

Redefining Digital Health Interfaces with Large Language Models

Fergus Imrie^{1*}, Paulius Rauba² and Mihaela van der Schaar^{2,3}

¹University of California, Los Angeles, USA.

²University of Cambridge, UK.

³The Alan Turing Institute, UK.

*Corresponding author(s). E-mail(s): imrie@ucla.edu;
Contributing authors: pr501@cam.ac.uk; mv472@cam.ac.uk;

Abstract

Digital health tools have the potential to significantly improve the delivery of healthcare services. However, their adoption remains comparatively limited due, in part, to challenges surrounding usability and trust. Large Language Models (LLMs) have emerged as general-purpose models with the ability to process complex information and produce human-quality text, presenting a wealth of potential applications in healthcare. Directly applying LLMs in clinical settings is not straightforward, however, with LLMs susceptible to providing inconsistent or nonsensical answers. We demonstrate how LLM-based systems can utilize external tools and provide a novel interface between clinicians and digital technologies. This enhances the utility and practical impact of digital healthcare tools and AI models while addressing current issues with using LLMs in clinical settings such as hallucinations. We illustrate LLM-based interfaces with the example of cardiovascular disease risk prediction. We develop a new prognostic tool using automated machine learning and demonstrate how LLMs can provide a unique interface to both our model and existing risk scores, highlighting the benefit compared to traditional interfaces for digital tools.

Introduction

Digital healthcare technologies represent a frontier in medicine. Despite a multitude of tools being developed [1, 2], clinical adoption of such methods faces significant hurdles [3, 4], with some even calling their use “infeasible” [5] and “substantially conceptual” [6]. One key issue is usability [7], which can result in errors associated with patient harm [8] and contribute to clinician frustration, jeopardizing patient safety [9, 10]. New tools employing artificial intelligence (AI) and machine learning offer substantial promise, with their impact expected to be felt across all areas of healthcare [11]. Yet these approaches face the same usability challenges as existing digital tools, while presenting additional questions around model trust [12, 13]. Consequently, these issues have limited the uptake and impact of AI technologies in clinical settings thus far [14–16].

To improve the usability of clinical predictive models, several approaches have sought to simplify or automate the process of obtaining a prediction. These include points-based scoring systems [17], web-based calculators [18, 19], or integration within electronic health records [20]. While this can make such tools easier to use, simply obtaining a prediction is frequently insufficient and more is required to build model trust with both clinicians [12] and regulators [21, 22].

Large Language Models (LLMs) offer a potential solution to the challenges faced by digital tools. LLMs have recently captured the imagination of both the research community and the general public, pushing the boundaries of human-machine interaction. Consequently, there is great interest in applying LLMs in healthcare, with potential applications including facilitating clinical documentation, summarizing research papers, or as a chatbot for patients [23].

Applying LLMs in safety-critical clinical settings is not straightforward. LLMs may provide inconsistent or nonsensical answers [24, 25] and have a tendency to hallucinate facts [26, 27]. This is clearly unacceptable in medicine when making high-stakes decisions. Additionally, LLMs can encounter difficulty with seemingly basic functionality, such as mathematical calculations or factual lookup [28, 29], and are unable to access up-to-date information by default [30]. These limitations constrain the clinical utility of LLMs.

In this paper, we explore a new application of LLMs in healthcare and propose using LLMs to facilitate clinician interactions with AI models and digital tools. Conceptually, this differs substantially from previous applications, such as training medicine-specific LLMs [31, 32] or using LLMs for prediction [33]. Instead, we construct LLM-based systems, offering an intuitive natural language interface that can streamline clinician interactions with multiple digital tools and sources of information, improving efficiency and usability (Fig. 1).

Dynamic interactions in the form of natural language dialogues have been identified as a key feature for practitioners to deploy machine learning models in healthcare [34]. The most appropriate interface or dialogue cannot be pre-specified; instead, it depends on the clinician and the patient. Thus, the

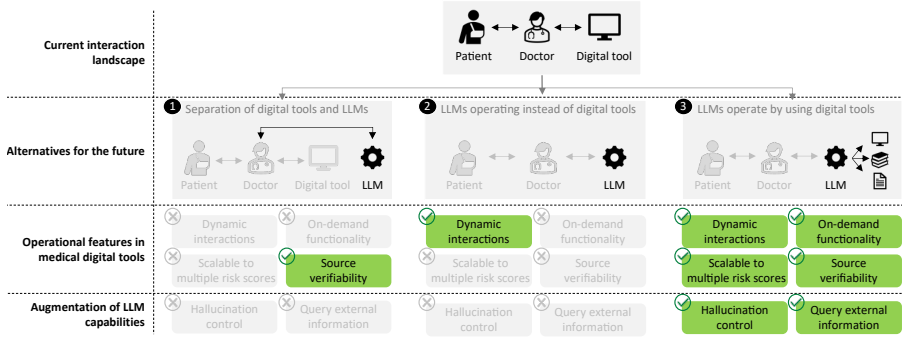


Fig. 1: Clinicians have previously needed to interact directly with digital tools, such as risk scores. While others have discussed LLMs replacing existing clinical tools (1,2), we envisage LLMs forming a novel interface by enabling dynamic interactions and facilitating deeper engagement with tools and related information, such as model explainability, medical papers, and clinical guidelines (3).

ability to adapt and tailor interactions represents a critical advance in the functionality of such tools that can be unlocked by LLMs.

By default, LLMs do not possess the ability to access external tools or information. We augment the base functionality of LLMs and enable them to access approved medical tools and other sources of information, thereby not solely relying on the inherent capabilities of a given LLM. This approach is scalable to multiple predictive models, unifying digital tools within a single, natural language-based interface. By adopting a systems approach, the LLM does not itself issue predictions and can access relevant domain-specific information, rather than needing to supply specific knowledge itself. Consequently, the potential for hallucinations is limited and we ensure actionable information is provided by approved clinical sources.

As an example of our approach, we specifically focused on risk scoring [35] and considered cardiovascular disease (CVD), the most common cause of mortality globally [36]. Primary prevention uses prognostic models to estimate the future risk of developing CVD. This allows high-risk individuals to be identified and their risk to be managed via interventions, such as lifestyle modifications or pharmaceuticals. Numerous CVD risk scores have previously been developed, for example the Framingham score [37] in the United States, SCORE2 [38] in Europe, and QRisk3 [18] in the UK. We begin by developing an improved risk score for CVD using AutoPrognosis 2.0 [19], an open-source automated machine learning framework for developing diagnostic and prognostic models. We then demonstrate a novel LLM-based interface for both our prognostic model and existing risk scores. We provide several illustrative examples of dynamic interactions that substantially extend the capabilities of existing fixed interfaces.

Table 1: Model validation on UK Biobank. C-indices, Brier scores, and Expected/Observed ratios for our approach (AP2), QRisk3, SCORE2, and Framingham score. We report results for the original existing risk scores and the performance following recalibration. Mean performance (95% CI).

