with greater trust if these are supported by 'contextual explanations' that let the practitioner connect system inferences to their context of use. However, their importance in improving model usage and understanding has not been extensively studied. Hence, we consider a comorbidity risk prediction scenario and focus on contexts regarding *the patients' clinical state, AI predictions about their risk of complications, and algorithmic explanations supporting the predictions.* We explore how relevant information for such dimensions can be extracted from Medical guidelines to answer typical questions from clinical practitioners. We identify this as a question answering (QA) task and employ several state-of-the-art Large Language Models (LLM) to present contexts around risk prediction model inferences and evaluate their acceptability. Finally, we study the benefits of contextual explanations by building an *end-to-end* AI pipeline including data cohorting, AI risk modeling, post-hoc model explanations, and prototyped a visual dashboard to present the combined insights from different context dimensions and data sources, while predicting and identifying the drivers of risk of Chronic Kidney Disease (`CKD`) - a common type-2 diabetes (`T2DM`) comorbidity. All of these steps were performed in deep engagement with medical experts, including a final evaluation of the dashboard results by an expert medical panel. We show that LLMs, in particular BERT and SciBERT, can be readily deployed to extract some

relevant explanations to support clinical usage. To understand the value-add of the contextual explanations, the expert panel evaluated these regarding actionable insights in the relevant clinical setting. Overall, our paper is one of the first end-to-end analyses identifying the feasibility and benefits of contextual explanations in a real-world clinical use case. Our findings can help improve clinicians' usage of AI models.

## 1. Introduction

Artificial Intelligence (AI) and Machine Learning (ML) have been applied to the medical and healthcare domains for decades [1, 2] but their adoption has been slow due to various aspects including the need for explaining the black-box nature of such methods. AI explainability (XAI) has tried to provide a rationale for model predictions so that subject matter experts (SMEs) can interpret their results [3, 4, 5]. Studies on XAI have shown that users in different contexts require explanations that match their different levels of expertise and focused on their particular goals and needs [6, 7, 8, 9]. Hence, there often is no single solution for XAI pipelines [10, 11, 9], but processes can be followed for catering to the specific needs of the use case and the intended users [12, 13]. The impact of AI in patient-facing applications will increase the need not only for XAI but also for contexts that subject matter experts (SMEs) in the clinical domain are familiar with [7, 14, 15].

**Definition 1.1** (Explanation)**.** An account of the system, its workings, the *implicit and explicit* knowledge it uses to arrive at conclusions in general and the specific decision at hand, that is *sensitive* to the end-user's *understanding, context, and current needs.* [16]

**Definition 1.2** (Contextual Explanation)**.** Explanations that contain context, are often explicit information [17] to characterize the situation of (an) entity(ies), wherein "an entity is a person, place, or object that is considered relevant to the interaction between a user and an application" [18].

2

In recent years, efforts to describe and formalize explanations [19, 9] have identified various dimensions and types for it. Specifically, 'contextual explanations' [20, 9] hold great promise to satisfy aforementioned clinical needs and can improve the adoption of AI methods among clinical workflows. Risk prediction is one of the most important tasks in clinical decision making, and an increasingly important in view of the move toward personalized medicine. [21, 22]. To interpret risk scores, clinicians often consult evidence from different levels of the scientific pyramid [23] to lookup associations that might impact the patient's treatment or future trajectory. For example, questions like those in Tab. 1, are often asked by clinicians when they are trying to understand or use AI model predictions in their practice. Additional contextual information, such as answers to these questions, can help clinicians interpret and trust predictions to take actions. However, current work in risk prediction has often narrowly focused on improving model's accuracy, ignoring the aforementioned needs. Interestingly, several researchers have posited contextual explanations [24, 25, 9], that go beyond post-hoc model explanations to frame the predictions in the context of the applied setting and decisions being made. However, the feasibility of extracting such contextual explanations and the added benefit in an end-to-end setting of clinical relevance has not been studied and forms the focus of this paper. Specifically, we consider how to derive and support contextual explanations from authoritative domain knowledge sources, not already considered by prediction models, that clinicians would typically use to reason through decisions presented to them when dealing with recommendations from learning health systems.

Table 1: Questions that could be asked in clinical use cases around model explanations / predictions, and which can benefit from contextual explanations in the context of use.

| Sample Question |
| --- |
| What treatment can be suggested for this patient who has an increased risk of cardiovascular disease? |
| What other conditions does this patient have that might impact this decision? |
| What was the patient's A1C value when this prediction was made? |
| Why are you telling me that this risk is important? |

Working closely with medical experts, we identify three specific focus ar-

eas on contexts such as *the patients' clinical state, AI predictions about their risk of complications, and algorithmic explanations supporting the predictions.* Multiple data sources are required to extract such contextual explanations, including patient medical records, AI model predictions, and authoritative information around clinical facts and best practices. Medical guidelines are one of the most trusted authoritative sources of information and can provide the required additional context. Here, we thus study the feasibility of extracting answers from guidelines to typical questions from clinical practitioners to satisfy their explainability needs. We can identify this problem to be a question answering (QA) task. In the natural language processing (NLP) domain, the efficacy and ready-availability of state-of-the art (SOTA) deep learning based QA modules, have rendered such QA problems solvable and is being increasingly productivized. In this paper, we thus aim to extract such contextual explanation using SOTA LLM methods. Especially, we aim to study the following questions relative to contextual explanations:

• Feasibility of extracting and generating contextual explanations from authoritative sources: *can we reliably extract contextual explanations from medical guidelines using state-of-the-art QA models? Can knowledge augmentation improve the QA performance?* We use a suite of readily available QA language models, with and without knowledge augmentations, and compare against manually annotated answers to evaluate the extracted contexts from authoritative clinical sources to explain decisions of post-hoc model explainers and risk prediction models.

•Understanding the added benefit of the derived contexts: *does the derived contextual explanations improve model usage by clinicians?* We evaluate usefulness of these contexts from two perspectives: (a) user persona needs and (b) model accuracy. Particularly, we discuss themes that emerged from our conversations with clinicians to understand if clinical contexts can better support model explanations and what more is desired of contextual explanations.

• Practical considerations in a clinical workflow: *what considerations and challenges might we face in implementing support for derived contexts in a setting of clinical relevance?* To conduct our study, we developed an end-to-end system including (i) data cohorting, (ii) AI models for risk prediction, (iii) post-hoc explainers to identify driver of risk, and (iii) a prototype dashboard to present the combined insights from the contextual explanations. Specifically, as a case-study, we considered the problem of predicting and identifying the drivers of risk of chronic kidney disease (CKD) - a common type-2 diabetes

4

(T2DM) comorbidity and extracted contextual explanations for 175 questions of 5 different types. An expert panel of 4 medical experts evaluated these explanations for 20 prototypical patients. Our end-to-end system enabled us to identify practical steps in creating such a pipeline and the steps needed to generalize this for other clinical settings. Also, this enabled us to answer the aforementioned two questions and derive a holistic understanding supporting contextual explanations in a risk prediction setting. We further identify scenarios within this clinical setting where contextual explanations would be most useful and discuss how they would fit into the clinical workflow.

The rest of the paper is organized as follows. First, we provide a brief description of the use case, the data sources, and the motivation and background for contextual explanations (Sec. 2). Next, we provide an overview of our methodology including the end-to-end AI pipeline and prototype dashboard used to conduct the experiments in (Sec. 3). In Sec. 4, we present comprehensive evaluations including quantitative performance of QA language models as well as qualitative analysis of extracted contextual explanations via a thematic analysis of expert panel sessions. We present a detailed analysis of these results in Sec. 5 and answer the aforementioned questions of interest. Finally, in Sec. 6, we present a summary of other related works and contrast our unique contributions, and conclude with general take-aways, including opportunities for future research in Sec 7.

## 2. Motivation and Background

In this section, we provide details on several assumptions and considerations that will be used to describe our methodology and will be crucial in analyzing the experimental results. Here, we describe the content we will use to support user-centered contextual explainability in chronic disease - comorbidity risk prediction settings. In Sec. 2.1, we present a high-level overview of the selected use-case. in Sec. 2.2, we introduce the entities along which we extract contextual explanations (or contextualized entities), which could provide additional information to help clinicians interpret the risk prediction scores and the factors influencing the scores in clinical settings. Finally, in Sec. 2.3, we provide an overview of the datasets used for our study. In particualr, we describe clinical practice guidelines (CPGs) in Sec. 2.3.2. CPGs are considered to be at the highest level of the evidence-based pyramid [23] and is our selected source to derive high-quality clinical context.
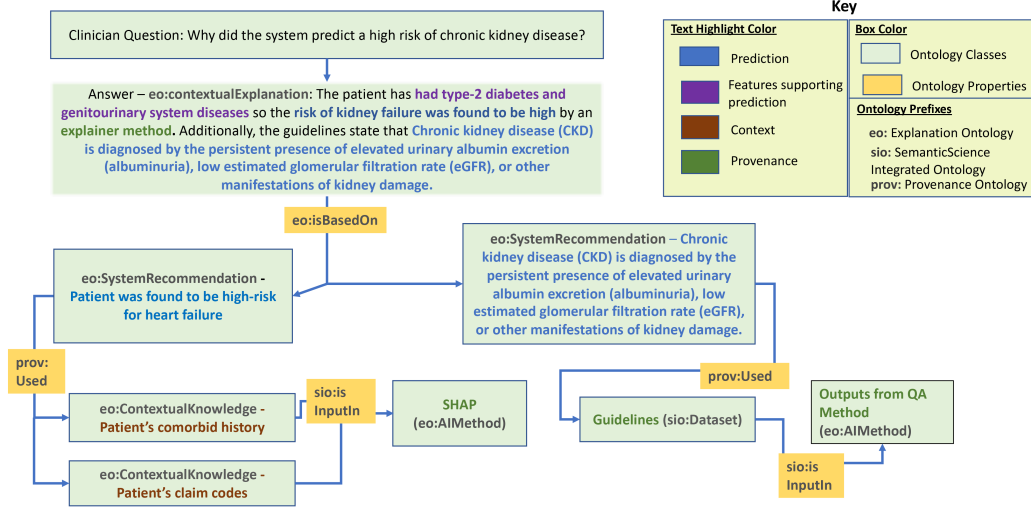
## 2.1. Use Case



Figure 1: Trace of dependencies of a contextual explanation example on system outputs from post-hoc explainer models and question-answering methods. In this example, we have used ontology classes and properties to annotate the data, but in this paper we mainly focus on a multi-method approach to support such contextual explanations.

AI models promises to help clinicians by providing tools for improved decision making. Risk prediction of patients is one of the key steps in a clinical decision scenario and models for such use cases can be consumed by a broad spectrum of clinicians with differing roles and experience, e.g. a specialist vs a primary-care physician. Depending on their roles, the needs, and thereby the desired functionality from a risk prediction contextualization standpoint, can be different. We worked closely with a clinical expert to to understand the clinical use-case and determine the context of AI tools. Crucially, we aimed to form an understanding of the unmet needs and identify relevant contexts that can benefit clinicans. We can further motivate this via Fig. 1 which shows an example question posited by a clinician while consuming the ouputs of a risk prediction model. In this case, the relevant response can be identified via contextual explanations [9] that is generated from multiple sources and via multiple extracted contexts. Our aim, was to thus scope the relevant contexts that will be the focus of our study. We followed a sequence of user-centric research principles [12], to (1) define the scope of our tool's capabilities, (2) identify the end-user/target persona who would most benefit

from our tool, and (3) scope the most relevant contexts. Through our interviews, we identified primary-care physicians (PCP), especially those with lesser years of experience, to be the persona who may most benefit from such contextualizations. We describe the covered context types in Section 2.2. Furthermore, to study our stated problem in a real-world setup, we identified the problem of risk prediction of CKD among new T2DM patients at their first diagnosis.

This is motivated by the fact that diabetes is one of the top five chronic diseases affecting the adult population in the US [26]. Diabetes management involves monitoring for and treating related comorbid conditions. Effective and timely prediction of such conditions can lead to an overall improvement in the quality of care and thus evaluating the impact of AI models in improving clinical decision workflow can have tremendous real-world impact. Especially, we focus on CKD, a commonly occurring micro-vascular complication of T2DM and one of the leading causes of death in the US [27], with an estimated 37 million cases in the US (who are mostly undiagnosed) and cost medicare in 2018 81.1$B$, and end stage renal disease an additional 36.6$B$. Typically, actions to prevent onset of CKD among T2DM patients revolve around proper disease control, including close disease monitoring, proper treatment adherence, and patient education. Incorporating accurate risk prediction of CKD in the clinical workflow can lead to more timely actions, potentially delaying the onset of CKD, and in some cases, preventing its progression. While such predictions could be of use along various time-points of the patients' T2DM prognosis, in this paper, we predict the risk of developing CKD within 360 days of T2DM onset. Under this use case, we explore strategies to provide context around interventions for particular patients, and explain their T2DM state and individual risk factors.

## 2.2. Selected Contextual Entities of Interest

To support the goal of providing user-centered, clinically relevant, and contextual explanations, in consultation with a medical expert on our team who is also a co-author, we identify three entities of interest to provide contextual explanations around predicting the risk of CKD among T2DM patients. Fig. 1 shows an example of contextual explanation that can answer a clinician's question around patient management. It can be seen that such explanations are usually composed of multiple entities and from multiple sources. In general, we identified and subsequently focused on extracting the following contexts:

7

- Contextualizing the patient by connecting their clinical history and indicators to treatments typically recommended for such patients, according to CPGs.
- Contextualizing risk predictions for the patient in terms of the prediction's impact on decisions, based on general norms of practice concerning potential complications, as evident from guidelines and other domain knowledge, including medical ontologies.
- Contextualizing details of algorithmic, post-hoc explanations, such as connecting features that were the most important to other information based on their potential medical significance, such as through connections to physiological pathways and CPGs.
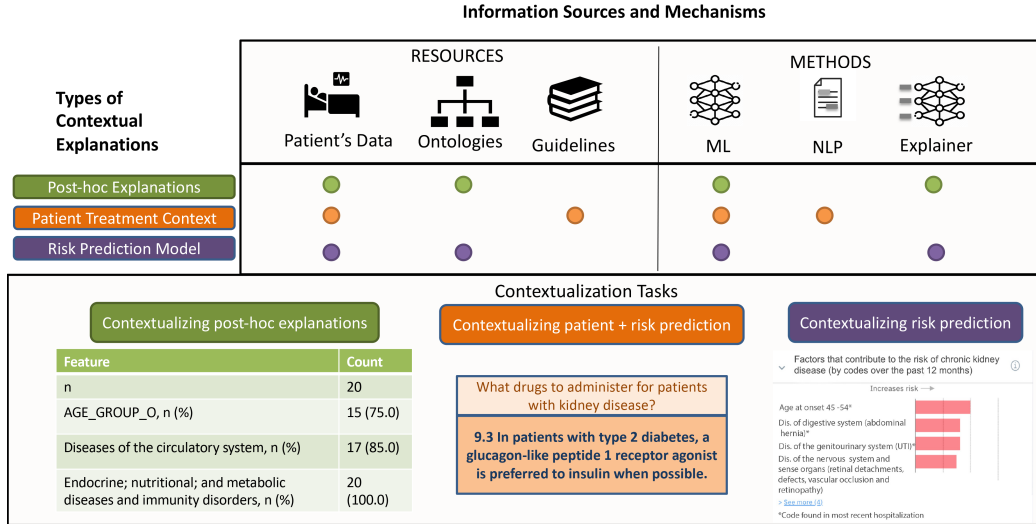


Figure 2: Different types of contextualizations supported by methods, that help provide additional context around patients, their risk predictions and features contributing to risk, via connections to different knowledge sources including patient data, medical ontologies and guidelines.

In Fig. 2, some examples of contexts that we support around the three entities of interest in the risk prediction setting can be seen. Also seen in the figure are the pathways in which answers providing context could borrow from different domain knowledge sources and methods. For example, the answer to the question, "What drugs to administer for chronic kidney disease?" provides context around the patient and risk prediction, borrows both from guidelines and patient data, and is supported by the risk prediction and

8

natural language modules which we describe in Sec. 3.

## 2.3. Data Sources

To conduct our real-world study, we focus on two specific sources of data as described below.

### 2.3.1. Patient Data

We conduct our analysis on and retrieve patient data from the claims sub-component of the Limited IBM MarketScan Explorys Claims-EMR Data Set (LCED), covering both administrative claims and EHR data of over 5 million commercially insured patients between 2013 and 2017. Medical diagnoses are encoded using International Classification of Diseases (ICD) codes. We selected only those `T2DM` patients (with ICD9 codes 250.*0, 250.*2, 362.0, and ICD10 code E11) that satisfied the following criteria as our cohort. Only `T2DM` patients with the following criteria are included:

- have had two or more visits with `T2DM` diagnosis,

- were enrolled continuously for 12 months prior to `T2DM` diagnosis,

- number of visits for `T2DM` is greater than those for other forms of diabetes such as `T1D`, and

- age at the initial `T2DM` diagnosis is between 19-64 years.

Among `T2DM` patients, we use the first diagnosis of chronic kidney disease (`CKD`) (ICD10 N18 or ICD9 585.*, 403.*) after the initial diagnosis of `T2DM` as the outcome to predict. At the time of the first `T2DM` diagnosis, we predict the risk of the patient developing `CKD` within 1 year using Clinical Classifications Software (CCS) codes, age group, and sex as features for the predictive model.

