# Human-AI Co-design for Clinical Prediction Models

Jean Feng[1,*,†], Avni Kothari[1,*], Patrick Vossler[1], Andrew Bishara[1], Lucas Zier[1], Newton Addo[1], Aaron Kornblith[1], Yan Shuo Tan[2], Chandan Singh[3]

[1]University of California, San Francisco.
[2]National University of Singapore.
[3]Microsoft Research.

[*]Equal contribution.
[†]Corresponding author: jean.feng@ucsf.edu.

### Abstract

Developing safe, effective, and practically useful clinical prediction models (CPMs) traditionally requires iterative collaboration between clinical experts, data scientists, and informaticists. This process refines the often small but critical details of the model building process, such as which features/patients to include and how clinical categories should be defined. However, this traditional collaboration process is extremely time- and resource-intensive, resulting in only a small fraction of CPMs reaching clinical practice. This challenge intensifies when teams attempt to reliably incorporate information from unstructured clinical notes, which can contain an essentially infinite number of concepts. To address this challenge, we introduce HACHI, an iterative human-in-the-loop framework that uses AI agents to accelerate the development of fully interpretable CPMs by enabling the exploration of concepts in clinical notes. HACHI alternates between (i) an AI agent rapidly exploring and evaluating candidate concepts in clinical notes and (ii) clinical and domain experts providing feedback to improve the CPM learning process. HACHI defines concepts as simple yes-no questions that are used in linear models, allowing the clinical AI team to transparently review, refine, and validate the CPM learned in each round. In two real-world prediction tasks (acute kidney injury and traumatic brain injury), HACHI outperforms existing approaches, surfaces new clinically relevant concepts not included in commonly-used CPMs, and improves model generalizability across clinical sites and time periods. Furthermore, HACHI reveals the critical role of the clinical AI team, such as directing the AI agent to explore entire categories of concepts that it had not previously considered, adjusting the granularity of concepts it considers, changing the objective function to better align with the clinical objectives, and identifying issues of data bias and leakage. Code for HACHI is available at http://github.com/jjfenglab/HACHI.