	C-Index \uparrow	Brier Score \downarrow	E/O
AP2	0.741 (0.740-0.743)	0.041 (0.041-0.041)	1.003 (1.001-1.006)
<i>Original Risk Scores</i>			
Framingham	0.705 (0.703-0.707)	0.047 (0.047-0.047)	2.495 (2.491-2.499)
SCORE2	0.710 (0.709-0.712)	0.042 (0.042-0.042)	1.606 (1.603-1.608)
QRisk3	0.722 (0.721-0.724)	0.044 (0.043-0.044)	1.941 (1.937-1.944)
<i>Recalibrated Risk Scores</i>			
Framingham	0.705 (0.703-0.707)	0.041 (0.041-0.041)	1.000 (0.998-1.002)
SCORE2	0.710 (0.709-0.712)	0.041 (0.041-0.041)	1.000 (0.998-1.002)
QRisk3	0.722 (0.721-0.724)	0.041 (0.041-0.041)	1.000 (0.998-1.002)

Results

We conducted experiments using data from the UK Biobank [39], a prospective population study of around half a million individuals from the UK, enrolled between 2006 and 2010. UK Biobank collected a broad set of information from participants at enrollment, including health and medical history, blood tests, physical examination, and socio-demographics, with ongoing linkage to healthcare records. The outcome of interest was the incidence of CVD within a 10-year horizon, where CVD was defined as the composite endpoint of myocardial infarction (ICD-10 codes: I21, I22), angina pectoris (I20), stroke (I63, I64), or transient cerebral ischaemic attacks (G45). The cohort contained 407,605 individuals and there were 17,600 CVD events within a 10-year horizon. A flow chart of the study cohort is provided in Fig. S.2 and a summary of patient characteristics is provided in Table S.1. Further details can be found in Section 3.

We used AutoPrognosis 2.0 [19] to automatically construct a prognostic model for incidence of CVD. We first validate the performance of our approach compared to existing risk scores and then demonstrate a novel interface for using such risk scores using LLMs.

Model performance

In Table 1, we compared our prognostic model (AP2) with three existing CVD risk scores, namely QRisk3 [18], SCORE2 [38], and Framingham score [37]. The ensemble model constructed using AutoPrognosis achieved a C-index of 0.741 (95% CI: 0.740-0.743), a Brier score of 0.041 (95% CI: 0.041-0.041), and an expected/observed (E/O) ratio of 1.003 (95% CI: 1.001-1.006), outperforming all three existing risk scores assessed.

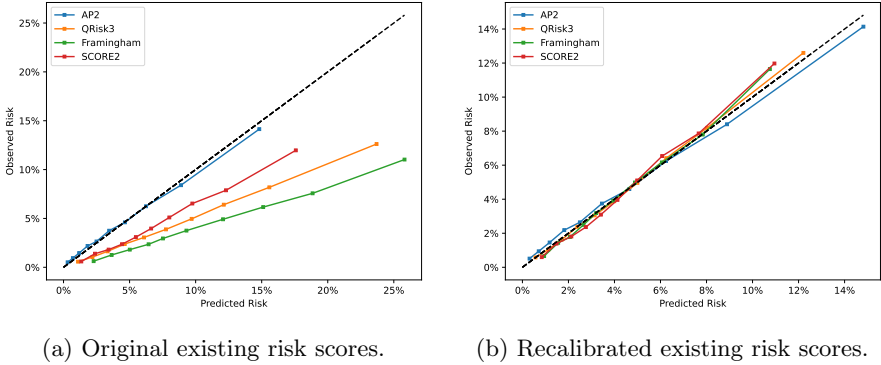


Fig. 2: Calibration curves. Calibration curves for our approach (AP2), QRisk3, SCORE2, and Framingham score, both before and after calibrating the existing risk scores to the UK Biobank cohort. Observed risk was calculated using Kaplan-Meier estimators [40].

QRisk3 and SCORE2 achieved C-indexes of 0.722 (95% CI: 0.721-0.724) and 0.710 (95% CI: 0.709-0.712), respectively, with similar Brier scores (QRisk3: 0.044, 95% CI: 0.043-0.044; SCORE2: 0.042, 95% CI: 0.042-0.42). Both approaches overpredicted risk in the UK Biobank population (E/O ratio QRisk3: 1.941, 95% CI: 1.937-1.944; SCORE2: 1.606, 95% CI 1.603-1.608) consistent with previous findings [41]. The Framingham score was the weakest performing risk score, achieving a C-index of 0.705 (95% CI: 0.703-0.707). This underperformance likely was a consequence of the Framingham score being developed on a substantially different population from the USA.

Calibration curves (Fig. 2) show that AutoPrognosis predictions were well calibrated across all risk deciles. In contrast, the existing risk scores over-predicted risk for all individuals in the UK Biobank cohort (Fig. 2a). We recalibrated the existing risk scores by scaling the predictions according to the overall incidence in the UK Biobank cohort and the mean predicted risk (Fig. 2b). Following recalibration, QRisk3 exhibited good calibration across all risk deciles, while SCORE2 and Framingham score both modestly overpredicted risk in the lower risk deciles and underpredicted risk in the highest risk decile.

To assess model performance in a manner that considers the effect on clinical decision making, we additionally performed decision curve analysis [42, 43], as recommended in the TRIPOD guidelines [44]. Decision curve analysis is a method to evaluate the clinical utility of a model by assessing the net benefit for a range of risk thresholds at which a decision maker would act or intervene. The model with the higher net benefit at any given threshold is preferred. We compared our model with the existing risk scores, as well as baseline strategies that assume all patients will develop CVD (All) or none will (None). Our model achieved a greater net benefit at all decision thresholds compared to the

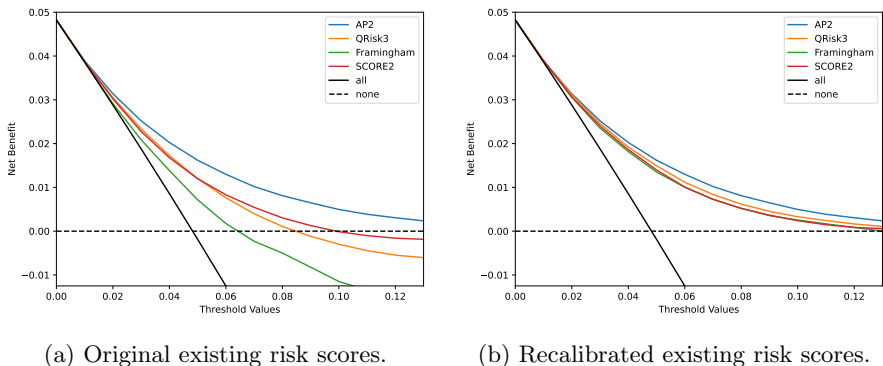


Fig. 3: Decision curve analysis. Our approach (AP2) provides greater net benefit at all thresholds than the existing risk scores (Framingham score, SCORE2, and QRisk3) and the baseline strategies (All and None).

existing risk scores (Fig. 3). Net benefit takes into account both discrimination and calibration. While recalibrating the existing risk scores improved their net benefit, the AutoPrognosis model still resulted in the highest net benefit at all thresholds.

To better understand the rationale for the predictions and drivers of the outperformance of our approach, we assessed feature importance by calculating SHAP values [45]. SHAP values for the 20 features included in the final version of our model are shown in Fig. 4. Consistent with clinical knowledge, age, sex, and blood pressure are three of the most important features and were included in all three existing risk scores considered in this work. However, a number of additional features, including both laboratory and nonlaboratory tests, also significantly contributed to model predictions. A number of these features have been shown to be risk factors for CVD, but have not been incorporated into other risk scores. For example, Apolipoprotein B (ApoB) is a primary component of several lipoproteins, including low-density lipoprotein (LDL), and transports cholesterol and other lipids in the blood. High levels of ApoB are the primary driver of atherosclerosis and there is evidence that ApoB is a more accurate predictor of CVD than total cholesterol or LDL [46, 47]. ApoB is deemed one of the most important variables in our approach (Fig. 4), along with Apolipoprotein A and Lipoprotein A which have also been shown to be risk factors for CVD [48]. Finally, “Overall health rating”, a non-laboratory, self-reported feature, was found to be predictive, mirroring the findings of Alaa et al. regarding the benefits of non-laboratory variables for CVD risk prediction [49].