### 2.3.2. Clinical Practice Guidelines

Clinical Practice Guidelines are position statements published by a board of experts in different disease areas [28]. These guidelines are updated often, latest summaries of updated evidence in the disease areas, and follow the highest standards of evidence appraisal (e.g., Grading of Recommendations, Assessment, Development and Evaluations (GRADE) evidence schemes [1]).

---

[1] https://www.gradeworkinggroup.org

Further, the guidelines are written to be comprehensive sources covering different aspects of treatment, management, and assessment of the disease and are often regarded as first-line lookup sources for clinicians and primary care physicians [29, 28]. Given their comprehensive and updated nature, CPGs provide a great resource for providing clinical contexts in various clinical settings. We utilize the 2021 edition of the American Diabetes Association (ADA) Standards of Care guidelines for our experiments.

## 3. Methods

To study the problem of risk prediction of `CKD` among `T2DM` patients, we created an end-to-end AI enabled system. Fig. 3 shows a conceptual overview of the components of this system. In general, to extract contextual explanations around our three identified entities of interest, we used a number of components including risk prediction models, post-hoc explanation models, and our multi-method, question-answering approach to provide context. Crucially, to analyze the importance of the supported contextual explanations, we prototyped a clinical-friendly dashboard and ran qualitative analysis. In this section, we provide high-level details of some of the key components involved in the process.
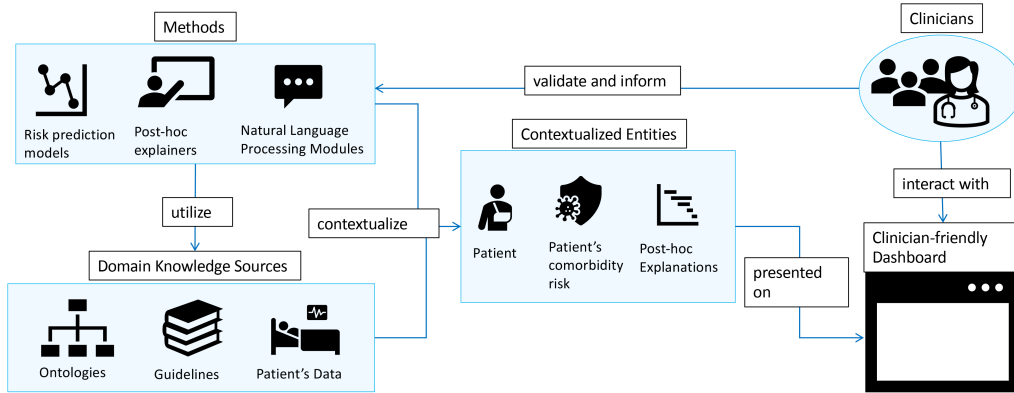


Figure 3: Overall view of the different methods in our pipeline and how they interact to provide risk prediction scores, factors contributing to the risk and contexts around the patient, their predicted risk and the factors contributing to the risk.

## 3.1. Risk Prediction Models

In the first step of our pipeline, we build risk prediction models from the constructed cohort and the aforementioned use-case (Sec. 2.1). In particular, we train a suite of machine learning models (ML), including both classical and deep-learning models, and select the best performing one based on the highest predictive accuracy and other appropriate metrics for the use case, such as favoring models with a higher recall. We used DPM 360 [30], an open-source, reusable, disease progression model training package, we compared a suite of classification models on the patients' demographic and diagnosis history to predict future complications. In this paper, we only used the demographic and diagnostic features to model risk. Furthermore, to handle the temporal features, for some of our models, such as Logistic Regression (LR) and Multi-layer perceptron (MLP), we used temporally aggregated features (summation). We also compared two state-of-the art Recurrent Neural Networks (RNN) where temporal history can be handled in a more natural manner, viz., Long-Short Term Memory (LSTM) and Gated Recurrent Units (GRU). All model implementations are available via DPM360 including classical ML models (backed by scikit-learn) and deep learning models (custom built for DPM). In this paper, we split the data according to a train-validation-test split (70-10-20). Using the best performing models on the validation set, we present our results on the hold-out test set. Since the data is imbalanced, we selected the models based on the best AUC-ROC and AUC-PRC from the validation set. We also evaluate the models based on precision, recall, and brier score [31]. Deep learning networks are known to be under-calibrated and the last metric measures how well the model is calibrated, i.e., it measures the probabilistic interpretation of risk prediction. In other words, if a model predicts a 0.7 risk for a patient, brier score measures how well that translates to a 70% chance of the patient developing the complication. The hyper-parameters for the deep-learning models were selected using a grid search strategy varying batch sizes $\{8, 16, 32, 128\}$, number of layers $\{1, 2, 3\}$, and dropout $\{0.0, 0.1, 0.2\}$ along with standard initialization and using ADAM as the optimizer of choice.

## 3.2. Post-hoc Explainer Models

While some of the classical algorithms considered in Section 3.1 are inherently interpretable with easy access to the features deemed important for the model (such as LR), several of the deep learning models are black-box models. To extract feature importances from such models, we used post-hoc

11

explainers which have been found to be favored by clinicians in past studies [15]. In particular, we used the well accepted SHAP algorithm [32] to find feature importance [2]. The algorithm uses game-theoretic principles to identify importance of features by ascertaining the dip in performance of the model with and without access to the feature at the personalized level. Such personalized feature importance is key so that our overall risk prediction presentations are more actionable for the clinicians by allowing them to focus on the particular attributes of the patients that are driving their risk.

Typically, clinician time is costly and hard to obtain. Thus to conduct the expert panel sessions and let them focus on some of the most 'interesting' patients, we apply Protodash [33] to select a subset of patients. Protodash is a post-hoc sample selection method used to obtain a set of prototypical or representative patients from the high risk category that naturally spans the varied set of patient characteristics for the selected sub-group. This also allows the clinician to build trust in using the AI models by inspecting the different patient modalities of the dataset without having to inspect the entire dataset.

### 3.3. Extracting Contextual Explanations from Clinical Guidelines

We intend to support a set of possible clinical questions around risk prediction setting for patients. The explanations to these questions can help provide more context around patient predicted risk, features contributing to it and data. Each of these question sets, or question types as seen in Tab. 2, can be addressed by multiple sources. Critically, we set up our problem of extracting context around entities of interest in a risk prediction setting from clinical guidelines to help clinicians make sense of comorbidity risk prediction scores of chronic disease as a question-answering (QA) task. Figure 4 shows a detailed overview of the steps involved in this 'guideline QA' task. In this section, we give a brief overview of the important steps. For a detailed description please refer to Appendix A.

Information Retrieval from Data Sources: We support the extraction of context from three domain sources in our QA approach, including patient data,

---

[2]In this paper, our primary goal is study the importance of contextual explanations and thus we chose SHAP as a well-known SOTA post-hoc explainer. However, we caution the readers about known criticism of SHAP, and in general explainability methods, that are still active area of research without a common consensus
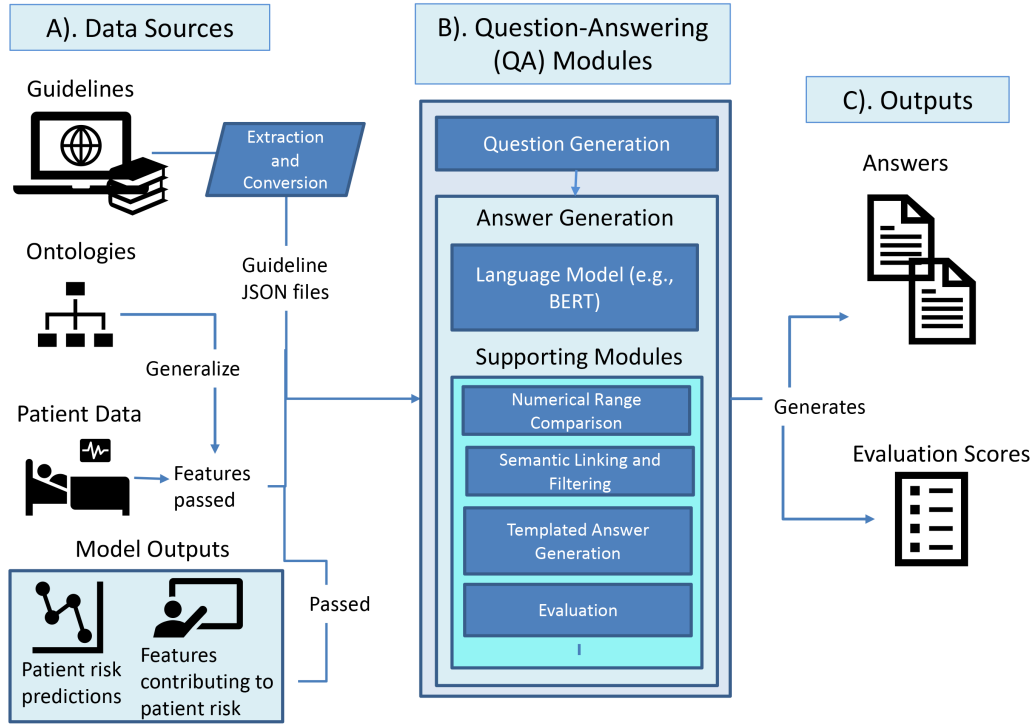
Figure 4: Visualization of different modules within our Question Answering (QA) pipeline including A). information extraction modules, B). QA modules and submodules and C). output modules.

Table 2: Question types currently supported by our QA module, we also indicate the data sources used to address questions of the type.

|   | Question Type | Contextualized Entity | Domain Knowledge Source |
|---|---|---|---|
| 1 | Patient's T2DM summary | Patient | Patient data |
| 2 | Patient's risk summary | Risk Prediction | Risk Prediction and population data |
| 3 | Features contributing to patient's CKD risk | Post-hoc Explanation | Feature importances and ADA guidelines |
| 4 | Patient's medication list | Patient and Risk Prediction | Patient Data and guidelines |
| 5 | Patient's lab values | Patient | Patient Data and guidelines |

13

medical ontologies like Clinical Classification Software (CCS) codes [3] and medical guidelines from ADA Standards of Care 2021 (as introduced in Sec. 2). We query patient data from Limited Claims Explorys Dataset (LCED) claim records (see Sec. 2.1) on-demand, either when we need to create questions based on patient parameters or when we need to include these patient values in answers to questions about the patient. We extract content from the HTML or web version of the 'Standards of Medical Care in Diabetes' [34] guidelines, published by the American Diabetes Association (ADA) [4] using a Python library, BeautifulSoup [35]. We also query patient risk predictions and feature importances using a unique identifier, the patient ID. Some of the extracted contexts are used for generating questions such as patient data, risk predictions and feature importances, and others are used to query against such as the extracted guidelines.

Question Answering Steps: Here we describe part B). of our QA architecture (Fig. 4), including the question and answer generation modules and their supporting submodules. In our QA setup we leverage SOTA LLMs and introduce knowledge augmentations to improve their performance on the ADA 2021 medical guidelines. Additionally, we introduce sub-modules to enhance the LLMs' capabilities to address question types 3 - 5 from Tab. 2, i.e., diagnosis codes, drugs and clinical indicators, which are run against our extracted guideline content. Below are the submodules in our QA setup:

*Question Generation:* The *question generation module* almost always creates templated questions using Python's native support for String Templates, [5] and does so based on patient data, more specifically from the patient's diagnoses codes, lab values, and medication list. We also support the creation of two standard, non-variant questions for each patient, i.e., whose values don't change from patient data, that can help clinicians easily interpret their predicted risk (question type 1) and their T2DM state (question type 2). Moreover, as can be seen from Tab. 2, each of the question types that we support on a per patient basis is populated from different data sources. Hence, we have developed different answering methods for each, including simple lookups and knowledge augmented language model capabilities, including combinations of either a LM + value range comparison or LM +

---

[3]https://www.hcup-us.ahrq.gov/toolssoftware/ccs10/ccs10.jsp
[4]https://care.diabetesjournals.org/content/44/Supplement_1
[5]https://docs.python.org/3/library/string.html

semantic filtering. We provide examples of questions and answers for each question type in Appendix C.

*Answer Generation:* For questions types 1 and 2 from our supported question types, whose context does not depend on guidelines as shown in Tab. 2, we query patient data and feature importances and use a similar templating approach for question generation to populate answer templates with the retrieved query results. For other question types 3 - 5 that are answered by guideline content, we pass them through our LLM and knowledge-augmented LLM setup that we describe next.

*Language Models for Generating Answers:* We use a LLM approach in order to find answers to our questions with the unstructured and natural language discussion and recommendation sentences of the ADA 2021 guidelines. We have applied the original Bidirectional Encoder Representations from Transformers (BERT) model or BERT [36] and other variants of the same retrained on clinical datasets, including SciBERT [37], BioBERT [38], BioBERT-ASQ [39] and BioClincalBERT-ADR [40]. All of the models we utilize are available on the HuggingFace [41] model repository, and we choose BERT models that were made available specifically for clinical question-answering. We built two other submodules to enhance the capabilities of the LLM approach, specifically to address questions with numerical comparisons of question type 5, and to improve the semantic match between the question and the answers returned by a LLM (question type 3 and 4). For details on standard processing steps to include more data types like numerical ranges, refer to Sec. Appendix A in the appendix.

*Augmenting Knowledge to LLM:* Transformer based LLM approaches like BERT and its variants, work on sequences of words that are often seen together and their surrounding words, but don't leverage the semantics of whether these words are diseases, medications, or biological processes. We found that in the absence of this semantic knowledge, we would often get answers from BERT that don't correlate on a semantic level with the question. To eliminate such answers, we explored options for a biomedical semantic mapper and zeroed in on the National Library of Medicine (NLLM)'s Metamap tool [42]. We choose Metamap because of its extensive coverage of biomedical semantic types and its ability to capture entity mentions within the ADA 2021 CPG. Within our pipeline, we have integrated a Python wrap-

per for Metamap[6] that can recognize biological entities within the guideline text and their semantic types (e.g., dsyn: disease or syndrome, phsu: pharmacolgic substance, etc. for a complete list of types returned by Metamap see: [7]). Additionally, given this ability to filter based on semantic types, we want to allow additional answers with mentions of related diseases. To provide more broad answers, we use the UMLS Concept Unique Identifier (CUI) codes from the Metamap returned outputs to map to Snomed-CT disease codes [43]. From the mapped Snomed-CT disease codes, we can traverse the Snomed-CT disease tree to identify how many hops apart question and answer disease codes are and if the answer codes are an ancestor of those in the question. We operate on the idea that answers about the parent disease code apply to children nodes. We use the outputs of these knowledge augmentation modules to both pre-filter and post-sort LLM model answers for question types 3 and 4 from Tab. 2. We report the accuracies for answers that use these knowledge augmentations in the results section (Sec. 4). For more details on how we used the Metamap and Snomed codes as knowledge augmentation methods within our QA setup, refer to Appendix A.

Below in Tab. 3 and 4, we present sample questions and answers for each question type to provide examples of questions and extracted answers supported by our QA approach. We intentionally don't show patient values in these examples to be compliant with HIPAA restrictions.

### 3.4. Prototype Dashboard and Expert panel sessions

To present the supported contextual explanations, we have adapted a question-driven design [12] for user-interface (UI) development and built a running prototype of a risk prediction dashboard (as seen in Fig. 5). The content we show on it is rendered on a per-patient basis chosen from a landing page not shown here. For each patient, we show multiple panes (or UI sections) at a high level, each of which displays content under a particular grouping. These panes include groupings of *patient details, history timeline of claim incidences, risk prediction scores, features contributing to risk, and questions in context*. In Fig. 5, we highlight the risk prediction, feature importance, and questions in context panes. The explanations pane serve as a section where our contextual explanations, that provide context around our

---

[6]PyMetamap: `https://github.com/AnthonyMRios/pymetamap`

[7]`https://lhncbc.nLLM.nih.gov/ii/tools/MetaMap/Docs/SemanticTypes_2018AB.txt`

Table 3: Sample questions and answers for each question type supported by our question-answering approach. Answers such as these serve as contextual explanations that provide information to interpret risk predictions better. We don't provide patient values here due to HIPAA restrictions.

| Question Type | Sample Question | Answer |
| --- | --- | --- |
| 1. Patient's `T2DM` summary | What is the patient's A1C value? What are their most frequent diagnoses codes? | Patient's A1C is A. Their most frequent diagnosis codes are essential hypertension, septicemia, etc. |
| 2. Patient's risk summary | How does the predicted risk of the patient compare against the population? | The predicted risk of chronic kidney disease the patient is X %. The population averages for the same condition are as follows: For Medicare patients: Y % For patients with Charlson Comorbidity Index (CCI) score of 3 : Z % |
| 3. Features contributing to patient's `CKD` risk | What can be done for Essential Hypertension? | 10.3 For patients with diabetes and hypertension, blood pressure targets should be individualized through a shared decision-making process that addresses cardiovascular risk, potential adverse effects of antihypertensive medications, and patient preferences. C |

Table 4: Sample questions and answers for each question type supported by our question-answering approach. Answers such as these serve as contextual explanations that provide information to interpret risk predictions better. We don't provide patient values here due to HIPAA restrictions.

| Question Type | Sample Question | Answer |
| --- | --- | --- |
| 4. Patient's lab values | What should be done for this patient, whose A1C levels are greater than 10 ? | The early introduction of insulin should be considered if there is evidence of ongoing catabolism (weight loss), if symptoms of hyperglycemia are present, or when A1C levels are greater than 10% [86 mmol/mol] or blood glucose levels greater than or equal to 300 mg/dL [16.7 mmol/L] are very high. |
| 5. Patient's medication list | What do the guidelines state about the GLP-1 RA drug the patient is taking? | Meta-analyses of the trials reported to date suggest that GLP-1 receptor agonists and SGLT2 inhibitors reduce risk of atherosclerotic major adverse cardiovascular events to a comparable degree in patients with type 2 diabetes and established AS-CVD (185). |

**The Full Risk Prediction Dashboard**

**365 day risk of CKD among T2DM patients: 66.35% High Risk**

Predicted Risk Bar

**Factors that Influenced the Risk Prediction**

Age at onset 45 - 54
HbA1C
Abdominal Hernia
Calculus of Urinary Tract
Retinal detachments,defects
Other skin diseases
Hemorhoids
Spondylosis, invertebral disc
Essential Hypertension
Other circulatory disorders

**Explanations**

What can be done for patients with Essential hyper...

**Question:** What can be done for patients with Essential hypertension feature, a Disease of the circulatory system?

The risk of atherosclerotic cardiovascular disease and heart failure (see section 10 " cardiovascular disease and risk management," [1]), chronic kidney disease staging (see section 11 " microvascular complications and foot care," [2]), presence of retinopathy, and risk of treatment-associated hypoglycemia (table 4. 3)should be used to individualize targets for glycemia (see section 6 " glycemic targets," [3]), blood pressure, and lipids and to select specific glucose-lowering medication (see section 9 " pharmacologic approaches to glycemic treatment," [4]), antihypertension medication, and statin treatment intensity.[sep]

Confidence Score

Source: ADA Standards of Medical Care Guidelines: Comprehensive Medical Evaluation and Assessment of Comorbidities

References:
[1]: https://doi.org/10.2337/dc21-s010
[2]: https://doi.org/10.2337/dc21-s011
[3]: https://doi.org/10.2337/dc21-s006
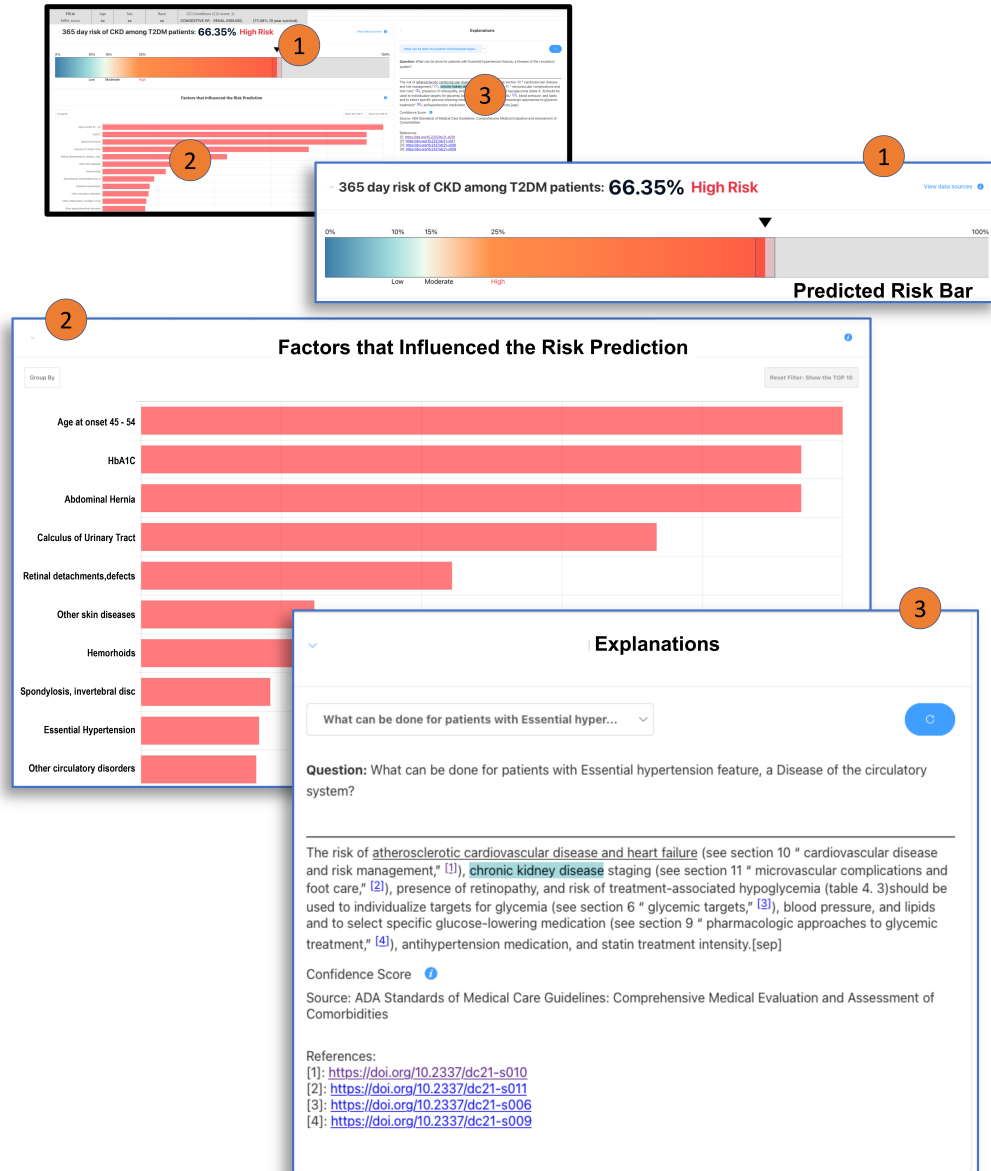[4]: https://doi.org/10.2337/dc21-s009

Figure 5: A screenshot of a running prototype of our risk prediction dashboard which includes: 1) the risk prediction score, 2) the features contributing to the predicted risk with the size of their impact on the model results, and 3) a "questions in context pane", in which the user can select and see answers to questions that provide additional contextual information about the patient, the predicted risk calculation, and individual features contributing to the risk.

19

identified entities of interest in the risk prediction setting - patients, their predicted risk, and the features contributing to risk - can be selected and browsed. Additionally, as we have described in Sec. 2, risk scores can be interpreted better in the context of use, i.e., by enabling connections to patient data, feature importance, and domain knowledge, hence, we had to support interactions between these panes (refer Appendix D for the details), which would make it easier for clinicians to establish the connections.

### 3.4.1. Expert Panel Sessions using Prototype Dashboard as an Aid

We used our risk prediction dashboard as an aid during our structured feedback sessions, where we walked clinicians through a live demonstration of our dashboard for a set of prototypical patients (see Tab. 5). We conducted sessions individually with four clinicians in our expert panel to understand whether the contextual explanations provided, patient predicted risk, and risk explanations, were helpful for clinical practice. We *explained that this dashboard would be available in addition to the clinician's regular EHR tools and the patient information they provide*, and is meant specifically to provide additional information related to the CKD Risk Prediction. To strike a balance between limited clinician time and the need for diverse feedback, we generated such reports from among 20 prototypical CKD high-risk patients from our T2DM cohort, identified by the Protodash algorithm [33].

During the sessions, we first familiarized each clinician with the different sections of the risk-prediction dashboard. We asked them to imagine that they would be meeting with the patient and had seen the CKD prediction, and stated we wanted to understand what information would be useful to them in understanding the prediction and its impact on their treatment decisions. We then presented the dashboard as it would appear for 3 randomly selected prototypical patients. We asked the panel members to imagine that they were preparing to treat a patient that was new to them. We navigated through the dashboard as instructed by the subjects, opening sections or clicking on items as they requested. We asked the clinicians to speak aloud as they were working with the dashboard. We also probed the relevance and usefulness of the different sections of the dashboard and the specific content shown in them. We asked if there was other information they would have liked to have been provided, or questions they would want answered. Sessions were recorded and transcribed, similar to the approach mentioned in [44].

Through these sessions, we wanted to understand the usefulness of our

supported patient contextualizations and the features contributing to their risk. Specifically, we showed clinicians the content on different panes of this dashboard and the supported interactions to understand what features were most important, both from a UI and informational perspective. We report the results of these interactions in Sec. 4.3.

## 4. Results and Evaluation Study

In this section, we present quantitative results for the guideline question-answering methods (Sec. 4.1 and Sec. 4.2). As a qualitative analysis, we also discuss themes and subthemes that we found from an analysis of discussions from our expert panel sessions (Sec. 4.3). For our risk prediction model we choose MLP, and derive important features for the prediction using the SHAP model. The distribution of important features identified by SHAP can be seen in Tab. 5. Results for these models can be browsed via the appendix Appendix C.

### 4.1. Data coverage and support



Figure 6: Overview of the evidence structure in the ADA Standards of Medical Care - Diabetes Guidelines 2021.

Table 5: Summary (generated using Tableone library [45]) of 20 prototypical patients highlighting the demographic and diagnoses counts. We report the disease diagnoses by their higher-level disease groupings (e.g. for T2DM the higher-level code is endocrine, nutritional and metabolic disorders). We highlight the conditions that are most prevalent amongst the patients ($> 50\%$).

| Feature | Overall counts (%) |
|---|---|
| Age at onset 45-54 | 4 (20.0) |
| **Age at onset $\geq 55$** | **15 (75.0)** |
| Age at onset $\leq 44$ | 1 (5.0) |
| SEX - FEMALE | 7 (35.0) |
| Mood disorders | 3 (15.0) |
| Diseases of the blood and blood-forming organs | 3 (15.0) |
| **Diseases of the circulatory system** | **17 (85.0)** |
| Diseases of the digestive system | 6 (30.0) |
| Diseases of the genitourinary system | 9 (45.0) |
| **Diseases of the musculoskeletal system and connective tissue** | **12 (60.0)** |
| Diseases of the nervous system and sense organs | 9 (45.0) |
| **Diseases of the respiratory system** | **11 (55.0)** |
| Diseases of the skin and subcutaneous tissue | 7 (35.0) |
| **Endocrine; nutritional; and metabolic diseases and immunity disorders** | **20 (100.0)** |
| Infectious and parasitic diseases | 10 (50.0) |
| Injury and poisoning | 4 (20.0) |
| Mental Illness | 3 (15.0) |
| Neoplasms | 6 (30.0) |
| Symptoms; signs; and ill-defined conditions and factors influencing health status | 10 (50.0) |

One of the aims of this manuscript is to extract contextual explanations from medical guidelines. To explore the feasibility of this task we first analyzed the coverage of the ADA 2021 guidelines used to extract contextual explanations of predictions in CKD comorbodity, T2DM risk prediction setup.

We extracted the recommendations and discussion sentences across the 16 chapters of the current ADA 2021 CPGs. These recommendation and discussion sentences are expressed in natural language (See Fig. 6). Table 6 shows a high-level overview of the coverage statistics. Specifically, the extracted sentence corpus can hence be analyzed for the total number of tokens, average token length per sentence, and their composition of Metamap semantic types, to understand the coverage of the guideline text in terms of volume and semantic diversity. For tokens, we report words recognized by both BERT's tokenizer model to be consistent with our QA approach. [46] report similar statistics for three other CPGs, neither of which are Diabetes focused, but it can be seen that the total number of tokens and sentences in the ADA 2021 CPG are more than the three guidelines reported in this paper. Hence, pointing to the comprehensiveness of our approach. Also, note that some of the recommendations were not captured by our guideline extraction script, and hence our statistics might be lesser than the actual count. [8].

Table 6: Coverage statistics from extracted content from the ADA Standards of Care - Diabetes Guidelines 2021. We report these statistics on the recommendations and discussion sentences we extracted across chapters.

| Field | Count |
|---|---|
| Chapters | 16 |
| No. of sentences | 2379 |
| Tokens from BERT | 118350 |
| Avg. BERT tokens per sentence | 49 |
| Metamap Semantic types covered | 116 / 126 |

The many semantic types (see Fig. 7 for 25 most populous semantic types) covered by the ADA 2021 CPG reaffirms that guidelines are a comprehensive source of evidence-based information in the clinical domain [28].

---

[8]In future, we aim to expand our coverage and update our methods to better capture these recommendation groups
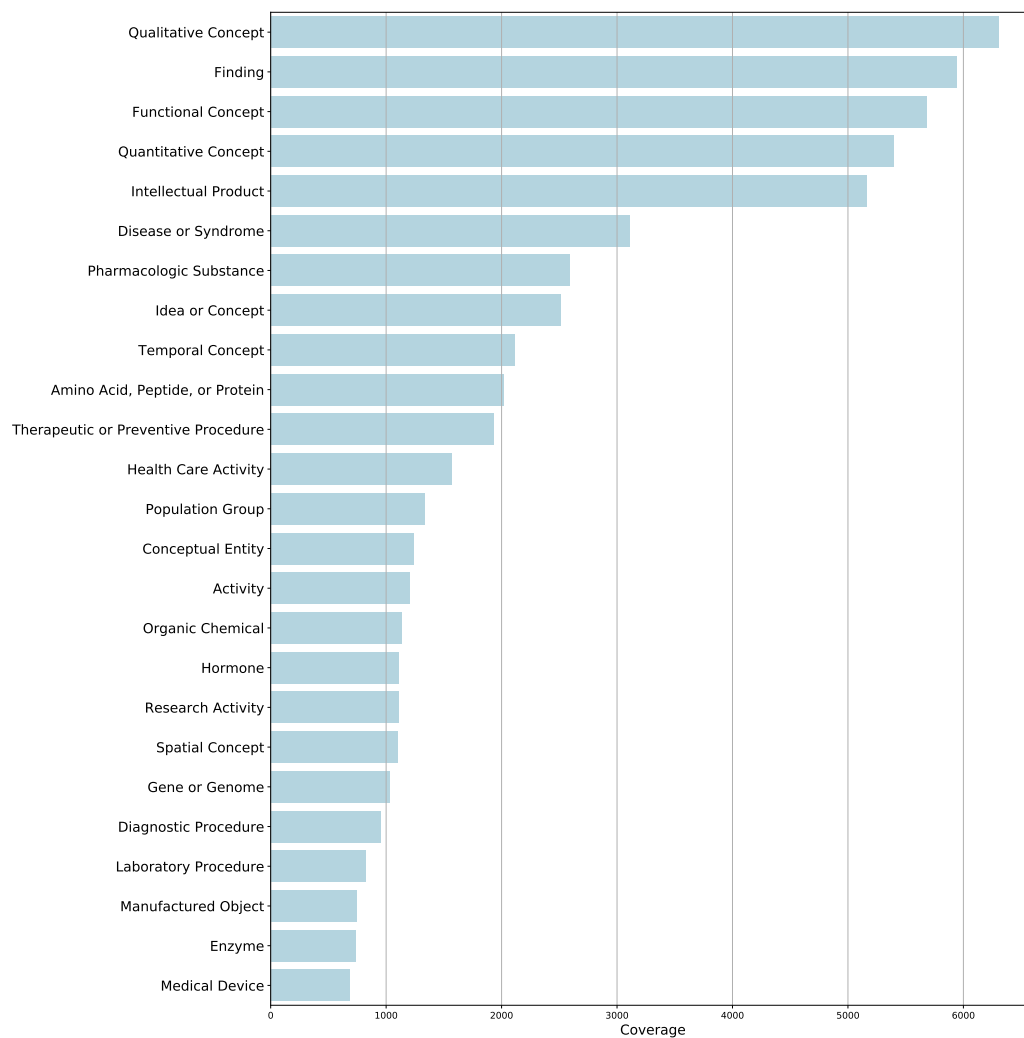
Figure 7: Frequency distribution for 25 / 116 of the top semantic types that were found in the extracted guideline text.

Further, in Sec. 5, we also discuss how well the ADA 2021 CPG alone can support the themes we analyzed from conversations with clinicians during our structured feedback sessions and hence discuss the CPG's ability to serve as a source of context in our risk prediction setting.

## 4.2. Quantitative Evaluation of Guideline QA

One of our key aim is to study whether SOTA LLM methods can be used to extract high quality contextual explanations. Thus, while we can address different question types in our QA approach, as seen in Tab. 2, we only evaluate those question types that are addressed by ML methods, i.e., through our knowledge and knowledge augmented LLM modules (as described in Sec. 3 and seen in Fig. 4). We evaluate feature importance questions of types 3 (of diagnostic importance), type 4 (of treatment importance), and type 5 (asking about clinical indicators and important for both diagnostic and management purposes). Additionally, we evaluate answers to questions of types 3 and 4 differently from question type 5 since questions of types 3 and 4 are served by the knowledge augmented LLM modules and questions of type 5 are addressed by the LLM with the numerical range comparison module. For answers to questions of type 3 and 4, we report standard and frequently used NLP QA metrics including mean average precision (MAP) [47], F1-score their contributors of precision and recall and BLEU scores. For answers to questions of type 5, we report the number of times the combination of the LLM and numerical range comparison module could correctly predict whether the answer outputted was in / out range for the numerical value being asked about in the question.