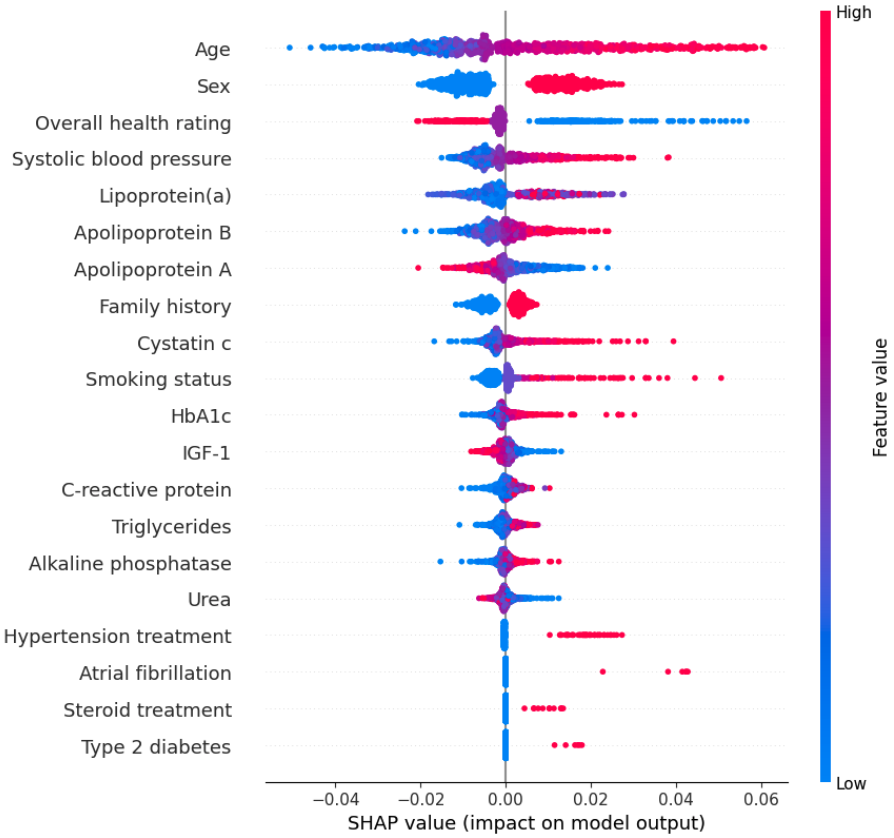


Fig. 4: Feature importance. SHAP values of the variables included in the AutoPrognosis model.

LLM-based Interfaces

While LLMs are powerful models for natural language processing, by default, they do not possess the ability to access external tools or information. Methods to extend the functionality of LLMs beyond text generation are in their infancy but can already be used to significantly expand the capabilities of LLMs [29, 50].

In this section, we construct LLM-based systems and demonstrate multiple examples of how such systems can provide a novel interface for digital health tools, in particular clinical risk scores. The LLM-based systems can incorporate numerous external tools, sources of information, and clinical data (illustrated in Fig. 5). Instead of using an LLM to issue predictions directly, we augmented LLMs and enabled them to access approved medical tools and other sources of information, thereby not solely relying on the inherent capabilities of a given LLM, while unifying multiple digital tools within a single, natural language-based interface.

Table 2: Representative questions that a clinician might have relating to a risk score at different stages of a patient encounter, together with whether existing interfaces for risk scores provide this information. All questions can be addressed using LLM-based interfaces.

Stage	Representative questions	Existing interfaces
Before Patient Encounter	Which features does the risk score use?	✗
	Why are these features included in the risk score?	✗
	How was the risk score validated?	✗
	What is the methodology underlying the risk score?	✗
Before Risk Scoring	When do clinical guidelines recommend risk scoring?	✗
	What is the recommended risk score?	✗
	Who is the risk score suitable for?	✗
During Risk Scoring	What is the risk for this patient?	✓
	What characteristics led to the patient’s risk?	✗
	What effect would changing this feature have on the risk?	✗
After Risk Scoring	What action is recommended by the guidelines based on the risk?	✗

To demonstrate the versatility of our approach, we took an off-the-shelf pretrained LLM, GPT-4 [51], and imbued it with a set of tools using in-context learning, requiring no additional training of the LLM. Additionally, we implemented a “source” functionality that allows the user to see whether the LLM used a tool or accessed specific information and, if so, which tool and with what input. This helps avoid hallucinations since it enables verification that the information was issued by an underlying clinical tool or source rather than being generated by the LLM. Additional implementation details are provided in Section 3.

AutoPrognosis. We constructed an LLM-based interface for the machine learning model developed in this work (Fig. 6). In our first example, we demonstrate the role such an interface could have upon model deployment (Fig. 6a). A clinician might have many questions about a risk score before using the predictions to inform recommendations or interventions (Table 2). Here, the LLM provided specific information about the variables used by the risk score, the underlying methodology, how the model was constructed, and the cohort of individuals used to train the model.

Building model trust is a crucial step for prognostic models, in particular for models that are not inherently interpretable [12]. A recent study found that medical decision-makers had a strong preference for interactive explanations and, in particular, for these interactions to take the form of natural language

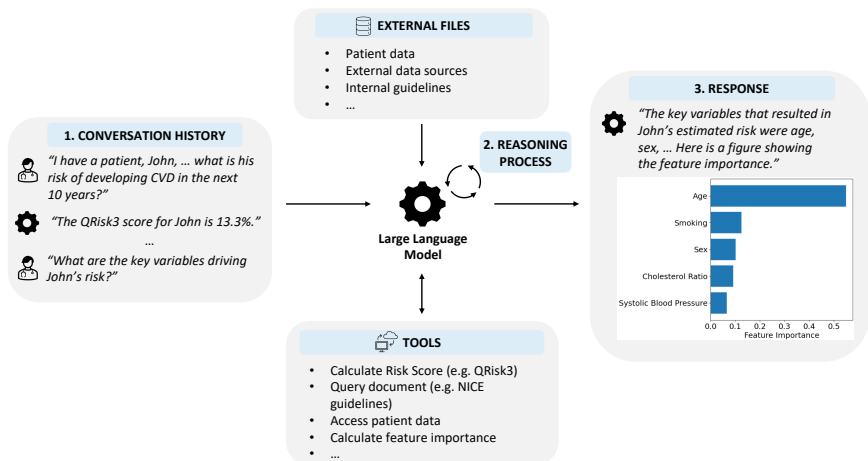


Fig. 5: Overview of an LLM-based system that enables clinicians to interface with digital tools using natural language inputs. (1) The LLM is provided with the history of the interaction, including the current request. (2) Using an iterative reasoning process, the LLM decides which, if any, tools are required and with what input. (3) The LLM provides a response to the user incorporating information provided by any tools that were used.

dialogues [34]. To enable this, in addition to allowing the LLM to access the model, we provided the LLM with the ability to apply methods from explainable AI (XAI) [52] to the underlying model, allowing users to investigate the rationale for predictions.

As shown in Fig. 6b, after calculating the patient’s risk using the Auto-Prognosis model, the clinician can query why this prediction was issued using XAI, facilitating a deeper understanding of the model. In this example, we used Shapley additive explanations (SHAP) [45] to explain the prediction. For this individual, their estimated 4.2% risk was primarily caused by their age and family history of CVD, mitigated by being a woman, not smoking, and low levels of lipoprotein (a). An additional interaction can be found in Fig. S.3. Due to their general pretraining, LLMs have knowledge of many topics; this can provide valuable additional information during interactions beyond the specific tools and external information sources provided to the LLM. For example, suppose the clinician is not familiar with the underlying XAI methodology, SHAP [45]. The LLM could explain how this approach works, in a variety of different ways and possibly over multiple interactions with the clinician, allowing specific queries or misunderstandings to be clarified. This demonstrates the utility of LLMs beyond simply using existing tools and the benefit of their underlying knowledge.