### 4.2.1. Evaluation Results of Feature Importance Questions of Diagnostic Importance

We first evaluated the quality of extracted contextual explanations for feature importance questions, type 3. As outlined in Sec. 3.3, our aim was to evaluate the readiness of SOTA LLM models for this task. Here we report the results for 71 questions covering relevant feature importances for patients' risk predictions and these questions cover 14 CCS LVL 1 diagnosis code types. As mentioned previously, we report the performance in terms of a number of standard metrics. However, given our information extraction setting, among these we are especially interested in the precision metrics. These metrics measure how many documents, among the ones retrieved, were relevant. Furthermore, 'MAP' and 'precision@k' measures the

same while considering the order of retrieval. This aligns closely with how clinicians evaluate presented informations where the presented answers are expected to accurate in order of most acceptable to least. It is to be noted, that we evaluated the predicted results to disease feature importance questions against candidate answers by manually inspecting the ADA 2021 CPG. The annotations were done by an author and some of these annotations were verified by a clinical expert on the team who is also a co-author in this paper. We report results on the expert validated subset in Appendix C. Table 7 reports results in comparison to the entire annotated dataset of 85 feature questions and 654 candidate answers for the out-of-the shelf LLM models under consideration. The results show that in terms of 'MAP' as well as 'precision at k' metrics, vanilla BERT outperforms the other LLM models. SciBERT is a close second with an improved recall and F1 score. We analyze the importance of these results further in Sec. 5.

Additionally, we also evaluated the results by augmenting the base LLM models with different strategies. Tab. 8 reports the result scores for answers to question type 3 using the best knowledge augmentation strategy (across the five different settings) for each LLM models. Overall, we can see a significant improvement in 'MAP' and 'precision at 5' for knowledge augmented BERT model (BERT-KA) over BERT, the best performing base language model. In terms of other metrics, knowledge augmented SciBERT-KA shows a consistent improvement for all metrics while being best/second-best overall.

While these evaluation numbers are reported from a small evaluation set of 71 questions and 654 candidate actual answers, we consider this as a somewhat comprehensive evaluation due to the diversity of diseases covered (1844 diseases) and in total semantic types covered within the answers (116 semantic types, see Tab. 6), both in the candidate and predicted sets. Overall, from these results, we see that knowledge augmentation can improve the base language model performance.

### 4.2.2. Evaluation Results for Drug Questions

To demonstrate the flexibility of the LLM setup on various settings that match the coverage of the T2D guideline data (see Fig. 7) as well as to test the generalizability of our results. We also report results for 6 anti-diabetic drug questions of question type 4. Table 9 show the results for out-of-the shelf language models as well as the knowledge augmented language models for the metrics of interest. Overall, we once again found the knowledge-augmented language models to be the best performing ones. In terms of 'MAP' and

Table 7: Performance of Guideline QA on different language model approaches reported at mean average precision (MAP), F1 and recall at top-10 answers and precision at top-1 and top-5 for 71 disease feature importance questions. The models are sorted by MAP values, to indicate an ordering of the best models.

| model | bleu | P@1 | P@5 | map | f1 | recall |
|---|---|---|---|---|---|---|
| BERT | 0.117 | 0.468 | 0.382 | 0.390 | 0.213 | 0.241 |
| BioBERT | 0.116 | 0.431 | 0.339 | 0.346 | 0.200 | 0.238 |
| BioBERT-BioASQ | 0.132 | 0.383 | 0.329 | 0.332 | 0.217 | 0.281 |
| BioClinicalBERT-ADR | 0.125 | 0.368 | 0.317 | 0.316 | 0.205 | 0.259 |
| SciBERT | 0.165 | 0.461 | 0.349 | 0.364 | 0.261 | 0.354 |

Table 8: Performance of Guideline QA of knowledge augmented language models reported at mean average precision (MAP), F1 and recall at top-10 answers and precision at top-1 and top-5 for 71 disease feature importance questions. Here we show the best knowledge augmentation approach per model to indicate highest gains over baseline performance for the native language model approaches we tried. Best and second-best values for each column is highlighted in Green and Blue color, respectively. Language model (e.g. BERT) suffixed with KA represents the corresponding knowledge augmented model (e.g. BERT-KA).

| model | bleu | P@1 | P@5 | map | f1 | recall |
|---|---|---|---|---|---|---|
| BERT-KA | 0.075 | 0.467 | 0.419 | 0.438 | 0.169 | 0.186 |
| BioBERT-KA | 0.127 | 0.434 | 0.348 | 0.353 | 0.215 | 0.254 |
| BioBERT-BioASQ-KA | 0.141 | 0.458 | 0.362 | 0.369 | 0.237 | 0.280 |
| BioClinicalBERT-ADR-KA | 0.121 | 0.406 | 0.321 | 0.330 | 0.202 | 0.242 |
| SciBERT-KA | 0.192 | 0.473 | 0.341 | 0.375 | 0.291 | 0.405 |

'precision at 5', BERT-KA comes out as the best performing model, where as SCIBERT-KA comes out as either the best or the second-best model for all metrics. It is to be noted that the overall results are significantly better than the ones for disease questions. One possible reason behind this effect may be related to the fact that drugs are referred directly in guidelines and thus QA models are able to pick these sentences with greater efficacy.

Table 9: Performance of Guideline QA with different knowledge augmentations of language model approaches reported at mean average precision (MAP), F1 and recall at top-10 answers and precision at top-1 and top-5 for 6 anti-diabetic drug feature questions. Best and second-best values for each column is highlighted in Green and Blue color, respectively. Language model (e.g. BERT) suffixed with KA represents the corresponding knowledge augmented model (e.g. BERT-KA).

| model | bleu | P@1 | P@5 | map | f1 | recall |
|---|---|---|---|---|---|---|
| BERT | 0.100 | 0.910 | 0.751 | 0.757 | 0.254 | 0.206 |
| BioBERT | 0.100 | 0.726 | 0.643 | 0.635 | 0.231 | 0.192 |
| BioBERT-BioASQ | 0.081 | 0.708 | 0.694 | 0.704 | 0.222 | 0.162 |
| BioClinicalBERT-ADR | 0.075 | 0.593 | 0.614 | 0.597 | 0.192 | 0.146 |
| SciBERT | 0.121 | 0.947 | 0.757 | 0.772 | 0.281 | 0.228 |
| BERT-KA | 0.099 | 0.900 | 0.863 | 0.821 | 0.281 | 0.213 |
| BioBERT-KA | 0.083 | 0.802 | 0.704 | 0.720 | 0.234 | 0.170 |
| BioBERT-BioASQ-KA | 0.117 | 0.711 | 0.725 | 0.716 | 0.272 | 0.221 |
| BioClinicalBERT-ADR-KA | 0.085 | 0.598 | 0.595 | 0.587 | 0.199 | 0.152 |
| SciBERT-KA | 0.128 | 0.912 | 0.823 | 0.794 | 0.298 | 0.232 |

### 4.2.3. Evaluation Results of Clinical Indicator Questions

In Tab. 10, we report the accuracy statistics for the performance of our rule augmentation / numerical range comparison module of our QA approach. We show granular breakdowns based on different numerical operators, greater than, lesser than, and equal to, to indicate which settings are being picked up the best by the rule augmentation and which others are harder. In this evaluation, we manually went through the outputted answers to ensure that they were within range of the numerical values in questions. The reason for this annotation approach is that the guidelines have few sentences for actions to be taken on clinical indicators. Hence, there is not much diversity in the answers that a LLM like BERT can output before passing the answer

Table 10: Results of Guideline QA with rule augmentation of language model approaches for numerical comparisons reported for 9 questions across the 20 prototypical patients identified from our predicted high-risk chronic kidney disease cohort. The split of question variations is equal across the different numerical range comparison operators of lesser than, equal to and greater than.

| Comparison | Accuracy | TP | TN | FP | FN | Total |
|---|---|---|---|---|---|---|
| Overall | 0.78 | 7 | 7 | 3 | 0 | 18 |
| Lesser Than | 0.84 | 2 | 3 | 1 | 0 | 6 |
| Equal To | 0.67 | 1 | 3 | 2 | 0 | 6 |
| Greater Than | 100 | 4 | 2 | 0 | 0 | 6 |

TP - True Positives, TN - True Negatives, FP - False Positives, FN - False Negatives. Accuracy computed as accuracy = (TP + TN)/ Total

to our range comparison module for validation. These accuracies highlight the value in combining syntactic parsing output with a LLM to improve its capabilities.

Overall, these quantitative evaluations points to the potential in applying scalable, augmented LLM-based approaches to extract content from authoritative guideline literature that can then be used to provide context to interpret model predictions, such as in our setting, risk prediction scores and their model explanations.

*4.3. Qualitative Evaluation with Clinicians*

We conducted a thematic analyses on the responses and feedback received during the expert panel sessions (Sec. 3.4), as follows. Three independent researchers, who are coauthors on this paper, reviewed the transcripts, flagged significant utterances, and characterized these utterances in terms of the major points and themes they expressed. The researchers then reviewed their sets of identified themes and utterances together, and grouped and combined them into a single agreed-upon set of overarching themes. We report this combined list of themes below (Tab. 11, 12, 13, 14). These themes reflect areas that clinicians prioritize and where the support of explanation-driven, AI risk prediction tools would be appreciated.

As seen in Tab. 11, 12, 13 and 14, we have grouped the discussions from the expert panel sessions into *four high-level themes* spanning different areas where clinicians would benefit the most in a chronic disease, comorbidity

risk-prediction setting such as ours. The high-level themes we found include 'Theme 1: Clinical Value of Explanations and Contextualizations', 'Theme 2: Highlighting Actionability', 'Theme 3: Connections to Patient Data' and 'Theme 4: Connections to External Knowledge'. We were further able to create sub-themes for more granular topics that came up during the discussions under each of these themes, bringing the theme and sub-theme total to *four high-level themes and twelve sub-themes*.

Table 11: *Clinical Value of Explanations and Contextualizations* - 1st theme that emerged during our expert panel interviews with clinicians where we walked through the risk prediction dashboard and the contextual explanations that we support. We attach a description for each sub-theme that we found and we also provide examples in quotes.

| Sub-theme | Description |
| --- | --- |
| Value of Contextual Information around CKD risk | All clinicians saw value in connecting the T2DM patient's CKD risks to *data on their other conditions*, and to, *relevant recommendations from the T2DM guidelines*. For example, some clinicians reasoned about "how the patient's CKD risk changes their dosage / treatment" and some others were interested about "connections to other conditions that patient has." |
| Value of Contextual Information around Individual Features | Clinicians found that information from T2DM guidelines and cited literature relevant to factors that contributed to the system's predicted CKD risk were helpful to understand how the *factors could be related to CKD or T2DM*, and how they might interact with *other factors shown*. For example, "how does a skull fracture elevate CKD risk?" This was particularly valuable when not previously known by the clinician, for example: "it is surprising, and I have learned something about celiac disease and abdominal pain connection" in patients with diabetes. |
| Value of Contextual Information around patient's T2DM | Besides the patient's CKD risk and its implications, the clinicians were interested to know about the *patient's T2DM state, their comorbid conditions and other parameters in relation to their T2DM diagnosis*. For example, the clinicians wanted to know "how long has the patient had their T2DM" or "what is their A1C progression?" |

More specifically, under 'Theme 1: Clinical Value of Explanations and Contexts', we group instances where clinicians could make sense of the risk predictions and post-hoc explanations by the additional context provided, or instances where clinicians would appreciate more context. Within 'Theme 2: Highlight Actionability,' we discuss instances where clinicians mentioned a need to depict actionable features and indicate actions for them concerning the patient's T2DM diagnoses or their elevated CKD risk. Under 'Theme 3: Connections to Patient Data,' we cover instances where clinicians looked for connections to patient history or their lab results while reasoning about the patient case. Finally, under 'Theme 4: Connections to External Knowledge,'

Table 12: *Highlighting Actionability* - 2ⁿᵈ theme that emerged during our expert panel interviews with clinicians, where we walked through the risk prediction dashboard and the contextual explanations that we support. We attach a description for each sub-theme that we found and we also provide examples in quotes.

| Sub-theme | Description |
|---|---|
| Highlight Actionable and Modifiable Factors | Most of the clinicians were interested in highlighting patient risk factors that could be controlled or acted upon, vs. those (such as age) that could not be influenced: "what factors can be changed?" |
| Highlight the Impact of CKD risk prediction on Treatment Decisions for Diabetes and other conditions | When shown information from treatment guidelines, the clinicians wanted to understand how they *reflect the patient's CKD risk and T2D diagnosis*:"are any of these proposed medications contraindicated?" |
| Suggest Specific Actions to Reduce CKD risk | Clinicians wanted to understand ways to reduce the CKD risk including ways of addressing risk factors and changes to medications: "do any of the patient's current medications increase risk of renal toxicity?' |

Table 13: *Connection to Patient Data* - 3ʳᵈ theme that emerged during our expert panel interviews with clinicians, where we walked through the risk prediction dashboard and the contextual explanations we support. We attach a description for each sub-theme that we found and we also provide examples in quotes.

| Sub-theme | Description |
|---|---|
| Connections to Patient's Clinical Indicators | The clinicians indicated that they want to see clinical indicators for diagnoses ( if available ), when interpreting the *factors that led to the risk* or the *patient's CKD risk score*: "Given the patient has essential hypertension, what was their lab systolic blood pressure reading," or "The patient's eGFR value will be important to show for CKD," or "what do the guidelines say about this patient's systolic and diastolic blood pressure readings." |
| Need for Information on Related Diagnoses | When shown factors which were diagnosis codes that influenced the risk prediction, clinicians wanted to *see what other diagnoses that the patients had, that might align or contribute*: "show COVID-19 answers for lower respiratory disorders?" or "what episodes of abdominal pain did the patient have?" |
| Connections to Patient's History | When shown certain factors, the clinicians wanted to know the when the patient had the diagnosis, if the condition was a current one, and about changes over time. for example, "when did the patient have a genitourinary diagnosis?", or"what does their eGFR progression look like?' |

Table 14: *Connection to External Medical Domain Knowledge* - 4<sup>th</sup> theme that emerged during our expert panel interviews with clinicians, where we walked through the risk prediction dashboard and the contextualization we support. We attach a description for each sub-theme that we found and we also provide examples in quotes.

| Sub-theme | Description |
|---|---|
| Links to Medication Databases | When deciding upon *treatment suggestions* for patients given the *knowledge of their CKD risk and T2DM diagnosis*,the clinicians wanted to understand how their medications interact : "what drugs they are on currently have a bad renal impact?" or "how does their current anti-diabetic drug interact with a CKD drug?" |
| Links to Published Articles | When connections between the *CKD risk prediction* and *the factors contributing to the risk* were unclear, clinicians mentioned they would look for published references: "what is the connection between CKD and respiratory disorders?" or "how does celiac disease mention from the guideline answer, affect CKD?" |
| Support familiar categorizations | Some clinicians were looking for more *provenance around the categorization schemes* we were utilizing to show higher-level physiological pathways for diagnoses codes, and were also hoping for connections to familiar schemes like "ICD-10": "how are hemorrhoids linked to the circulatory system?" |

we describe instances where clinicians mentioned a need to make connections to the latest literature, other medication databases, or other clinical schemes they utilize. In addition, these themes have an order among them, to address the clinical value of explanations from Theme 1 and the actionability aspects from Theme 2, the content requests from Themes 3 and 4 would be contributing features, which are the connections to patient data from Theme 3 and the links to external knowledge from Theme 4. We present deeper breakdowns in the form of sub-themes and descriptions for each of these four themes in Tab. 11, 12, 13 and 14.

As the expert panel sessions were conducted mid-way through our current implementation, some of these themes served as a means for further refinements, such as features to support a 'need for more related diagnoses,' and 'connections to patient's clinical indicators,' To address these themes, we refined different modules of our pipeline, including the user interface, the question-answering, and post-explanation modules. In future, to address the requests under 'Theme 4: Connections to External Knowledge,' we plan to support connections to external resources, which clinicians may find valuable.

In summary, these themes and sub-themes from the expert panel sessions validate our hypothesis about the need for additional clinical context to situate risk predictions and span requests for better connections and pre-

sentations of domain knowledge that clinicians are familiar with in these settings.

## 5. Discussion

In this section we analyze both the quantitative and qualitative results presented in Sec. 4. We analyze both the feasibility of supporting contextual explanations from authoritative sources such as CPGs, and the usefulness of providing contextual explanations from an analysis of the themes derived from our expert panel sessions.

### 5.1. Feasibility of extracting and generating contextual explanations from authoritative sources

For our feasibility analysis, we further analyze the results for type 3 questions with respect to the disease groups to gain a deeper understanding of the QA performance. Specifically, we want to understand whether SOTA LLM backed QA methods, potentially augmented with knowledge, are ready for real-world use for our states use-case as well as identify patterns that might apply to other CPGs in different disease areas. We pose a number of questions around this idea as follows:

Is CPG guideline suitable for T2DM contexts: *How well are the disease subgroups among T2DM patients covered in the guidelines?* Here, we attempt to understand the applicability of ADA 2021 CPG in our use-case. From Fig. 8, we can interpret that the guidelines cover a smaller number of disease groups [9] than the patient data. Since CPGs are authoritative literature in their disease fields, their coverage is mainly limited to the primary disease area. Thus ADA 2021 CPG focuses on Diabetes, an Endocrine, Nutritional and Metabolic Disorder, and its comorbid conditions (mainly spanning diseases of the circulatory and genitourinary systems). Unsurprisingly, these patterns are seen in the Fig 8 as well where the Endocrine, Nutritional and Metabolic Disorder have the largest coverage in the guideline data, a 66%. In contrast, patients might have other conditions that do not arise from the T2DM diagnosis alone, and hence we can deduce that we see more diversity in disease groups in the patient data.

---

[9]Disease groups are derived by rolling up disease codes both in the patient data and guidelines to their higher-level CCS LVL 1 groups.
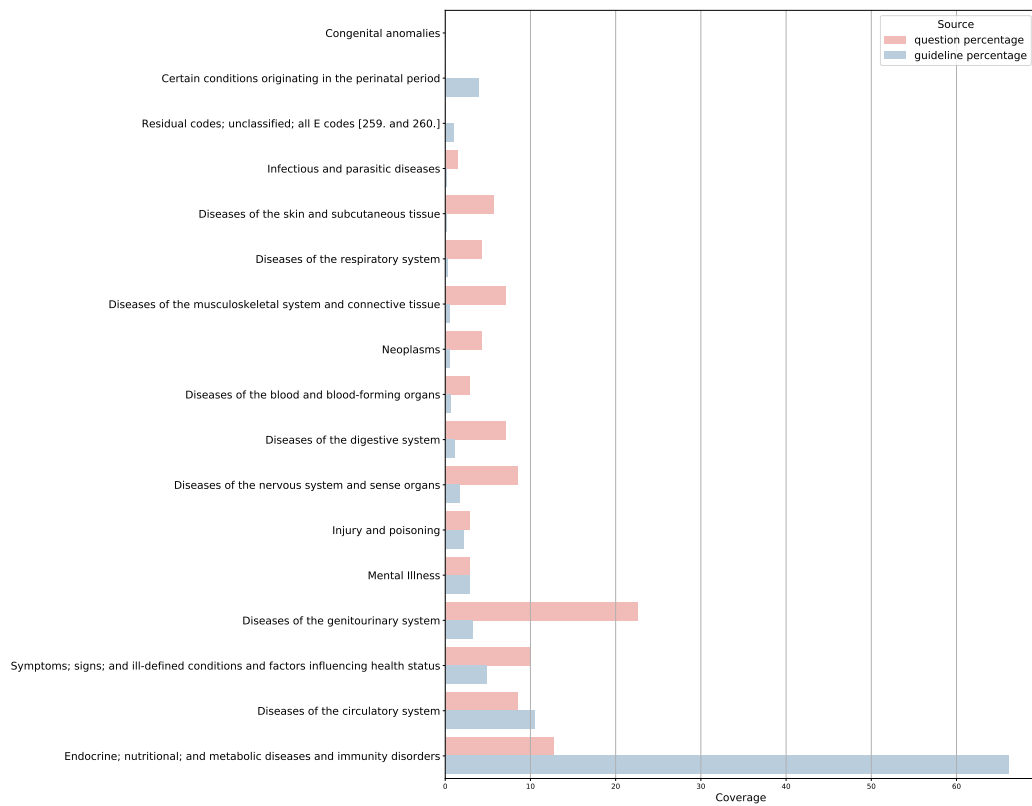
Figure 8: Comparing disease group occurrences in the ADA CPG 2021 versus those in the feature importance questions from our chosen prototypical patients.
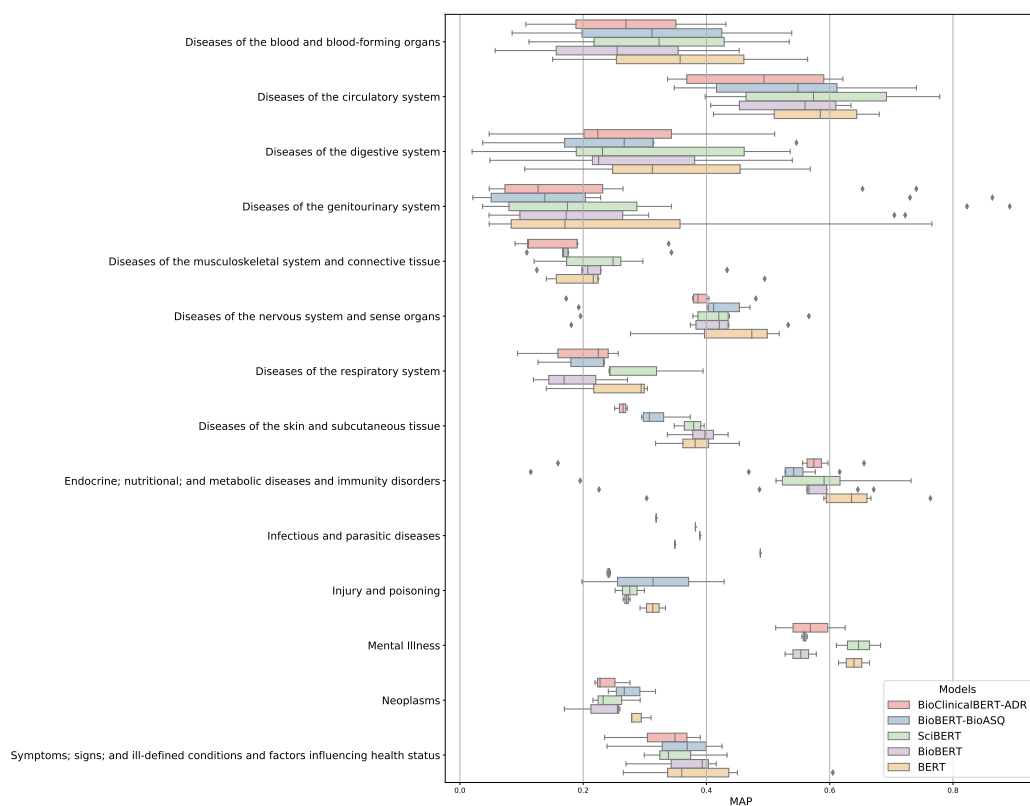
Figure 9: Overall view of the performance of all language models in our experiments over the 14 disease groups covered in feature importance questions.
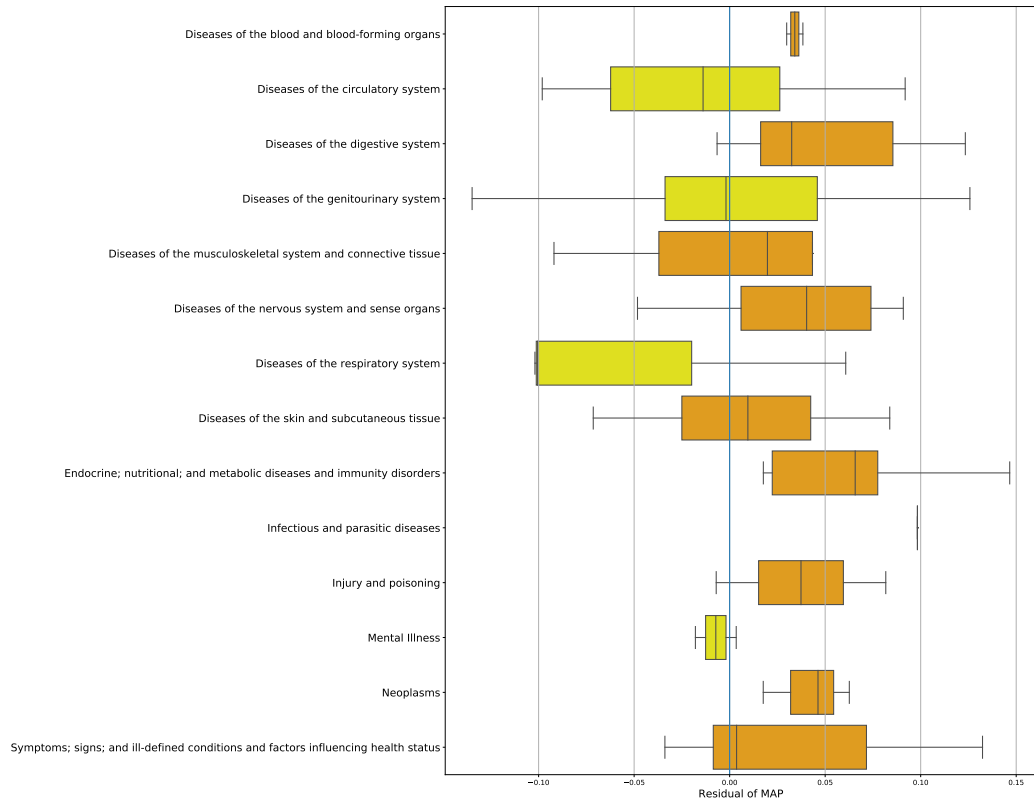
Figure 10: Comparative performances of top-performing native language models, BERT vs. SciBERT. Plotting the residuals under equal performance hypothesis for the 140 feature importance questions that span 14 disease groups. Orange box indicates BERT performs better on average while Yellow indicates SciBERT is the better choice for the disease group.

Can a single SOTA LLM method be used to extract the contexts? We attempt to understand if the LLMs are inherently better at certain disease groups over others. Fig 9 shows the distribution of base LLM models over the disease groups. We can see that there are a few disease groups which have a higher MAP performance than others (towards the right end of the plot), some have their the box centers in the middle of the plot and others who are not doing as well since they are in the first quadrant of the plot. While the results are not strikingly decisive and statistically significant everywhere, in concordance with our overall results, we note that SciBERT and BERT models have better performance over most of the disease groups. Thus to further discern between these top 2 performing models, we conducted a point-wise analysis of relative performance difference between BERT and SciBERT (distribution of residual values between the MAP performances of BERT and Scibert under equal performance hypothesis). Fig. 10 shows the outcomes of the analysis where orange box indicates that BERT performs better on average while yellow indicates the same for SciBERT. We see that BERT is better for most disease groups, especially for 'Disease of the blood and bone forming organs', 'Diseases of the digestive system', 'Diseases of the nervous system and sense organs', 'Endocrine, nutritional, and metabolic diseases and immunity disorders', 'Injury and poisoning', and 'Neoplasms' (0 not contained in the inter-quartile range). SciBERT is only doing better on 'Diseases of the respiratory system' (and marginally better for 'Mental Illness'). These results, in addition to the quantitative results, indicate that LLM models are better at addressing some disease groups than others. While vanilla BERT is a defensible choice, the results point to the need for domain adaptation for LLM for this problem. However, considering the limited availability of data, novel ML methods such as one-short learning and as weak supervised techniques may be required to improve the performances of these LLMs reliably across multiple disease groups.

Does knowledge augmentation reliably improve QA methods? We proposed 4 possible strategies for Knowledge-augmentation (See Appending Appendix A.2). Among these settings, the best knowledge augmentation strategies reflected in Table. 8 originated from a composite of strategies. As seen from Tab. 7, 8 and 9, the guideline QA's *best MAP score of 0.82* is obtained in a BERT + knowledge augmented setting on drug questions and among disease features, the guideline QA's *best MAP score is 0.438*. The *recall with its highest value of 0.405* is obtained in a post-filtering knowledge augmentation
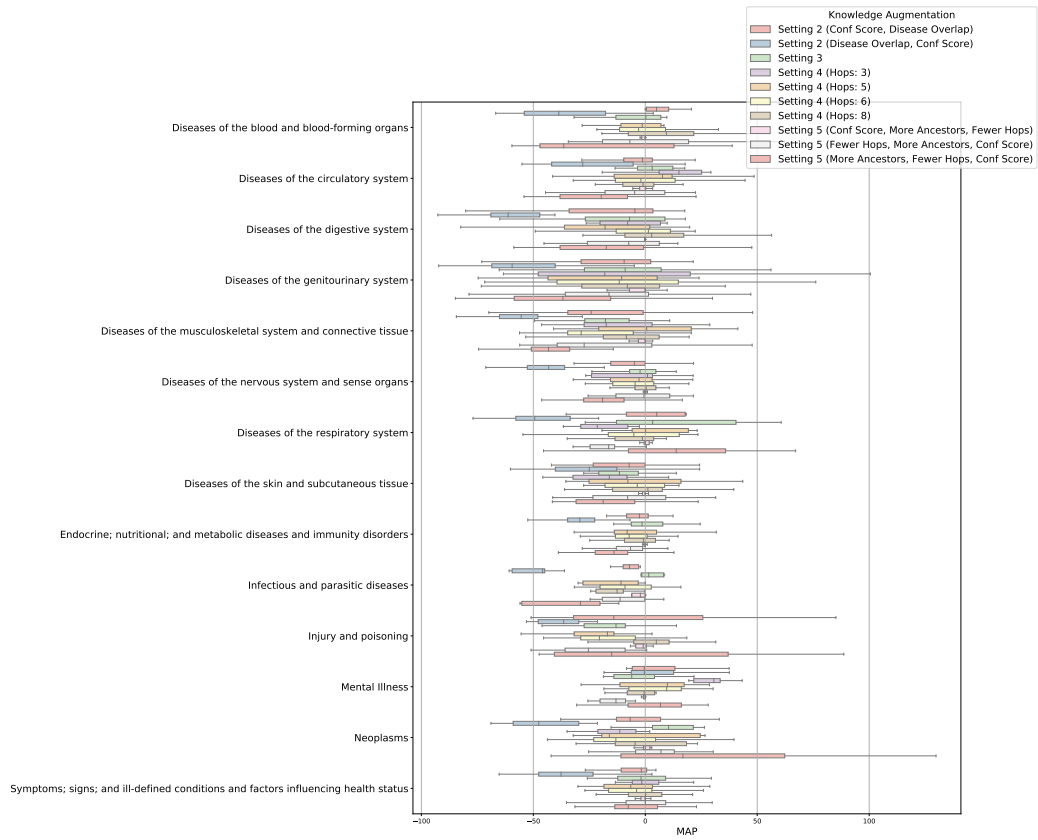
Figure 11: Lift in performance over native language models by different knowledge augmentation strategies over the 14 disease groups covered in feature importance questions.

setting 5 of SciBERT (Tab. 8), where we sort answers by disease overlap between question and answer and we also see the best *BLEU score of 0.19* in this setting. These points to the fact that there is value in either filtering the answers to be passed to a LLM or sorting the answers from it, using aids from known domain knowledge sources that data sources like guidelines are expected to adhere to. We further analyzes these at disease sub-group level in Fig. 11 where we plot the lift in performance over corresponding base LLM using any particular strategy. While any one strategy is not found to be dominating the others, for most disease groups we can find one or more strategy that improves the performance (median lift greater than 0). These results support our previous insight that there is a value in augmenting domain knowledge. However, finding a single universal strategy is difficult and may need further research.

Overall, how feasible it is to extracting contexts from guidelines? Which strategies are beneficial? Are the methods scalable? We have addressed 175 clinically relevant questions that provide context around 20 prototypical patients, their predicted risk, and the factors influencing their risk. We have implemented logical adaptations given what we know about the guideline data to improve the LLM model's capabilities and performance. These adaptations include knowledge augmentation from well-used medical ontologies like Metamap and Snomed to improve semantic overlap and rule augmentation to address numerical range questions. Our baseline LLM, BERT itself, has a variable performance that does well on some questions and not on others. Similarly, the best performances on the LLM + knowledge augmentation approaches varies across pre-filtering settings 3 and 4, that filter by Metamap disease codes and Snomed disease hops and post filtering-setting 5 that sorts by Snomed disease hops. From our result evaluations, we see that the order of introducing the knowledge augmentation outputs impacts the accuracy scores, namely the MAP and recall. Mainly, pre-filtering the answer set before passing to a LLM can help it output more precise answers. In the best case, pre-filtering settings provide a gain of 4% over the baseline LLMs both for disease and drug questions (BERT-KA from Tab. 7 and BERT-KA from Tab. 9). Similarly, post sorting the answers from a LLM can improve the recall, and in the best case (SciBERT-KA from Tab. 8), we see a gain of 5% from the baseline LLM.

Our result numbers also indicate that unsupervised adaptations can only reach a certain accuracy and point to the need for domain adaptations to

medical guidelines. Additionally, since we were dealing with a setting with little or no annotations on the ADA 2021 CPG, we had to create our own annotations. Currently, we are dealing with a relatively small annotated corpora (85 questions and 654 candidate sentences), and we consulted with a medical expert on our team to review these annotations. Even for this small corpora, we find that it is time-consuming for a clinical expert (s) to review the annotations or create them. We are exploring techniques like weak supervision to scale and improve the coverage of the annotations. In summary, our guideline QA results depict incremental gains in adding knowledge and rule augmentations to enhance LLMs' performance and capabilities in domain applications and point to the need for supervised and semi-supervised approaches to improve these gains.

We are, to the best of our knowledge, the first to report any QA performance numbers on the ADA CPG 2021 dataset. Additionally, we are the first few who have tried a LLM approach for more scalable upstream tasks on medical guidelines like question answering than the current more time-consuming and dataset-dependent task of converting guidelines to rules and applying logical reasoning techniques over these rules [48, 49]. Our approach to guideline extraction and question answering (Fig. 4) is a step towards providing a more flexible way ( [46, 50]) to swap in guideline text from different diseases as needed. Our enhanced LLM approach (Fig. 4) can be applied to any medical text corpus like medical guidelines extracted into a machine-comprehensible format and can address different question types (as seen in Tab. 2) relevant in risk prediction settings.