QRisk3. We also show how LLMs can incorporate existing tools and information for CVD risk prediction (Fig. 7). We provided the LLM access to

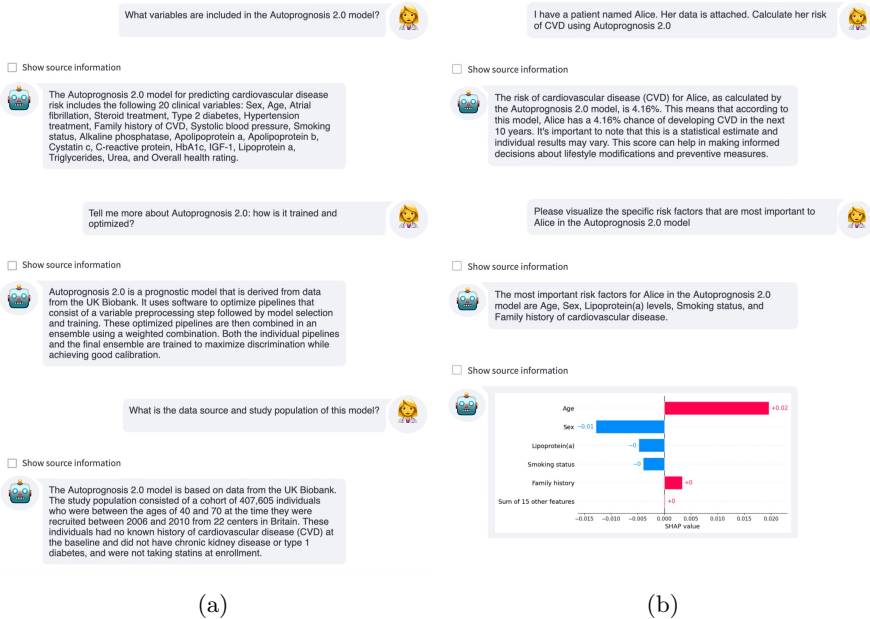




Fig. 6: Example interactions with LLM-based interfaces. (a) Clinicians can ask the LLM questions about the AutoPrognosis risk score, such as what variables are included, the underlying methodology, and the cohort on which the model was derived. (b) After calculating the patient’s risk using the AutoPrognosis model, the clinician can query why this prediction was issued using explainable AI to improve understanding of the model predictions.


QRisk3 [18], a risk prediction tool that assesses the likelihood of developing CVD within 10 years. Additionally, we provided the LLM with access to the academic paper describing QRisk3 [18] and the National Institute for Health and Care Excellence (NICE) clinical guidelines for CVD [53]. An example interaction between a physician and the LLM-based system is shown in Fig. 7. An illustration of the reasoning process by which the LLM uses external tools is provided in Fig. S.1.

In this example, when asked for the patient’s 10-year risk of developing CVD, the LLM used QRisk3 to estimate the patient’s risk, providing this to the user. The LLM then summarized the QRisk3 paper to explain the inclusion of certain features before providing the recommended action for this patient from the NICE clinical guidelines. Finally, the LLM used QRisk3 to recalculate the patient’s counterfactual risk assuming that they were able to reduce their systolic blood pressure to within normal ranges. This allows both the clinician and patient to understand the potential impact of changes to modifiable variables on the patient’s risk. While the clinician could have used


I have a patient named John who is a 56-years-old male with no family history of cardiovascular disease. His data is uploaded. Can you calculate his risk of developing cardiovascular disease in the next 10 years? 

Show source information


 The risk of your patient, John, developing cardiovascular disease in the next 10 years is 13.3%, as calculated by the QRisk3 model. This model is the recommended method for cardiovascular disease risk assessment in the UK and takes into account a wide range of variables.

I noticed the use of corticosteroids is a new variable in the QRisk3 calculation. Why is this included in QRisk3 and what are the implications of this for the predictive model? 


Show source information

 The inclusion of corticosteroid use in the QRisk3 risk score for cardiovascular disease is based on evidence that corticosteroids are associated with an increased risk of cardiovascular events. The QRisk3 study found that, in the study population, corticosteroid use was linked to an 82% increased risk of cardiovascular disease in women and a 58% increased risk in men. This finding is consistent with other research indicating that corticosteroids can have adverse effects on cardiovascular health.


The model extracted the information from a document titled Development and validation of QRISK3 risk prediction algorithms to estimate future risk of cardiovascular disease. Press [here](#) for the original document.

Thank you. Given John's risk, what action is recommended by the guidelines for addressing the CVD risk? 

Show source information

 In managing a patient with a 10-year predicted risk of 13% for cardiovascular disease, the primary recommendations focus on lifestyle modifications and lipid modification therapy. Patients are advised to adopt a healthy lifestyle, encompassing a cardioprotective diet with reduced saturated fats and increased mono-unsaturated fats, regular exercise, weight management, and cessation of smoking and alcohol consumption. Concurrently, if lifestyle adjustments prove insufficient, initiating statin therapy, particularly atorvastatin 20 mg, is advocated for those with a 10-year QRISK3 score of 10% or more. Also remember to communicate risks and provide regular monitoring.

The model extracted the information from a document titled Cardiovascular disease: risk assessment and reduction, including lipid modification. Press [here](#) for the original document.

I wonder if this is directly tied to John's systolic blood pressure. Can you please check what would be his risk if John's systolic blood pressure was 110 mm Hg? 

Show source information


 If John's systolic blood pressure was 110 mm Hg, his risk of developing cardiovascular disease in the next 10 years would decrease from 13.3% to 8.9%. This is a decrease of 4.4 percentage points.

Fig. 7: Example interaction with an LLM-based interface for QRisk3. The LLM uses QRisk3 to calculate the patient’s 10-year risk of developing CVD, explains why certain features are included in the risk score using the QRisk3 paper [18], and provides recommendations for this patient from the NICE clinical guidelines.

the underlying resources to obtain this information, the LLM interface made the interaction simpler and more efficient, which has been identified as a key limitation of digital tools [4, 7].

Discussion

Large Language Models hold substantial promise for the medical domain, particularly in augmenting digital workflows and improving efficiency in healthcare delivery. The ability to integrate external tools and functionality with LLMs paves the way for innovative applications and can overcome limitations of LLMs, such as hallucinations. Doing so offers a potential transformation for how clinicians interact with digital tools and could help overcome the challenges of deploying clinical AI models.

We have demonstrated how LLMs can provide a unique interface between healthcare professionals and clinical predictive models, such as risk scores. In particular, we developed a novel risk score for incidence of CVD using automated machine learning and developed LLM-based interfaces for our model and QRisk3 [18], the current recommended risk score in the UK for CVD. Our approach is scalable and does not require any additional training of the language model, although approaches that improve with use could be yet more powerful. Additionally, we aim to mitigate the problem of hallucination by ensuring that actionable advice is anchored in approved clinical resources, contrasting several previous applications of LLMs in medicine that focused exclusively on the knowledge and information learned by LLMs.

Numerous considerations exist before, during, and after assessing an individual's prognosis (Table 2). Currently, clinicians must access these tools via fixed user interfaces or application programming interfaces (APIs), with existing interfaces typically only calculating risk. Through an LLM-based interface, practitioners can obtain substantial additional information about the risk score, its development and methodology, the prediction issued, and related medical guidelines, in a manner that specifically addresses their needs or questions without providing superfluous information.