*5.2. Understanding the added benefit of the derived contexts*

What were the takeaways and feedback from clinicians about the supported contextual explanations? The four major themes - *Clinical Value of Explanations and Contextualizations, Highlighting Actionability, Connections to Patient Data, and Connections to External Knowledge Sources* - that we found during the expert panel interviews to evaluate our contextualization approach, mainly point to the overall value of supporting different types of contexts, both from literature and patient data, and the need to better present connections between these contexts. Many of the contexts the clinicians on our expert panel were looking for were around the post-hoc explanations of the factors contributing to the risk. This finding corroborates a recent study that reports that post-hoc explanations themselves are insufficient to provide reasoning that clinicians can interpret and act upon [7], and

also add to the well-accepted belief that risk scores are insufficient.

| Theme (# of subthemes) | Sources | Current Coverage |
|---|---|---|
| 1. Clinical Value of Explanations and Contextualizations (3) | Guidelines, Feature Importances, Medical Ontologies, Patient Data and Published Literature | ✓ 3/ 3 |
| 2. Highlighting Actionability ** (3) | Guidelines, Published Literature | ✓ 1 / 3 |
| 3. Connections to Patient Data (3) | History, Past Diagnoses and Clinical Indicators | ✓ 1 / 3 |
| 4. Connections to External Medical Domain Knowledge (3) | Guidelines, Medical Ontologies Medication Databases, Published Literature and Familiar Categorizations | ✓ 1 / 3 |

*Sources we currently support in our QA approach to provide context are shown in green
** Not currently supported by our methods

Figure 12: Summarizing the coverage of our current data sources to support the themes we found from our expert panel discussions.

Do the supported data sources address the clinicians needs? Through further analysis, we find that the contexts the clinicians were looking for and discussing can be addressed either by connections to patient history and data - patients' diagnoses, medications, and lab values - or through published literature. Specifically, we find that the different questions and question types (Tab. 2) that we support from the T2DM guidelines can address 6 **of the 12 sub-themes** (Fig. 12), i.e., providing contextual information around patient's T2DM state, their CKD risk and the individual features (Theme 1), highlighting the impact of CKD risk on treatment decisions for T2DM (Theme 2), providing links to published articles (Theme 4), and showcasing connections to patient clinical indicators (Theme 3) where mentioned. Some other themes can be easily addressed by enabling connections from the CKD risk scores and the features contributing to them, to patient timelines for diagnoses and lab values. Other themes - support for familiar categorizations (Theme 4) and the need for information on related diagnoses (Theme 3) - benefit from connections to medical ontologies that support either drill-downs to more specific diagnoses or abstracting up to higher-level pathways. We currently only support abstractions to higher level physiological pathways on the prototype dashboard (e.g., all disease of the circulatory system can be filtered from the patient's feature importances) and are investigating how to support drill-downs based off of these pathways more broadly.
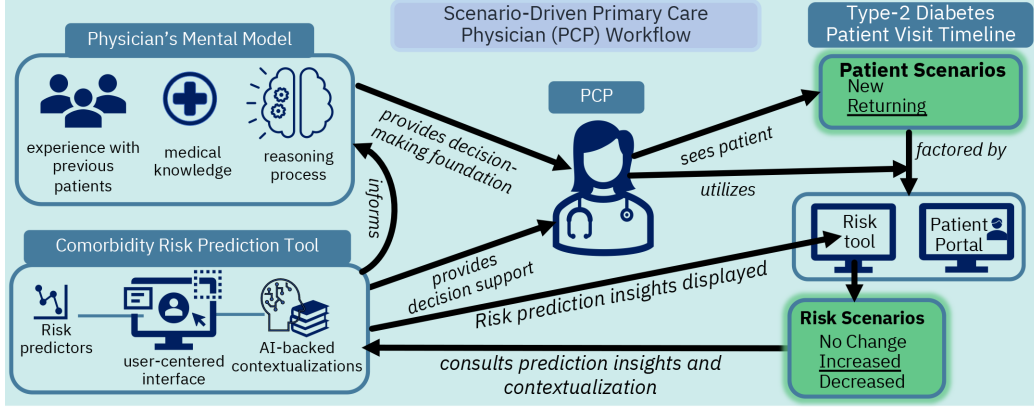
41

Figure 13: Illustration of findings for a Primary Care Physician (PCP) - one target persona among clinicians - whose workflow is dependent on the patient context and the clinician-patient history. We show an example scenario where the predicted risk for a returning patient has increased since the last visit. This context and the PCP's mental model drive the PCP's following actions, such as differential treatment decisions made by probing reasons for the increased risk.

What can be the impact of our contextual explanations beyond the comorbidity risk prediction setting? To identify specific scenarios in which our contextualizations might be most impactful, we discussed with a clinical expert typical situations of T2DM patient care by clinician type and patient characteristics to understand where risk predictions could help clinicians improve patient care planning. While current literature [13] on designing healthcare models points to a user-centered approach, from this understanding of clinician workflows, our discussion showed the importance of grounding such user centered work in specific clinical scenarios (see Fig 13). For example, it became clear that dashboards containing contextual explanations around risk prediction could be used by clinicians in different ways. For example, a clinician seeing a patient for the first time and/or for the first diagnosis would be interested in creating their mental model of the patient's diagnosis and understanding the causes for the risk, whereas a clinician seeing a returning patient with a previously established diagnosis (where clinicians might be more interested in understanding any changes to the risk prediction over time and the effectiveness of various interventions). Hence, based on these understandings, we formalized our use case to provide contextual explanations to a PCP around the predicted risk of CKD among new T2DM patients at their first diagnosis.

Additionally, while we focus our approach in the risk prediction of CKD among T2DM patients, the contextualization approach can be applied to other comorbid risk prediction settings given access to authoritative guidelines in the disease area, and likewise, the themes that we analyze from our expert panel interviews are also general enough to be considered applicable in other disease settings. These themes indicate a larger need for AI systems to support insights from multiple data and knowledge sources and present them as actionable and contextual explanations [9] (also pointed out in the self-explanation scorecard from [51]). In summary, our approach is a step towards extracting clinically relevant context from different data and knowledge sources, including guidelines, patient data, and medical ontologies, and using these contexts to augment and explain answers to a list of clinically relevant questions that can help clinicians reason and interpret the risk predictions for patients.

What are some future directions that emerge from the clinician discussions?
Some subthemes under *Theme 2: Highlighting Actionability* provide future directions, highlighting actionable factors and suggesting specific actions to reduce CKD risk are not currently addressed by our contextualization approach. These sub-themes require more investigation and development of methods to identify actionable, most relevant factors to CKD risk. Another point which we observed is that some of the factors that the model picked up on are not covered by the T2DM guidelines, and could either be factors only relevant to CKD, or are those that are not considered to be well-known enough to be covered in position statements like CPGs. We are also investigating how to combine insights from multiple guidelines (also mentioned in [52]) and if that would be useful. In summary, while some of these themes provide validation for the modules we currently support in our multi-method approach to provide context, others offer directions for us to build towards, such as enabling connections to external medication databases, supporting temporality in post-hoc explanations of risk, and efforts to better present answers in terms of relevance and actionability. We are also considering interviewing more user groups within the clinical domain to strengthen an understanding of where such a risk prediction tool would be most impactful. Future steps would also include a practical study in a clinical setting to further assess the utility of our method.

## 6. Related Work

Our methods build on both expert feedback and past efforts to leverage clinical domain knowledge for generating explanations within AI assistants. Some notable and relevant past works include: MYCIN [1], where domain literature was encoded as rules and trace-based explanations, which addressed 'Why,' 'What,' and 'How,' were provided for the treatment of infectious diseases; the DESIREE project [53], where case, experience, and guideline-based knowledge was used to generate insights relevant to patient cases; and a mortality risk prediction effort [54] of cardiovascular patients, where a probabilistic model was utilized to combine insights from patient features, and domain knowledge, to ascertain patient conformance to the literature. However, these approaches are either not flexible nor scalable for the ingestion of new domain knowledge [1, 53], or are narrowly focused in their approach to explanations along limited dimensions [54]. We attempt to allow clinicians to probe the supporting evidence systematically and thoroughly while asking holistic questions about the supporting evidence(s) to understand their patients better.

On the guideline QA front, there have been several efforts on representation formats for guidelines and more recent work on applying machine learning and language model approaches on guidelines for upstream tasks other than QA [50, 46, 55]. Guideline representation efforts attempt to model guidelines as rules that can then be checked against patient data for conformance. While rule engineering is more accurate than applying machine learning models, it is not scalable without human effort. In a more scalable effort, Schlegel et al. [55], have shown how a standard NLP pipeline of tools like named entity recognizers, syntactic, semantic, and dependency parsers can be applied to convert guideline text into annotated text snippets. However, their system, ClinicalTractor, is not available for reuse yet, and hence we could not use it for the semantic annotation portion of our QA pipeline. Another similar effort is from Hussain et al. [50], where they use heuristic patterns to identify different composition patterns in guideline sentences. These guideline natural language understanding efforts while useful, still require significant effort to be used upstream by QA approaches and could instead be used to augment QA approaches such as ours with additional information that can help improve semantic and syntactic coherence of answers.

On the other hand, with the rise of LLM [36], several papers have been published on adapting language models to the biomedical and clinical do-

mains by the pre-training of these models on large biomedical corpora [38, 56]. There are currently very few efforts on the applications of these domain adapted language models to medical guidelines [46]. Hussein and Woldek [46] have applied language models to identify condition-action statements from three medical guidelines, and they find that the combination of syntactic and semantic features from Metamap can boost the performance of language models like BioBERT. However, it is not immediately clear how the extraction of condition-action statements can be used in a QA setting where the range of question types like those we support goes beyond condition-action pairs. For example, questions asking about diagnoses features don't always have a condition to be searched against. Contrarily, Sarrouti et al. [57], have designed a semantic biomedical question answering system that achieves the state-of-the-art results on the BioASQ challenge by using UMLS similarity scores and a novel passage retrieval algorithm to find candidate answers from Pubmed documents. In their future work, they list the needs for large training samples as a limiting factor to use a deep learning algorithm. While we agree, we have shown how knowledge augmentation algorithms can improve the performance of deep learning language models in new settings such as unseen guidelines.

Several studies have also tried to utilize patient data to query the literature for applicable evidence or treatment suggestions. However, very few of these studies combine multiple modalities and sources of data and knowledge for querying. Agosti et al. [58], conducted an analysis of query reformulation techniques for precision medicine. Natarajan et al. [59], conducted an analysis of clinical queries in an electronic health record search utility and found that queries on diseases and lab results were most searched. Patel et al. [60], matched patient records to clinical trials using ontologies and used a purely logical A-box and T-box approach to query literature.

Finally, several studies have hinted that model explanations alone are not sufficient and indicate that context can be an important dimension to make these explanations more useful. We summarize a few studies which either contextualize model explanations with links to knowledge or can provide context around risk prediction scores to make them more useful. Rieger et al. [61], found that interpretations are useful by penalizing explanations to align neural networks with prior knowledge. Zhang et al. [62], presented context-aware and time-aware attention-based model for Disease Risk Prediction with Interpretability. They used disease code hierarchies as context in RNN network's attention layer. Weber et al. [63], attempted a Knowledge-

based XAI through CBR and found that there is more to explanations than models can tell. Yao et al. [64], refined Language Models with Compositional Explanations by align LLM and post-hoc output with human knowledge and Tonekaboni et al. [15], analyzed what clinicians want and found that Contextualizing Explainable Machine Learning for Clinical End Use.

## 7. Conclusion

Contextual explanations have been posited to be useful for clinicians for real-world usage of AI models. In this paper, we have developed an end-to-end AI systems and studied the feasibility and usability of extracted contextual explanations from medical guidelines using state-of-the art QA methods. We have focused our study in a risk prediction use-case for CKD among T2DM patients and have conducted both quantitative and qualitative analysis. Upon conversations with clinicians, we have selected three entities of interest in the risk prediction setting to provide contextual explanations along - the patient, their predicted risk, and the model explanations of their risk. Crucially, we have identified several themes covering the explainability needs of clinicians. The supported contextual explanations support some of these themes and thus improves clinician's confidence is using AI supported systems. We also found state-of-the art large language models to be effective in extracting such contexts, especially for certain disease sub-groups. While our results support the hypothesis that presenting contextual explanations to clinicians is both feasible and usable, the performance and requirement gaps points to the need for further research in this field. For example, while we have considered three domain sources for the contexts in this paper, the themes from the expert panel interviews also indicate that there may be value to connecting to other sources, including extracting additional guideline details from tables and flowcharts, and potentially involving multiple layers of the evidence pyramid to include such sources as systematic reviews, randomized clinical trials, cohort studies, and expert opinions. Similarly, novel machine learning techniques such as weak-supervision or one-shot learning may be need to improve the quality of extracted contextual explanations. A combination of both may also enable other approaches such as 'prompt engineering' whereby patient data can be used to seed the QA model questions and get richer response. Our future research will be directed at overcoming the aforementioned opportunities. Overall, by closely working with clinical experts and adopting inter-disciplinary approaches, from the use case crystal-

lization and methods development, to the evaluation stages, we have shown the promise in supporting clinically relevant contexts to help clinicians better interpret risk prediction scores and their model explanations.

## Acknowledgments

## References

[1] E. H. Shortliffe, Mycin: a rule-based computer program for advising physicians regarding antimicrobial therapy selection., Tech. rep., Dept. of Computer Sci., Stanford University Stanford (1974).

[2] G. Briganti, O. Le Moine, Artificial intelligence in medicine: today and tomorrow, Frontiers in medicine 7 (2020) 27.

[3] W. Swartout, C. Paris, J. Moore, Explanations in knowledge systems: Design for explainable expert systems, IEEE Expert 6 (3) (1991) 58–64.

[4] D. Gunning, Explainable artificial intelligence (xai), Defense Advanced Research Projects Agency (DARPA), nd Web 2 (2017).

[5] S. Chari, D. Gruen, O. W. Seneviratne, D. L. McGuinness, Foundations of explainable knowledge-enabled systems, in: Knowledge Graphs for eXplainable Artificial Intelligence, Vol. 47, IOS Press, 2020, pp. 23 – 48.

[6] S. Dey, P. Chakraborty, B. C. Kwon, A. Dhurandhar, M. Ghalwash, F. J. S. Saiz, K. Ng, D. Sow, K. R. Varshney, P. Meyer, Human-centered explainability for life sciences, healthcare, and medical informatics, Patterns 3 (5) (2022) 100493.

[7] M. Ghassemi, L. Oakden-Rayner, A. L. Beam, The false hope of current approaches to explainable artificial intelligence in health care, The Lancet Digital Health 3 (11) (2021) e745–e750. doi:https://doi.org/10.1016/S2589-7500(21)00208-9.
URL https://www.sciencedirect.com/science/article/pii/S2589750021002089

[8] P. Chakraborty, B. C. Kwon, S. Dey, A. Dhurandhar, D. Gruen, K. Ng, D. Sow, K. R. Varshney, Tutorial on human-centered explainability for healthcare, in: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2020, pp. 3547–3548.

[9] S. Chari, O. Seneviratne, D. M. Gruen, M. A. Foreman, A. K. Das, D. L. McGuinness, Explanation ontology: A model of explanations for user-centered ai, in: International Semantic Web Conference, Springer, 2020, pp. 228–243.

[10] L. H. Gilpin, D. Bau, B. Z. Yuan, A. Bajwa, M. Specter, L. Kagal, Explaining explanations: An overview of interpretability of machine learning, in: 2018 IEEE 5th International Conference on data science and advanced analytics (DSAA), IEEE, 2018, pp. 80–89.

[11] V. Arya, R. K. Bellamy, P.-Y. Chen, A. Dhurandhar, M. Hind, S. C. Hoffman, S. Houde, Q. V. Liao, R. Luss, A. Mojsilović, et al., One explanation does not fit all: A toolkit and taxonomy of ai explainability techniques, arXiv preprint arXiv:1909.03012 (2019).

[12] Q. V. Liao, D. Gruen, S. Miller, Questioning the ai: informing design practices for explainable ai user experiences, in: Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, 2020, pp. 1–15.

[13] D. Wang, Q. Yang, A. Abdul, B. Y. Lim, Designing theory-driven user-centric explainable ai, in: Proceedings of the 2019 CHI conference on human factors in computing systems, 2019, pp. 1–15.

[14] F. Doshi-Velez, B. Kim, Towards a rigorous science of interpretable machine learning, arXiv preprint arXiv:1702.08608 (2017).

[15] S. Tonekaboni, S. Joshi, M. D. McCradden, A. Goldenberg, What clinicians want: contextualizing explainable machine learning for clinical end use, in: Machine learning for healthcare conference, PMLR, 2019, pp. 359–380.

[16] S. Chari, D. Gruen, O. W. Seneviratne, D. L. McGuinness, Directions for explainable knowledge-enabled systems, in: Knowledge Graphs for eXplainable Artificial Intelligence, Vol. 47, IOS Press, 2020, p. 245.

[17] H. Lieberman, T. Selker, Out of context: Computer systems that adapt to, and learn from, context, IBM systems journal 39 (3.4) (2000) 617–632.

[18] A. K. Dey, G. D. Abowd, A. Wood, Cyberdesk: A framework for providing self-integrating context-aware services, Knowledge-based systems 11 (1) (1998) 3–13.