In this paper, we have focused on clinicians interacting with digital tools. However, there are numerous stakeholders in healthcare in addition to clinicians, such as patients, regulators, and administrators, each with different goals and requirements [52]. Our framework and approach could be applied to improve digital health interfaces for these alternate stakeholders. While this could have additional challenges, there are potentially even more substantial benefits for such individuals, given the differences in requirements, knowledge, and familiarity with digital health technology, among other factors.

While LLMs have general capabilities, they can lack domain-specific knowledge. This has led to the development of medical-focused LLMs, either by training new LLMs from scratch [31, 32, 54] or by adapting existing general-purpose LLMs [24]. While we showed using such specialist LLMs is not required, they could be readily incorporated due to the modularity of our approach. This could further enhance the functionality of LLM-based interfaces.

As AI in medicine continues to advance, further research into LLMs and their potential applications in healthcare could provide significant benefits.

For example, LLMs could help to alleviate the data burden that is contributing to clinician burnout, as well as streamline patient management processes. Furthermore, studies have demonstrated high usability of LLMs, even with limited experience [55], which is critical for successful clinical deployment. While we believe this paper represents an important first step, we are only scratching the surface of the potential of LLMs in healthcare. Ultimately, this line of work may significantly change the digital health landscape, enhancing the capabilities of clinicians and the quality of patient care.

Methods

Data source and study population

We developed and validated our prognostic model using data from UK Biobank [39], a large prospective cohort of individuals from the UK. Participants in UK Biobank were enrolled between 2006 and 2010 and aged between 40 and 70 at the time of recruitment. We extracted a cohort of individuals with no known history of CVD at baseline. We excluded individuals who had been diagnosed with chronic kidney disease or type 1 diabetes and individuals who were being treated with statins. This resulted in a cohort of 407,605 individuals (Fig. S.2).

Clinical Predictors and Outcomes

We began by considering a number of predictors previously validated as CVD risk factors. A subset of features were selected based on their predictive power and feature importance. We included 20 features in our model, namely sex, age, atrial fibrillation, steroid treatment, type 2 diabetes, hypertension treatment, family history of CVD, systolic blood pressure, smoking status (never-smoker, ex-smoker, light smoker, moderate smoker, heavy smoker, alkaline phosphatase, apolipoprotein a, apolipoprotein b, cystatin c, C-reactive protein, HbA1c, IGF-1, lipoprotein a, triglycerides, urea, and overall health rating, which was self-reported by participants as excellent, good, fair or poor.

The primary outcome was incidence of CVD within a 10-year horizon. This was defined using the following ICD-10 codes: myocardial infarction (I21, I22), stroke (I63, I64), angina pectoris (I20), or transient cerebral ischaemic attacks (G45). There were 17,600 CVD events within a 10-year horizon in the UK Biobank cohort.

Model Derivation

We trained a prognostic model using AutoPrognosis 2.0, an open-source automated machine learning software package [19]. AutoPrognosis has been validated in several applications in medicine, for example to determine eligibility for lung cancer screening [56]. AutoPrognosis was used to optimize pipelines consisting of a variable preprocessing step followed by model selection and training. The optimized pipelines were subsequently combined in an ensemble using a weighted combination. Individual pipelines and the final

ensemble were trained to maximize discrimination, measured using the area under the receiver operating curve. Predictions from the final ensemble were calibrated in a piecewise linear manner using the observed risk calculated using a Kaplan-Meier estimator [40].

Statistical Analysis

Model discrimination was assessed using the concordance index (C-index) [57] and calibration via calibration curves and the ratio between the expected and observed number of events. We also assessed Brier scores [58]. Calibration curves were constructed by first dividing the population into 10 risk deciles based on each individual’s predicted 10-year risk and then comparing the mean estimated risk to the observed risk in each decile. Kaplan–Meier estimators [40] were used to calculate observed risk. Additionally, we conducted decision curve analysis, assessing the net benefit for a range of risk thresholds. In addition to the existing risk scores, when assessing net benefit, we also compared our model with baseline strategies that assume all patients develop CVD (All) or none do (None). Results are based on 5-fold cross-validation. We used multiple imputation with chained equations (MICE) [59] to impute missing values. We generated five imputed datasets and combined results using Rubin’s rules [60].

Model Comparisons

We compared our model with three existing risk scores, namely Framingham score [37], SCORE2 [38], and QRisk3 [18]. All three models assess the likelihood of developing CVD within a 10 year horizon.

The Framingham score [37] was recommended by the 2010 American College of Cardiology/American Heart Association (ACC/AHA) guidelines [61] and is based on eight risk factors: sex, age, systolic blood pressure, treatment for hypertension, smoking status, history of diabetes, total cholesterol, and high-density lipoprotein (HDL) cholesterol.

SCORE2 [38] was developed using data from 45 cohorts in 13 countries across Europe, including the United Kingdom, and has been included in the European Society of Cardiology guidelines. SCORE2 employs seven risk factors: sex, age, systolic blood pressure, smoking status, history of diabetes, total cholesterol, and high-density lipoprotein (HDL) cholesterol.

QRisk3 [18] was developed on data from 981 practices in England comprising almost 8 million patients and is currently the recommended risk score for CVD in England [53]. QRisk3 uses 22 variables, including demographics (such as sex, age), lifestyle factors (e.g. smoking), medical and family history (e.g. history of diabetes, family history of CVD), medical examination (e.g. blood pressure), and blood tests (total cholesterol/HDL cholesterol ratio).

LLM-based interface

Here, we describe the implementation details for the LLM-based systems demonstrated above.

LLM framework

By default, LLMs will provide responses in the form of text based on the provided context, such as a prompt or conversation history. To construct interfaces for digital tools using LLMs, we instead viewed the LLM as an agent that can interact with an environment to solve tasks. Formally, at each step $t \in T$, the agent receives observation $o_t \in O$ from the environment and subsequently takes action $a_t \in A$ according to policy $\pi(a_t|h_t)$, where $h_t = (o_0, a_0, \dots, o_{t-1}, a_{t-1}, o_t)$ is the history. To enable the agent to both reason and use external tools, we used the ReAct framework [62] which decomposes the action space as $\hat{A} = A \cup L$, where $a \in A$ are actions using specific tools and an action $a \in L$ represents not using an external tool but instead allows the model to reason over the history about what action to take next.

Since we will provide the agents tasks in the form of natural language, and actions in the language space L are (essentially) infinite, we chose to benefit from strong language priors and use a pretrained LLM. We implemented our LLM-based interfaces using GPT-4 [51], with access through the OpenAI API. Interactions with external tools were implemented using LangChain [63].

Frameworks such as Toolformer [29] and WebGPT [50] trained LLMs to use basic tools, such as calculators, calendars, and search engines, via self-supervised fine-tuning and fine-tuning using behavior cloning and reinforcement learning, respectively. In contrast, following ReAct [62], we used in-context learning [64] in the form of prompting, providing the LLM with sufficient information about possible actions and using the underlying reasoning capabilities of the LLM. This removes the need for further training of the LLM, which might be challenging in the medical domain without suitable examples. This also enabled us to use multiple tools. We built user interfaces using StreamLit [65].

External Tools

This section describes the external tools and sources of information made available to the LLM in our examples.

LLM-based interface for AutoPrognosis. To demonstrate an LLM-based interface, we equipped the LLM with the CVD risk prediction model developed in this work on the UK Biobank cohort using AutoPrognosis 2.0 [19]. Building model trust is a crucial step for prognostic models, in particular for models that are not inherently interpretable [12, 13]. A recent study found that medical decision-makers had a strong preference for interactive explanations and, in particular, for these interactions to take the form of natural language dialogues [34]. To enable this, we enabled the LLM to use explainable AI (XAI) methods [52] on the underlying model, allowing users to investigate the rationale for predictions, both in general and for the specific patient. In particular, we used SHAP [45] to interpret model predictions. We additionally provided the LLM with information about the variables used by the risk score, the underlying methodology and how the model was constructed, and details

regarding the cohort used to develop the model. Example interactions are shown in Fig. 6 and Fig. S.3.

LLM-based interface for QRisk3. As a second example, we showed how LLMs can incorporate existing tools and information for CVD risk prediction. We provided the LLM access to QRisk3 [18] and enabled the LLM to use the risk score either using the provided data or, if requested by the user, to modify a variable and assess the impact of such a change on the patient’s risk. Additionally, we provided the LLM with access to the academic paper describing QRisk3 [18] and the National Institute for Health and Care Excellence (NICE) clinical guidelines for CVD [53]. An example interaction is shown in Fig. 7.

Code availability

Code for the LLM-based interfaces can be accessed at https://github.com/pauliusrauba/LLMs_interface. AutoPrognosis is an open-source package available on GitHub at <https://github.com/vanderschaarlab/AutoPrognosis> and on PyPI at <https://pypi.org/project/autoprognois/>.

Data availability

This research has been conducted using the UK Biobank resource under application number 105160. Data from UK Biobank is accessible through a request process (<https://www.ukbiobank.ac.uk/enable-your-research/register>). The authors had no special access or privileges when accessing the data.

Acknowledgments. This study received no funding.

Competing interests. All authors declare no financial or non-financial competing interests.

Author contributions. F.I. and M.vdS. conceptualized the manuscript. F.I. and P.R. designed and performed the experiments. F.I. wrote the original draft of the manuscript, and all authors contributed to editing and revising it.

References

- [1] Sutton, R. T. *et al.* An overview of clinical decision support systems: benefits, risks, and strategies for success. *npj Digit. Med.* **3** (1), 17 (2020) .
- [2] Dunn, J., Runge, R. & Snyder, M. Wearables and the medical revolution. *Per. Med.* **15** (5), 429–448 (2018) .
- [3] Eichler, K., Zoller, M., Tschudi, P. & Steurer, J. Barriers to apply cardiovascular prediction rules in primary care: A postal survey. *BMC Fam. Pract.* **8**, 1–7 (2007) .
- [4] Mathews, S. C. *et al.* Digital health: A path to validation. *npj Digit. Med.* **2** (1), 38 (2019). <https://doi.org/10.1038/s41746-019-0111-3> .

- [5] Müller-Riemenschneider, F. *et al.* Barriers to Routine Risk-Score Use for Healthy Primary Care Patients: Survey and Qualitative Study. *Arch. Intern. Med.* **170** (8), 719–724 (2010). <https://doi.org/10.1001/archinternmed.2010.66> .
- [6] Abernethy, A. *et al.* The promise of digital health: Then, now, and the future. *NAM perspect.* (2022) .
- [7] Ratwani, R. M., Reider, J. & Singh, H. A Decade of Health Information Technology Usability Challenges and the Path Forward. *JAMA* **321** (8), 743–744 (2019). <https://doi.org/10.1001/jama.2019.0161> .
- [8] Howe, J. L., Adams, K. T., Hettinger, A. Z. & Ratwani, R. M. Electronic Health Record Usability Issues and Potential Contribution to Patient Harm. *JAMA* **319** (12), 1276–1278 (2018). <https://doi.org/10.1001/jama.2018.1171> .
- [9] Shanafelt, T. D. *et al.* Relationship between clerical burden and characteristics of the electronic environment with physician burnout and professional satisfaction. *Mayo Clin. Proc.* **91** (7), 836–848 (2016). <https://doi.org/10.1016/j.mayocp.2016.05.007> .
- [10] Gardner, R. L. *et al.* Physician stress and burnout: The impact of health information technology. *J. Am. Med. Inform. Assoc.* **26** (2), 106–114 (2019) .
- [11] Bajwa, J., Munir, U., Nori, A. & Williams, B. Artificial intelligence in health-care: Transforming the practice of medicine. *Future Health. J.* **8** (2), e188–e194 (2021). <https://doi.org/10.7861/fhj.2021-0095> .
- [12] Rajpurkar, P., Chen, O., Emmaand Banerjee & Topol, E. J. AI in health and medicine. *Nat. Med.* **28** (1), 31–38 (2022). <https://doi.org/10.1038/s41591-021-01614-0> .
- [13] Asan, O., Bayrak, A. E. & Choudhury, A. Artificial intelligence and human trust in healthcare: focus on clinicians. *J. Med. Internet Res.* **22** (6), e15154 (2020) .
- [14] Goldfarb, A. & Teodoridis, F. Why is AI adoption in health care lagging? (2022). URL <https://www.brookings.edu/articles/why-is-ai-adoption-in-health-care-lagging/>.
- [15] Davenport, T. & Kalakota, R. The potential for artificial intelligence in health-care. *Future Health. J.* **6** (2), 94–98 (2019). <https://doi.org/10.7861/futurehosp.6-2-94> .
- [16] Kelly, C. J., Karthikesalingam, A., Suleyman, M., Corrado, G. & King, D. Key challenges for delivering clinical impact with artificial intelligence. *BMC Med.* **17** (1), 195 (2019). <https://doi.org/10.1186/s12916-019-1426-2> .
- [17] Gage, B. F. *et al.* Validation of clinical classification schemes for predicting stroke results from the national registry of atrial fibrillation. *JAMA* **285** (22), 2864–2870 (2001). <https://doi.org/10.1001/jama.285.22.2864> .

- [18] Hippisley-Cox, J., Coupland, C. & Brindle, P. Development and validation of QRISK3 risk prediction algorithms to estimate future risk of cardiovascular disease: Prospective cohort study. *BMJ* **357** (2017). <https://doi.org/10.1136/bmj.j2099> .
- [19] Imrie, F., Cebere, B., McKinney, E. F. & van der Schaar, M. AutoPrognosis 2.0: Democratizing diagnostic and prognostic modeling in healthcare with automated machine learning. *PLoS Digit. Health* **2** (6), 1–21 (2023). <https://doi.org/10.1371/journal.pdig.0000276> .
- [20] Rothman, M. J., Rothman, S. I. & Beals, J. Development and validation of a continuous measure of patient condition using the electronic medical record. *J. Biomed. Inform.* **46** (5), 837–848 (2013). <https://doi.org/10.1016/j.jbi.2013.06.011> .
- [21] Food and Drug Administration and others. Proposed regulatory framework for modifications to artificial intelligence/machine learning (AI/ML)-based software as a medical device (SaMD) (2019) .
- [22] Mourby, M., Ó Cathaoir, K. & Collin, C. B. Transparency of machine-learning in healthcare: The GDPR & European health law. *Comput. Law Secur. Rev.* **43**, 105611 (2021) .
- [23] Moor, M. *et al.* Foundation models for generalist medical artificial intelligence. *Nature* **616** (7956), 259–265 (2023). <https://doi.org/10.1038/s41586-023-05881-4> .
- [24] Singhal, K. *et al.* Large language models encode clinical knowledge. *Nature* (2023). <https://doi.org/10.1038/s41586-023-06291-2> .
- [25] Lecler, A., Duron, L. & Soyer, P. Revolutionizing radiology with GPT-based models: Current applications, future possibilities and limitations of ChatGPT. *Diagn. Interv. Imaging* **104** (6), 269–274 (2023). <https://doi.org/10.1016/j.diii.2023.02.003> .
- [26] Maynez, J., Narayan, S., Bohnet, B. & McDonald, R. On faithfulness and factuality in abstractive summarization. *In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* 1906–1919 (2020). <https://doi.org/10.18653/v1/2020.acl-main.173> .
- [27] Ji, Z. *et al.* Survey of hallucination in natural language generation. *ACM Comput. Surv.* **55** (12) (2023). <https://doi.org/10.1145/3571730> .
- [28] Patel, A., Bhattamishra, S. & Goyal, N. Are NLP models really able to solve simple math word problems? *In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* 2080–2094 (2021). <https://doi.org/10.18653/v1/2021.naacl-main.168> .
- [29] Schick, T. *et al.* Toolformer: Language models can teach themselves to use tools. *Adv. Neural Inf. Process. Syst.* **36** (2023) .