[19] S. Chari, M. Qi, N. N. Agu, O. Seneviratne, J. P. McCusker, K. P. Bennett, A. K. Das, D. L. McGuinness, Making study populations visible through knowledge graphs, in: International Semantic Web Conference, Springer, 2019, pp. 53–68.

[20] D. W. Challener, L. J. Prokop, O. Abu-Saleh, The proliferation of reports on clinical scoring systems: issues about uptake and clinical utility, Jama 321 (24) (2019) 2405–2406.

[21] I. A. Videha Sharma, S. van der Veer, G. Martin, J. Ainsworth, T. Augustine, Adoption of clinical risk prediction tools is limited by a lack of integration with electronic health records, BMJ Health & Care Informatics 28 (1) (2021).

[22] M. Banning, A review of clinical decision making: models and current research, Journal of clinical nursing 17 (2) (2008) 187–195.

[23] A. L. Rosner, Evidence-based medicine: revisiting the pyramid of priorities, Journal of Bodywork and Movement Therapies 16 (1) (2012) 42–49.

[24] H. Lakkaraju, D. Slack, Y. Chen, C. Tan, S. Singh, Rethinking explainability as a dialogue: A practitioner's perspective, arXiv preprint arXiv:2202.01875 (2022).

[25] K. Främling, Decision theory meets explainable ai, in: International Workshop on Explainable, Transparent Autonomous Agents and Multi-Agent Systems, Springer, 2020, pp. 57–74.

[26] About chronic diseases — cdc, `https://www.cdc.gov/chronicdisease/about/index.htm`, accessed: 2021-05-26.

[27] Chronic kidney disease basics — chronic kidney disease initiative — cdc, `https://www.cdc.gov/kidneydisease/basics.html`, accessed: 2021-05-25.

[28] M. H. Murad, Clinical practice guidelines: a primer on development and dissemination, in: Mayo Clinic Proceedings, Vol. 92, Elsevier, 2017, pp. 423–433.

[29] R. Graham, M. Mancher, D. M. Wolman, S. Greenfield, E. Steinberg, et al., Trustworthy clinical practice guidelines: Challenges and potential, in: Clinical Practice Guidelines We Can Trust, National Academies Press (US), 2011.

[30] P. Suryanarayanan, P. Chakraborty, P. Madan, K. Bore, W. Ogallo, R. Chandra, M. Ghalwash, I. Buleje, S. Remy, S. Mahatma, P. Meyer, J. Hu, Disease Progression Modeling Workbench 360, arXiv preprint arXiv:2106.13265 (2021).

[31] K. Rufibach, Use of brier score to assess binary predictions, Journal of clinical epidemiology 63 (8) (2010) 938–939.

[32] S. M. Lundberg, B. Nair, M. S. Vavilala, M. Horibe, M. J. Eisses, T. Adams, D. E. Liston, D. K.-W. Low, S.-F. Newman, J. Kim, et al., Explainable machine-learning predictions for the prevention of hypoxaemia during surgery, Nature Biomedical Engineering 2 (10) (2018) 749.

[33] K. S. Gurumoorthy, A. Dhurandhar, G. Cecchi, C. Aggarwal, Efficient data representation by selecting prototypes with importance weights, in: 2019 IEEE International Conference on Data Mining (ICDM), IEEE, 2019, pp. 260–269.

[34] D. Care, Standards of medical care in diabetes 2021, Diabetes Care 44 (Suppl 1) (2021).

[35] L. Richardson, Beautiful soup documentation, Dosegljivo: https://www. crummy. com/software/BeautifulSoup/bs4/doc/.[Dostopano: 7. 7. 2018] (2007).

[36] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, ACL Anthology (2018).

[37] A. Otegi, J. A. Campos, G. Azkune, A. Soroa, E. Agirre, Automatic evaluation vs. user preference in neural textual QuestionAnswering over COVID-19 scientific literature, in: Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020, Association for Computational Linguistics, Online, 2020. `doi:10.18653/v1/2020.nlpcovid19-2.15`.
URL `https://www.aclweb.org/anthology/2020.nlpcovid19-2.15`

[38] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, J. Kang, Biobert: a pre-trained biomedical language representation model for biomedical text mining, Bioinformatics 36 (4) (2020) 1234–1240.

[39] W. Yoon, J. Lee, D. Kim, M. Jeong, J. Kang, Pre-trained language model for biomedical question answering, in: PKDD/ECML Workshops, 2019.

[40] Huggingface — bioclinicalbert-adr, `https://huggingface.co/anindabitm/sagemaker-BioclinicalBERT-ADR`, accessed: 2022-8-11.

[41] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al., Huggingface's transformers: State-of-the-art natural language processing, arXiv preprint arXiv:1910.03771 (2019).

[42] A. R. Aronson, F.-M. Lang, An overview of metamap: historical perspective and recent advances, Journal of the American Medical Informatics Association 17 (3) (2010) 229–236.

[43] K. Donnelly, et al., Snomed-ct: The advanced terminology and coding system for ehealth, Studies in health technology and informatics 121 (2006) 279.

[44] T. Knoll, F. Moramarco, A. P. Korfiatis, R. Young, C. Ruffini, M. Perera, C. Perstl, E. Reiter, A. Belz, A. Savkov, User-driven research of medical note generation software, arXiv preprint arXiv:2205.02549 (2022).

[45] T. J. Pollard, A. E. Johnson, J. D. Raffa, R. G. Mark, tableone: An open source python package for producing summary statistics for research papers, JAMIA open 1 (1) (2018) 26–31.

[46] H. Hematialam, W. W. Zadrozny, Identifying condition-action statements in medical guidelines: Three studies using machine learning and domain adaptation (2021).

[47] S. Teufel, An overview of evaluation methods in trec ad hoc information retrieval and trec question answering, Evaluation of text and speech systems (2007) 163–186.

[48] R. Gatta, M. Vallati, C. Fernandez-Llatas, A. Martinez-Millana, S. Orini, L. Sacchi, J. Lenkowicz, M. Marcos, J. Munoz-Gama, M. Cuendet, et al., Clinical guidelines: a crossroad of many research areas. challenges and opportunities in process mining for healthcare, in: International Conference on Business Process Management, Springer, 2019, pp. 545–556.

[49] D. Riaño, M. Peleg, A. Ten Teije, Ten years of knowledge representation for health care (2009–2018): Topics, trends, and challenges, Artificial intelligence in medicine 100 (2019) 101713.

[50] M. Hussain, J. Hussain, T. Ali, S. I. Ali, H. S. M. Bilal, S. Lee, T. Chung, Text classification in clinical practice guidelines using machine-learning assisted pattern-based approach, Applied Sciences 11 (8) (2021) 3296.

[51] S. T. Mueller, E. S. Veinott, R. R. Hoffman, G. Klein, L. Alam, T. Mamun, W. J. Clancey, Principles of explanation in human-ai systems, arXiv preprint arXiv:2102.04972 (2021).

[52] D. F. Sittig, A. Wright, J. A. Osheroff, B. Middleton, J. M. Teich, J. S. Ash, E. Campbell, D. W. Bates, Grand challenges in clinical decision support, Journal of biomedical informatics 41 (2) (2008) 387–392.

[53] B. Seroussi, J.-B. Lamy, N. Muro, N. Larburu, B. D. Sekar, G. Guézennec, J. Bouaud, Implementing guideline-based, experience-based, and case-based approaches to enrich decision support for the management of breast cancer patients in the desiree project., in: EFMI-STC, 2018, pp. 190–194.

[54] A. Raghu, J. Guttag, K. Young, E. Pomerantsev, A. V. Dalca, C. M. Stultz, Learning to predict with supporting evidence: applications to clinical risk prediction, in: CHIL '21: Proceedings of the Conference on Health, Inference, and Learning, ACM, 2021, pp. 95–104.

[55] D. R. Schlegel, K. Gordon, C. Gaudioso, M. Peleg, Clinical tractor: A framework for automatic natural language understanding of clinical practice guidelines, in: AMIA Annual Symposium Proceedings, Vol. 2019, American Medical Informatics Association, 2019, p. 784.

[56] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, H. Poon, Domain-specific language model pretraining for biomedical natural language processing, ACM Transactions on Computing for Healthcare (HEALTH) 3 (1) (2021) 1–23.

[57] M. Sarrouti, S. O. El Alaoui, Sembionlqa: a semantic biomedical question answering system for retrieving exact and ideal answers to natural language questions, Artificial intelligence in medicine 102 (2020) 101767.

[58] M. Agosti, G. M. Di Nunzio, S. Marchesin, An analysis of query reformulation techniques for precision medicine, in: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, 2019, pp. 973–976.

[59] K. Natarajan, D. Stein, S. Jain, N. Elhadad, An analysis of clinical queries in an electronic health record search utility, International journal of medical informatics 79 (7) (2010) 515–522.

[60] C. Patel, J. Cimino, J. Dolby, A. Fokoue, A. Kalyanpur, A. Kershenbaum, L. Ma, E. Schonberg, K. Srinivas, Matching patient records to clinical trials using ontologies, in: The Semantic Web, Springer, 2007, pp. 816–829.

[61] L. Rieger, C. Singh, W. Murdoch, B. Yu, Interpretations are useful: penalizing explanations to align neural networks with prior knowledge, in: International Conference on Machine Learning, PMLR, 2020, pp. 8116–8126.

[62] X. Zhang, B. Qian, Y. Li, S. Cao, I. Davidson, Context-aware and time-aware attention-based model for disease risk prediction with interpretability, IEEE Transactions on Knowledge and Data Engineering (2021).

[63] R. Weber, M. Shrestha, A. J. Johs, Knowledge-based xai through cbr: There is more to explanations than models can tell, arXiv preprint arXiv:2108.10363 (2021).

[64] H. Yao, Y. Chen, Q. Ye, X. Jin, X. Ren, Refining language models with compositional explanations, Advances in Neural Information Processing Systems 34 (2021).

[65] W. McKinney, et al., Pandas: a foundational python library for data analysis and statistics, Python for high performance and scientific computing 14 (9) (2011) 1–9.

[66] Y. Chen, A. Subburathinam, C.-H. Chen, M. J. Zaki, Personalized food recommendation as constrained question answering over a large-scale food knowledge graph, in: Proceedings of the 14th ACM International Conference on Web Search and Data Mining, 2021, pp. 544–552.

[67] E. W. Dijkstra, et al., A note on two problems in connexion with graphs, Numerische mathematik 1 (1) (1959) 269–271.

[68] A. Hagberg, P. Swart, D. S Chult, Exploring network structure, dynamics, and function using networkx, Tech. rep., Los Alamos National Lab.(LANL), Los Alamos, NM (United States) (2008).

[69] J.-B. Lamy, A. Venot, C. Duclos, Pymedtermino: an open-source generic api for advanced terminology services, in: Digital Healthcare Empowering Europeans, IOS Press, 2015, pp. 924–928.

## Appendix  A.  QA Architecture

We discuss in this section the extraction methods for information retrieval from the different data sources that we support as context, the sub-modules which help in question and explanation generation, and finally, the standard evaluation modules which can output scores for some of the question types.

### Appendix  A.1.  Information Retrieval

Here we describe the information extraction portion of our QA pipeline, as seen in part A). of Fig. 4. We support the extraction of context from three domain sources in our QA approach, including patient data, medical ontologies like Clinical Classification Software (CCS) codes [10] and medical guidelines from ADA Standards of Care 2021 (as introduced in Sec. 2). We

---

[10]https://www.hcup-us.ahrq.gov/toolssoftware/ccs10/ccs10.jsp

query patient data from Limited Claims Explorys Dataset (LCED) claim records (see Sec. 2.1) on-demand, either when we need to create questions based on patient parameters or when we need to include these patient values in answers to questions about the patient. As for extracting content from CCS codes, we download a static version of the 'CCS for ICD-10-PCS Tool' file and use four fields from this downloaded file - 'Field 2: CCS Category, Field 3: Code Description, Field 7: Multi-level 2 Category, and Field 8: Multi-level 2 Category Description'. Field 8 in particular, connects the lower level disease codes we see in the patient data to their higher level 2 groupings (e.g., Essential Hypertension's Multi-level 2 Category Description is Disease of the Circulatory System).

We extract content from the HTML or web version of the 'Standards of Medical Care in Diabetes' [34] guidelines, published by the American Diabetes Association (ADA) [11] using a Python library, BeautifulSoup [35]. The ADA CPGs, are updated annually and are an example of a well-maintained CPG, whose format has been fairly consistent over the past decade. The ADA CPG is released both in web-friendly formats like HTML and in PDF formats. The content in the ADA CPG is split across chapters, where each chapter focuses on a different aspect of diabetes management or treatment. Furthermore, each chapter contains different recommendation groups, each containing recommendations, discussions, and references for a sub-area within the chapter (Fig. 6) [19]. The recommendations themselves are supported by different grades of evidence and are graded accordingly. For example, a recommendation supported by a systematic review or meta-analysis is assigned a grade A, whereas only an expert opinion is graded as E. Also, while the references for the recommendations are not made available as direct associations, they can be found within the discussion supporting the recommendations. Mainly, within each chapter in the ADA guidelines, we extract the content of different sections, including the recommendations, supporting discussions, and references within these sections. We then write the output of this extraction, mirroring the structure of the original guidelines (see Fig. 6), to a semi-structured JSON format, which is then used within our QA modules.

---

[11]https://care.diabetesjournals.org/content/44/Supplement_1

*Appendix  A.2.  QA Architecture*

Here we describe the processes and modules, seen in part B) of Fig.  4, that are involved in generating questions based on question types for each patient (see Tab.  2) and the answers from the different domain sources of context in our risk prediction setting.

***Question  Generation*** The *question generation module* almost always creates templated questions using Python's native support for String Templates, [12] and does so based on patient data, more specifically from the patient's diagnoses codes, lab values, and medication list. The patient's diagnoses codes are sometimes abstracted from their higher-level disease groupings supported by the CCS scheme. A subset of these diagnosis codes can be included in the features that the post-hoc explanation module found were contributing to the patient's predicted risk. In an attempt to provide more context around these features, we create instances of the type 3 question , e.g., "What can be done for this patient's essential hypertension?" We also support the creation of two standard, non-variant questions for each patient, i.e., whose values don't change from patient data, that can help clinicians easily interpret their predicted risk (question type 1) and their `T2DM` state (question type 2).

Moreover, as can be seen from Tab.  2, each of the question types that we support on a per patient basis is populated from different data sources. Hence, we have developed different answering methods for each, including simple lookups and knowledge augmented language model capabilities, including combinations of either a LLM + value range comparison or LLM + knowledge augmentation. We provide examples of questions and answers for each question type in Appendix  C.

***Answer  Generation*** In our answer generation module of our QA approach, we support different submodules that can output answers to questions related to the question types. The answer generation module is capable of inputting questions generated by the previous question generation module and interacting with extracted content from our supported data sources.

*Template-based Answer Generation:* Question types 1 and 2 from Tab.  2 can be addressed by simple query lookups of our supported data sources. We populate the Python String Template object with the results of the queried components retrieved by using the widely-used Python Pandas library [65].

---

[12]https://docs.python.org/3/library/string.html

This process of creating natural language templates that can then be populated with values on a per-patient basis is supported by the Template-based Answer Generation module of our QA pipeline (Fig. 4). The results of these questions can be summarized or built from structured datasets, like patient data, their model outputs, like risk predictions, and features contributing to their predicted risk and population averages. Hence, there is no fuzziness in the results, which is why we don't evaluate the accuracy of this submodule.

This submodule is also leveraged in combination with other answer generation submodules when there is a pattern in the answers and slots to be filled. We discuss these details shortly after we set up our knowledge augmented language model capabilities and their usage.

*Numerical Range Comparison:* BERT LLMs cannot currently determine if a question that has a numerical value comparison, e.g., "What can be done for patients, whose Hemoglobin A1C > 10?", falls in the range of the answer returned. However, clinicians often look for recommendations that match patients' lab values in clinical settings such as ours. Further, within the ADA 2021 guidelines, there are mentions for suggestions based on different ranges for lab values, e.g., a recommendation from the Pharmacological Chapter of these guidelines has a recommendation with the mention of "when A1C levels (> 10% [ 86 mmol/mol ]." Hence, for question type 3 from Tab. 2, we need to determine if the patient's lab values are in the range of the answer returned by the LLM module. To address this requirement of performing numerical range comparison between the question and answer produced by LLM, we leverage syntactic parsing capabilities (similar to [66]) to identify numerical phrases in the question and answer and then determine if each numerical phrase from the question is in range of the same in the answer.

We use Natural Language Toolkit (NLTK) chunking and parsing functionalities to identify noun phrases, comparatives, and numerical mentions within both the question and answer. We then write regular expression (regex) rules to identify the patterns of the positional tags returned by NLTK that can constitute numerical phrases. For each of the numerical phrases, we convert them into a tuple of "(noun phrase, [upper bound, lower bound])." This tuple representation allows us to go through the phrases between the question and answer iteratively, and for those that match on the noun phrase dimension, identify if the ranges are in agreement. With these steps, we can then populate an answer using the Template-based Answer generation module, which says if the answer outputted by LLM is in/out of the range of the question. Hence, in this manner, we enhance the capabilities of BERT LLMs

for numerical range comparisons via rule-based syntactic methods, which is also why we consider this step a rule augmentation of LLMs' capabilities.