- [30] Komeili, M., Shuster, K. & Weston, J. Internet-augmented dialogue generation. *In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* **1**, 8460–8478 (2022) .
- [31] Luo, R. *et al.* BioGPT: Generative pre-trained transformer for biomedical text generation and mining. *Brief. Bioinformatics* **23** (6), bbac409 (2022). <https://doi.org/10.1093/bib/bbac409> .
- [32] Yang, X. *et al.* A large language model for electronic health records. *npj Digit. Med.* **5** (1), 194 (2022). <https://doi.org/10.1038/s41746-022-00742-2> .
- [33] Jiang, L. Y. *et al.* Health system-scale language models are all-purpose prediction engines. *Nature* (2023). <https://doi.org/10.1038/s41586-023-06160-y> .
- [34] Lakkaraju, H., Slack, D., Chen, Y., Tan, C. & Singh, S. Rethinking explainability as a dialogue: A practitioner’s perspective. *arXiv preprint arXiv:2202.01875* (2022) .
- [35] Moons, K. G., Royston, P., Vergouwe, Y., Grobbee, D. E. & Altman, D. G. Prognosis and prognostic research: what, why, and how? *BMJ* **338**, b375 (2009). <https://doi.org/10.1136/bmj.b375> .
- [36] Muthiah, V., A., M. G., Varieur, T. J., Valentin, F. & A., R. G. The global burden of cardiovascular diseases and risk. *J. Am. Coll. Cardiol.* **80** (25), 2361–2371 (2022). <https://doi.org/10.1016/j.jacc.2022.11.005> .
- [37] D’Agostino, R. B. *et al.* General cardiovascular risk profile for use in primary care. *Circulation* **117** (6), 743–753 (2008). <https://doi.org/10.1161/CIRCULATIONAHA.107.699579> .
- [38] SCORE2 working group & ESC Cardiovascular risk collaboration. SCORE2 risk prediction algorithms: new models to estimate 10-year risk of cardiovascular disease in Europe. *Eur. Heart J.* **42** (25), 2439–2454 (2021). <https://doi.org/10.1093/eurheartj/ehab309> .
- [39] Sudlow, C. *et al.* UK Biobank: An open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLOS Med.* **12** (3), 1–10 (2015). <https://doi.org/10.1371/journal.pmed.1001779> .
- [40] Kaplan, E. L. & Meier, P. Nonparametric estimation from incomplete observations. *J. Amer. Statist. Assoc.* **53** (282), 457–481 (1958) .
- [41] Parsons, R. E. *et al.* Independent external validation of the QRISK3 cardiovascular disease risk prediction model using UK Biobank. *Heart* **109** (22), 1690–1697 (2023). <https://doi.org/10.1136/heartjnl-2022-321231> .
- [42] Vickers, A. J. & Elkin, E. B. Decision curve analysis: A novel method for evaluating prediction models. *Med. Decis. Making* **26** (6), 565–574 (2006). <https://doi.org/10.1177/0272989X06295361> .

- [43] Vickers, A. J. Decision analysis for the evaluation of diagnostic tests, prediction models, and molecular markers. *Am. Stat.* **62** (4), 314–320 (2008). <https://doi.org/10.1198/000313008X370302> .
- [44] Moons, K. G. *et al.* Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): Explanation and elaboration. *Ann. Intern. Med.* **162** (1), W1–W73 (2015). <https://doi.org/10.7326/M14-0698> .
- [45] Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. *Adv. Neural Inf. Process. Syst.* **30**, 4765–4774 (2017) .
- [46] Marston, N. A. *et al.* Association of apolipoprotein B–containing lipoproteins and risk of myocardial infarction in individuals with and without atherosclerosis: Distinguishing between particle concentration, type, and content. *JAMA Cardiol.* **7** (3), 250–256 (2022). <https://doi.org/10.1001/jamacardio.2021.5083> .
- [47] Behbodikhah, J. *et al.* Apolipoprotein b and cardiovascular disease: Biomarker and potential therapeutic target. *Metabolites* **11** (10), 690 (2021). <https://doi.org/10.3390/metabo11100690> .
- [48] Erqou, S. *et al.* Apolipoprotein(a) isoforms and the risk of vascular disease: Systematic review of 40 studies involving 58,000 participants. *J. Am. Coll. Cardiol.* **55** (19), 2160–2167 (2010). <https://doi.org/10.1016/j.jacc.2009.10.080> .
- [49] Alaa, A. M., Bolton, T., Di Angelantonio, E., Rudd, J. H. F. & van der Schaar, M. Cardiovascular disease risk prediction using automated machine learning: A prospective study of 423,604 UK Biobank participants. *PLoS One* **14** (5), 1–17 (2019). <https://doi.org/10.1371/journal.pone.0213653> .
- [50] Nakano, R. *et al.* WebGPT: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332* (2021) .
- [51] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774* (2023) .
- [52] Imrie, F., Davis, R. & van der Schaar, M. Multiple stakeholders drive diverse interpretability requirements for machine learning in healthcare. *Nat. Mach. Intell.* **5** (8), 824–829 (2023). <https://doi.org/10.1038/s42256-023-00698-2> .
- [53] National Institute for Health and Care Excellence. *Cardiovascular disease: risk assessment and reduction, including lipid modification* (2014). URL <https://www.nice.org.uk/guidance/cg181>. Last updated: 24 May 2023.
- [54] Taylor, R. *et al.* Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085* (2022) .
- [55] Skjuve, M., Følstad, A. & Brandtzaeg, P. B. The user experience of ChatGPT: Findings from a questionnaire study of early users. In: *Proceedings of the 5th International Conference on Conversational User Interfaces* (2023). <https://doi.org/10.1145/3571884.3597144> .

- [56] Callender, T. *et al.* Assessing eligibility for lung cancer screening using parsimonious ensemble machine learning models: A development and validation study. *PLOS Medicine* **20** (10), e1004287 (2023). <https://doi.org/10.1371/journal.pmed.1004287> .
- [57] Uno, H., Cai, T., Pencina, M. J., D’Agostino, R. B. & Wei, L. J. On the C-statistics for evaluating overall adequacy of risk prediction procedures with censored survival data. *Stat. Med.* **30** (10), 1105–1117 (2011). <https://doi.org/10.1002/sim.4154> .
- [58] Brier, G. W. Verification of forecasts expressed in terms of probability. *Mon. Weather Rev.* **78** (1), 1–3 (1950) .
- [59] van Buuren, S. & Groothuis-Oudshoorn, K. mice: Multivariate imputation by chained equations in R. *J. Stat. Softw.* **45** (3), 1–67 (2011). <https://doi.org/10.18637/jss.v045.i03> .
- [60] Rubin, D. B. Multiple imputation. *Flexible Imputation of Missing Data, Second Edition* 29–62 (2018) .
- [61] Greenland, P. *et al.* 2010 ACCF/AHA guideline for assessment of cardiovascular risk in asymptomatic adults: A report of the American College of Cardiology Foundation/American Heart Association Task Force on practice guidelines. *J. Am. Coll. Cardiol.* **56** (25), e50–e103 (2010). <https://doi.org/10.1016/j.jacc.2010.09.001> .
- [62] Yao, S. *et al.* ReAct: Synergizing reasoning and acting in language models. In: *The Eleventh International Conference on Learning Representations* (2023) .
- [63] Chase, H. Langchain. <https://github.com/hwchase17/langchain> (2022). Last accessed: 2023-06-23.
- [64] Dong, Q. *et al.* A survey on in-context learning. *arXiv preprint arXiv:2301.00234* (2023) .
- [65] Streamlit. <https://streamlit.io/>.