*Knowledge Augmentations to LLMs:* Transformer based LLM approaches like BERT work on sequences of words that are often seen together and their surrounding words, but don't leverage the semantics of whether these words are diseases, medications, or biological processes. We found that in the absence of this semantic knowledge, we would often get answers from the LLM that don't correlate on a semantic level with the question. For example, a sentence from the Comorbidities chapter of the ADA 2021 CPG on Dementia was returned as a valid answer to a question asking about an Abdominal Hernia. To eliminate such answers, we explored options for a biomedical semantic mapper and zeroed in on the National Library of Medicine (NLLM)'s Metamap tool [42]. We choose Metamap because of its extensive coverage of biomedical semantic types and its ability to capture entity mentions within the ADA 2021 CPG. Within our pipeline, we have integrated a Python wrapper for Metamap[13] that can recognize biological entities within the guideline text and their semantic types (e.g., dsyn: disease or syndrome, bpoc: biological processes, etc. for a complete list of types returned by Metamap see: [14]).

We run Metamap on question types 3 and 4 from Tab. 2 to only output answers from BERT when there is a valid semantic match between the question and answer. Specifically, using this knowledge augmentation module, for question type 3, we only output answers whose matched term is a noun and is recognized as a disease term by Metamap, and similarly, for question type 5, we only output answers whose matched term is a noun and is recognized as medication by Metamap. We have observed that depending on the mention of a biological entity in the text, a disease term can be recognized as a disease, biological process, or a finding by Metamap. Hence, we allow for flexibility among semantic types, in filtering disease matches for question type 3. For example, we want to allow answers with the mention of the term 'hypertensive' for a question on hypertension, although hypertensive is identified as a finding by Metamap.

Additionally, given this ability to filter based on semantic types, we want

---

[13]PyMetamap: `https://github.com/AnthonyMRios/pymetamap`

[14]`https://lhncbc.nLLM.nih.gov/ii/tools/MetaMap/Docs/SemanticTypes_2018AB.txt`

to allow additional answers with mentions of related diseases. To provide more broad answers, we use the UMLS Concept Unique Identifier (CUI) codes from the Metamap returned outputs to map to Snomed-CT disease codes [43]. From the mapped Snomed-CT disease codes, we can traverse the Snomed-CT disease tree to identify how many hops apart question and answer disease codes are and if the answer codes are an ancestor of those in the question. We operate on the idea that answers about the parent disease code apply to children nodes. For example, a question about "What can be done for Asthma" can borrow from an answer on "What can be done for respiratory diseases?" Conversely, if disease codes in the question and answer are far apart in the Snomed tree, it would signify that they are semantically less related. In addition to Metamap codes, we append to the candidate guideline sentences the hop distances from each question computed by applying the Dijkstra's algorithm [67] on an uploaded Snomed graph in Python package NetworkX [68] and ancestor values derived from using Python library PyMedTermino's [69] is_ancestor function.

We use the outputs of these knowledge augmentation modules to both pre-filter and post-sort the LLMs answers. The LLM and LLM + post sorting settings 1, 2 and 5, were run against 410 passage chunks of guideline text, of average length 267 tokens, since BERT has a 512 token limit for an answer passage. The LLMs on the pre-filtering settings 3 and 5 were run on passage chunks of variable length, depending on the number of filtered sentences to be passed to the LLM model. In the pre-filtering settings, we varied the values of the features that we were filtering by to understand which feature values improve accuracy. In essence, the pre-filtering settings can be thought of as algorithmic knobs to control the set of answers that the LLM has to process. In contrast, in the post-filtering settings we sorted the LLM's answers by feature values, and here we could control the ordering of answers to be outputted. In the pre-filtering setting 2, we filter the guideline sentences by length of disease overlap with the question. In pre-filtering setting 4, we have more possibilities in the feature column because the number of Snomed disease hops between a question and answer can range between a continuous range of integer values. We report if restricting the number of hops to allow for more general yet precise answers improves accuracy. Similarly, in the post-sorting settings, 2 and 4, we use the feature values from the knowledge augmentation modules in addition to the LLM's own confidence scores to rank answers. Specifically, in setting 2, we sort the LLMs answerset on variations to a combination of length of disease overlap between question

and answer Metamap phrases and the LLM confidence scores. In setting 5, we sort the LLMS answerset based on variations to a combination of sum of hops between question and answer Snomed disease codes, number of Snomed ancestors in the answer and LLM confidence scores.

We report the accuracies for answers that address questions of question types 3 from Tab. 2, that use these knowledge augmentations in the results section (Sec. 4). We have written functions that use the NLTK toolkit in our evaluation submodule to generate standard, natural language processing (NLP), accuracy scores like F1, precision, recall, and BLEU. Overall, the integration of a semantic mapping tool helps us enhance the capabilities of BERT LLMs for more precise and better semantic matches via knowledge-driven methods.

## Appendix  B. True Label Annotations for Guideline Questions

We generated annotations for creating the gold standard dataset for feature importance questions (question type 3 from question types supported by our QA approach, see Tab. 2), whose results are presented in the main manuscript of the paper. These questions included disease feature importances of diagnostic value. The annotations were done by the first author by reading the ADA 2021 CPG and looking for answers to these questions. We also ran our annotations by a medical expert on our team, who is also a co-author on this paper, to verify if they were clinically meaningful. The medical expert validated 47 questions out of the 71 questions annotated by the first author and we report results on this expert validated set in Sec. Appendix  C. To the best of our knowledge, we are the first to report a dataset of questions with annotated answers on the ADA 2021 CPG. It is to be noted that these annotations are not endorsed by the ADA. Nevertheless, we hope that our annotations can serve as a valuable resource for academic advances in the clinical informatics community.

## Appendix  C. Results

Here we present additional material to support the results of our risk prediction and question-answering module described in Sec. 4.2.

*Appendix C.1. Model Performance: Risk-Prediction Model and Post-Hoc Explainers*

We present the performance of the risk prediction models in Tab. C.15. As the table shows, while GRU performs the best overall, depending on the use case, we may want to prefer other models. For the purposes of this paper, we chose MLP as our risk prediction model to benefit from the higher recall (such that the probability of false negatives is low) and high brier-score (to allow a more natural interpretation of our model outputs for clinicians), while still achieving an acceptable level of overall performance (AUC-ROC = 0.59).

Table C.15: Results of CKD risk from different prediction models

| Method | Precision | Recall | AUC-ROC | AUC-PRC | Brier |
|--------|-----------|--------|---------|---------|-------|
| LR | 0.333 | 0.023 | 0.582 | 0.215 | 0.127 |
| MLP | 0.139 | **0.977** | 0.587 | 0.224 | **0.621** |
| LSTM | **0.242** | 0.442 | **0.678** | 0.263 | 0.208 |
| GRU | 0.240 | **0.605** | 0.677 | **0.311** | 0.220 |



Figure C.14: Feature importance for `CKD` prediction among 20 prototypical patients using SHAP (left), showing absolute importance, and (right) showing feature impact on model prediction w.r.t. presence/absence of features

Feature importance for risk factors found by algorithmic explainers can be further contextualized, as mentioned previously. The left hand column of Fig. C.14 shows the top 20 features for the set of 20 prototypical patients

under investigation. These prototypical patients are all found to be at high-risk for CKD and hence, would be interesting to clinicians within the scope of this T2DM and CKD use case. For these prototypical patients, we present aggregated feature importance, as seen in Tab. 5, to account for HIPAA restrictions. We can see that demographic features, such as age and the presence of other disorders, such as 'other skin disorders,' were found to be important for the CKD risk prediction. The right side of Fig. C.14 shows an alternate view of the same, providing a view into the spread of individual importance. From this deeper view, we can see that features such as 'calculus of urinary tract' could be the most important drivers of risk for some patients. Such results further support our need to personalize features found to be important for the risk predictions.

While insights about the importance of such features are helpful, such clinical and patho-physiological features may need further contextualization for clinicians. Our structured feedback sessions found that clinicians found the contextual explanations we support around these features helpful, and we cover some of this next.

As seen in Tab. C.16 and Tab. C.17, we provide results numbers on a small set of expert validated answers from our guideline annotations of 12 questions and 47 candidate answers. We are considering methods like weak supervision to increase the expert validation coverage of our annotations. We find that generally the results on the expert validated answers (Tab. C.16 and C.17) follow the trend of accuracy values in the larger annotation set, in that the precision is highest in the pre-filtering by number of Snomed disease hop settings (setting 4) - 0.29 and that the recall and bleu is high in post-filtering by number of disease overlaps between question and answer (SciBERT + KA in Tab. C.17), setting 3 - 0.6 and 0.2 respectively. However, contrary to the larger set of results the F1 are highest in both the vanilla BERT and BioBERT-BioASQ + KA settings, 0.22.

In Fig. C.15, we also report the distribution of the number of hops between disease code pairs found in the question and candidate answers to provide an idea of how far or close the questions are to sentences within the guidelines. As can be seen, most question and answer pairs are between 5 - 15 hops away. We find that the results are best when we filter answer codes less than $3 - 6$ hops away from the question disease codes.

Table C.16: Performance of Guideline QA with different language model approaches reported at mean average precision (map), F1 and recall at top-10 answers and precision at top-1 and top-5 for 12 expert validated questions. Best and second-best values for each column is highlighted in Green and Blue color, respectively. Language model (e.g. BERT) suffixed with KA represents the corresponding knowledge augmented model (e.g. BERT-KA).

| model | bleu | P@1 | P@5 | map | f1 | recall |
|---|---|---|---|---|---|---|
| BERT | 0.155 | 0.363 | 0.262 | 0.267 | 0.224 | 0.365 |
| BioBERT | 0.121 | 0.296 | 0.200 | 0.222 | 0.186 | 0.342 |
| BioBERT-BioASQ | 0.131 | 0.259 | 0.200 | 0.205 | 0.192 | 0.363 |
| BioClinicalBERT-ADR | 0.112 | 0.227 | 0.178 | 0.179 | 0.171 | 0.351 |
| SciBERT | 0.153 | 0.366 | 0.216 | 0.244 | 0.235 | 0.463 |

Table C.17: Performance of Guideline QA with different language model approaches + knowledge augmentations reported at mean average precision (map), F1 and recall at top-10 answers and precision at top-1 and top-5 for 12 expert validated questions. Best and second-best values for each column are highlighted in Green and Blue color, respectively. Language model (e.g. BERT) suffixed with KA represents the corresponding knowledge augmented model (e.g. BERT-KA).

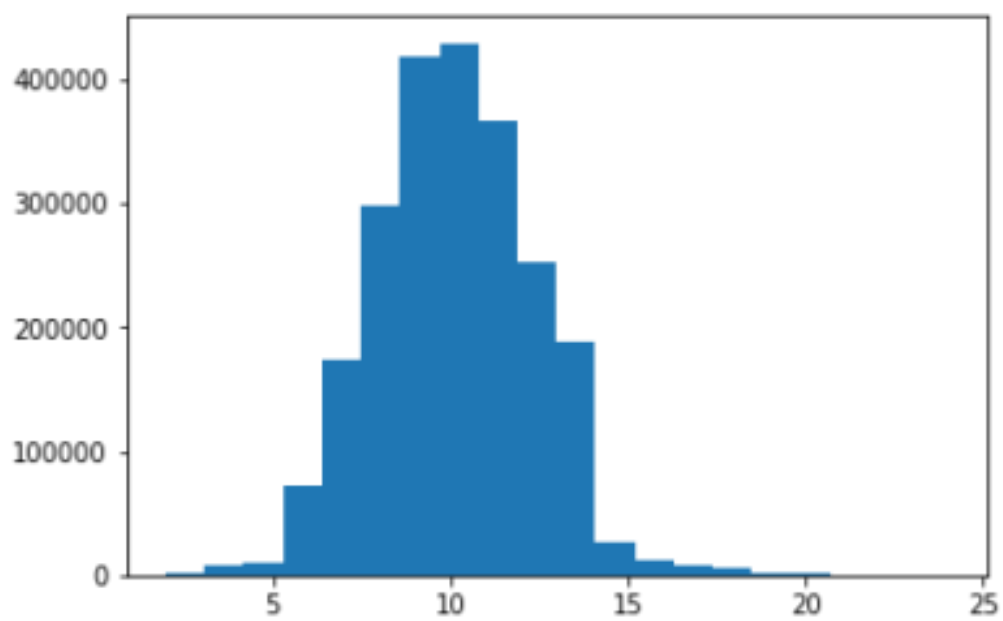| model | bleu | P@1 | P@5 | map | f1 | recall |
|---|---|---|---|---|---|---|
| BERT-KA | 0.021 | 0.296 | 0.296 | 0.296 | 0.127 | 0.081 |
| BioBERT-KA | 0.143 | 0.440 | 0.215 | 0.258 | 0.222 | 0.319 |
| BioBERT-BioASQ-KA | 0.147 | 0.366 | 0.249 | 0.272 | 0.227 | 0.335 |
| BioClinicalBERT-ADR-KA | 0.123 | 0.321 | 0.209 | 0.221 | 0.201 | 0.384 |
| SciBERT-KA | 0.201 | 0.356 | 0.209 | 0.284 | 0.297 | 0.600 |

Figure C.15: Distribution of number of hops between Snomed disease pairs from questions and candidate guideline answers. As can be seen most disease pairs are between 5 - 15 hops apart with 20 being the maximum number of hops.

## Appendix D. Risk Prediction Dashboard Description

We provide additional views of the different panes in our risk prediction dashboard that hosts our supported contextual explanations alongside the different entities that we contextualize, the patient, their risk prediction and the important features found to contribute to the risk prediction. We also support interactions between each of these panes which can help clinicians easily find the content that contextualizes the entities. These interactions or brushing capabilities include (see Fig. D.16, D.17, D.18):



Figure D.16: The risk prediction pane of our risk prediction dashboard, wherein the risk score is displayed alongside a severity of the score on a threshold scale.

- Clicking on the patient details pane brings up questions in the questions in context pane asking about the patient's diabetes state

- When a month is chosen on the history timeline, questions about a commonly accepted diabetes indicator, Hemoglobin A1C (HbA1C), are brought up or filtered in the questions in context pane

- Clicking anywhere in the risk prediction pane brings up questions about the patient's predicted risk

- In the feature importances pane, features can be filtered by the selection of their corresponding, higher-level disease grouping

- Clicking on a feature in the feature importance chart brings up questions about the feature

- Finally, clicking on a higher-level disease grouping also brings up questions about the disease grouping
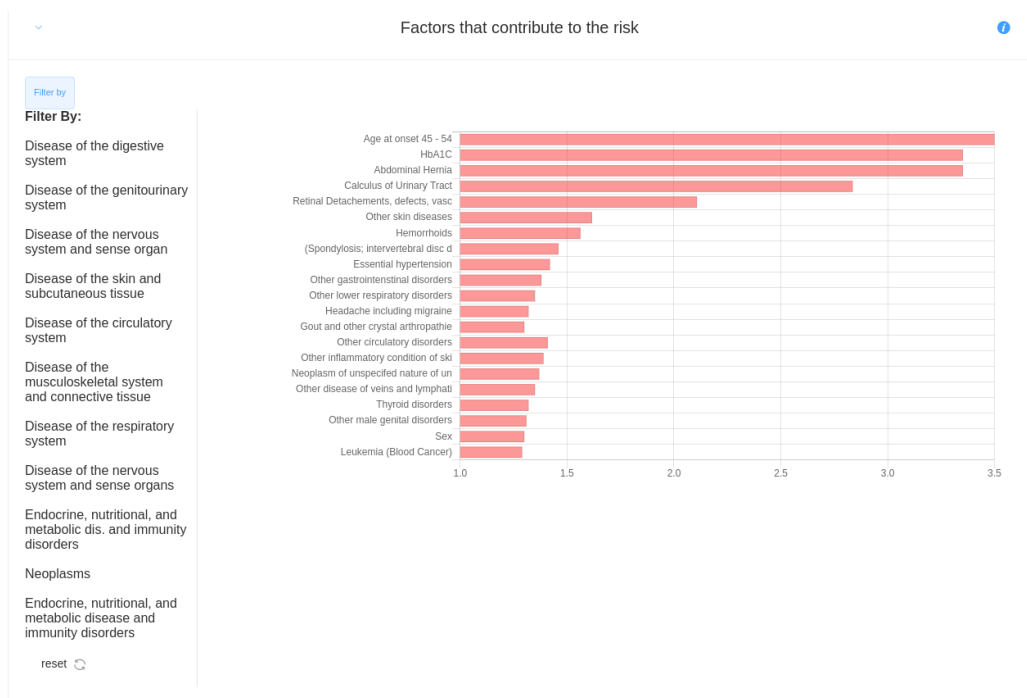
Figure D.17: The feature importances pane of our risk prediction dashboard shows in order of importance the features that contributed to the patient's predicted `CKD` risk. The diagnostic features can also be filtered by their higher-level disease groupings and can be selected from the filter by column seen on the left of this figure.

Figure D.18: Here is the questions in context pane, which has a list dropdown option as seen on the top of this figure, where clinicians can browse through the question list we support for the patient being shown. These questions span the different question types we support, and the length of the question list is variable depending on how many diagnostic features contributed to the patient's predicted risk. Also seen here is the detail we support for each question and their answer, including provenance details like the confidence score and the data source for the predicted answer.