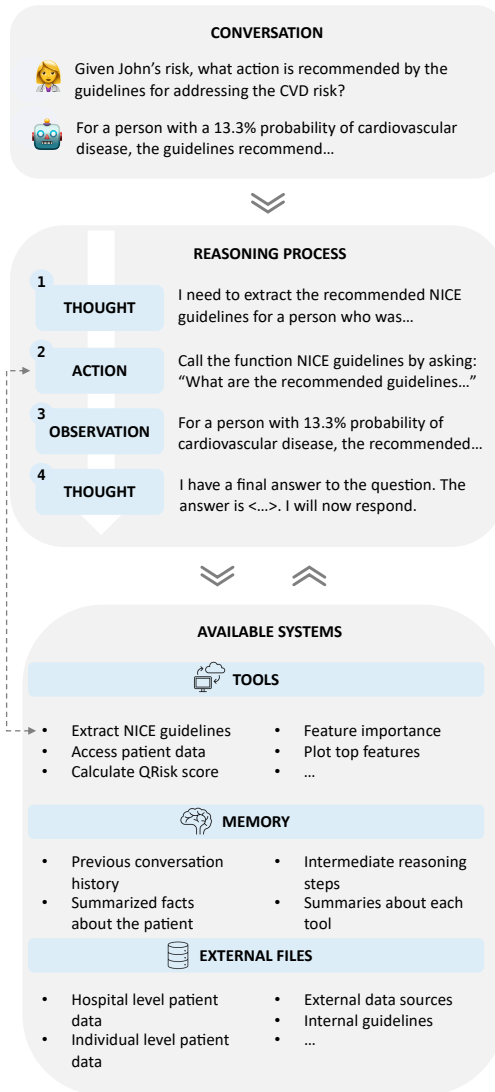


Fig. S.1: Illustration of the process by which the LLM uses external tools. The LLM is provided with the history of the interaction, including the current request. Using an iterative reasoning process, the LLM decides which, if any, tools are required and with what input (“Thought”). The LLM then uses the external tool (“Action”) and receives the output (“Observation”). Finally, the LLM decides to answer the question (“Thought”), providing a response to the user.

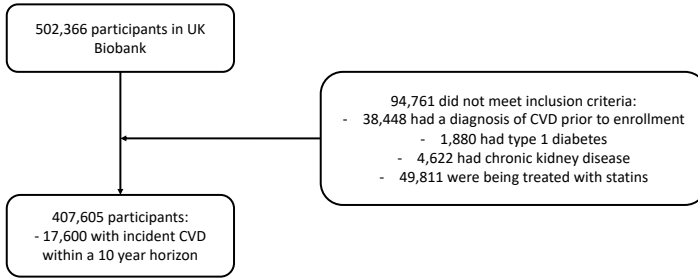




Fig. S.2: Flow diagram of UK Biobank participants.


Table S.1: Descriptive characteristics of the UK Biobank cohort.

	All (n=407,605)	CVD events (n=17,600, 4.32%)	No CVD event (n=390,005, 95.68%)
Age			
<50	110,750 (27.17%)	1,876 (10.66%)	108,874 (27.92%)
50-60	142,784 (35.03%)	5,272 (29.95%)	137,512 (35.26%)
≥60	154,071 (37.80%)	10,452 (59.39%)	143,619 (36.82%)
<i>Missing</i>	0 (0.00%)	0 (0.00%)	0 (0.00%)
BMI			
<18.5	2,403 (0.59%)	85 (0.48%)	2,318 (0.59%)
18.5-25	144,179 (35.37%)	4,514 (25.65%)	139,665 (35.81%)
25-30	171,140 (41.99%)	7,978 (45.33%)	163,162 (41.84%)
30-35	87,529 (21.47%)	4,880 (27.73%)	82,649 (21.19%)
<i>Missing</i>	2,354 (0.58%)	143 (0.81%)	2,211 (0.57%)
Systolic Blood Pressure			
<120	89,725 (22.01%)	2,027 (11.52%)	87,698 (22.49%)
120-140	164,263 (40.30%)	6,157 (34.98%)	158,106 (40.54%)
>140	128,984 (31.64%)	8,209 (46.64%)	120,775 (30.97%)
<i>Missing</i>	24,633 (6.04%)	1,207 (6.86%)	23,426 (6.01%)
Total Cholesterol / HDL			
<4.5	224,495 (55.08%)	7,282 (41.38%)	217,213 (55.69%)
4.5-6	98,739 (24.22%)	5,747 (32.65%)	92,992 (23.84%)
>6	25,215 (6.19%)	2,075 (11.79%)	23,140 (5.93%)
<i>Missing</i>	59,156 (14.51%)	2,496 (14.18%)	56,660 (14.53%)
Sex			
Men	173,818 (42.64%)	10,360 (58.86%)	163,458 (41.91%)
Women	233,787 (57.36%)	7,240 (41.14%)	226,547 (58.09%)
<i>Missing</i>	0 (0.00%)	0 (0.00%)	0 (0.00%)
Smoking			
Never	238,095 (58.41%)	8,608 (48.91%)	229,487 (58.84%)
Previous	125,829 (30.87%)	6,202 (35.24%)	119,627 (30.67%)
Current	41,420 (10.16%)	2,646 (15.03%)	38,774 (9.94%)
<i>Missing</i>	2,261 (0.55%)	144 (0.82%)	2,117 (0.54%)


I have a patient named John. His data is attached. Calculate his risk of CVD using AutoPrognosis 2.0 

Show source information

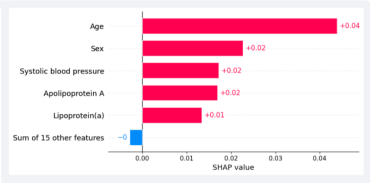
 The risk of cardiovascular disease (CVD) for your patient, John, as calculated by the AutoPrognosis 2.0 model, is 15.14%. This is considered a high risk. This means that John has a 15.14% chance of developing CVD in the next 10 years. It's important to consider this information in the context of John's overall health and lifestyle. High-risk individuals may need to consider lifestyle modifications, medications, or other interventions to reduce their risk. Please consult with a healthcare professional for the best course of action.

Please visualize the specific risk factors that are most important to John in the AutoPrognosis 2.0 model 

Show source information

 The most important risk factors for John in the AutoPrognosis 2.0 model are Age, Sex, Systolic blood pressure, Apolipoprotein A, and Lipoprotein(a).

Show source information



Risk Factor	SHAP Value
Age	+0.04
Sex	+0.02
Systolic blood pressure	+0.02
Apolipoprotein A	+0.02
Lipoprotein(a)	+0.01
Sum of 15 other features	-0.01

Fig. S.3: Example interaction with an LLM-based interface. After calculating the patient’s risk using the AutoPrognosis model, the clinician can query why this prediction was issued using explainable AI.