**Keywords:** Large language models, Electronic health records, Concept Bottleneck, Human-AI Interaction

## 1 Introduction

Clinical prediction models (CPMs) translate routinely collected clinical information into structured predictions, enabling clinicians to apply shared expertise in a more consistent manner across patients. While there is growing interest in highly complex black-box CPMs, simple fully-interpretable CPMs, such as rule-based scores like SOFA for ICU mortality [1] or PECARN and NEXUS for traumatic brain injury [2], remain the most popular in clinical practice due to a number of major advantages. First and foremost, these models can be easily implemented by clinicians at the bedside using simple mental arithmetic or a quick chart evaluation, which facilitates their integration into the clinical workflow [3, 4]. Furthermore, fully-interpretable models can be thoroughly audited by domain experts and, when necessary, modified to improve reliability, which promotes user trust [5, 6]. Consequently,
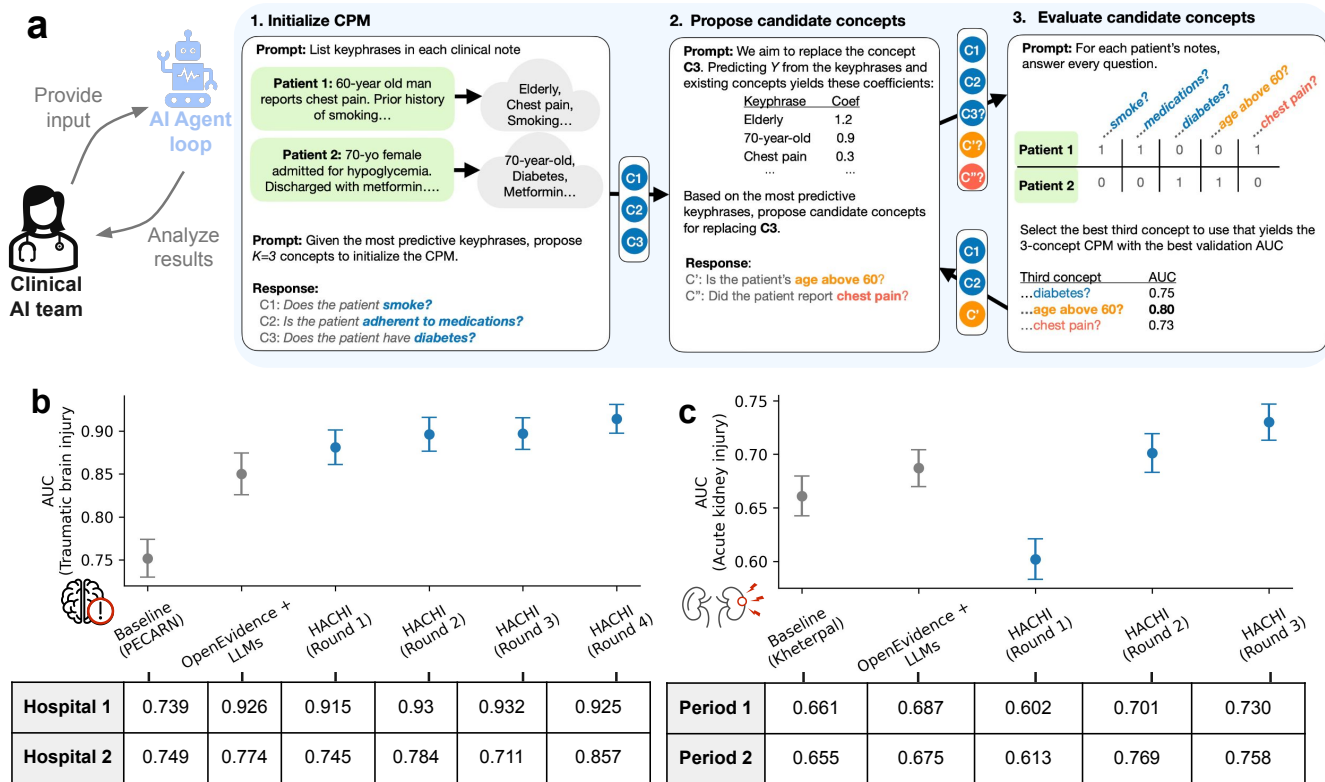
**Fig. 1**: **The HACHI framework uses LLMs with humans-in-the-loop to build an effective clinical prediction model (CPM).** (a) HACHI is composed of an outer loop where the clinical AI team provides guidance and feedback to the AI agent on how to learn a CPM, and an inner loop where the AI agent follows instructions from the clinical AI team to find $k$ concepts (formally defined as yes/no questions) that maximize the CPM's predictive accuracy. The inner AI-guided CPM learning procedure is broken down into three steps: 1. Initialize the CPM by brainstorming clinical concepts from keyphrases extracted out of clinical notes; 2. Propose candidate concepts by analyzing which keyphrases are most associated with the outcome of interest; 3. Evaluate the candidate concepts by annotating each concept and selecting the best-performing concept(s). Steps 2 and 3 are repeated until convergence. Each round, the clinical AI team analyzes the results from the AI-guided CPM learning procedure and provides feedback on how the procedure can be improved, such as by modifying the prompts and clarifying which concepts are and are not of interest. HACHI improves over baselines and across rounds for two real-world clinical prediction tasks: (b) diagnosis of traumatic brain injury (TBI), where performance is evaluated in terms of AUC with respect to the overall population (shown in plot) and stratified across two sites (shown in table). (c) development of acute kidney injury (AKI), where performance is evaluated in an internal validation set (Period 1, shown in plot and table) and a later, temporally disjoint test dataset (Period 2, shown in table). Error bars show standard errors.

when simple, fully-interpretable CPMs both meet clinical standards and achieve high performance, these tools are widely used by the clinical community [7–10].

Despite their utility, very few such CPMs have achieved the level of performance and clinical sophistication needed for widespread clinical adoption [11–13]. The main challenge is that "the devil is in the details"—a CPM is determined by the *numerous* details and decisions made in the CPM learning process, including which learning procedure was used, how the data was generated and compiled, which features are included, and how the model was evaluated [14–16]. All of these factors must be carefully chosen, which generally requires a lengthy and laborious collaboration between clinicians and data scientists to iterate on the model learning procedure so that the final CPM both satisfies complex clinical requirements and meets desired levels of performance [1, 17]. To make matters worse, maximizing model performance often requires extending beyond tabular data in the Electronic Health

Record (EHR) and analyzing unstructured clinical notes, as the latter often contains the most relevant and predictive features [18–20]. However, clinical notes cover an essentially *infinite* number of concepts; for example, the concept of *smoking* can be characterized by whether a patient currently smokes, has a history of smoking, is trying to quit smoking, and many more. As such, it is infeasible for clinical AI teams to manually explore all the possible concepts for inclusion in the CPM learning process.

Recent developments in artificial intelligence (AI) suggest that AI agents can efficiently navigate and optimize many of these critical design decisions for learning fully-interpretable CPMs. Indeed, recent research has shown that because large language models (LLMs) can now mimic complex clinical reasoning and accurately extract concepts from clinical notes [21–23], AI agents can now analyze clinical notes, brainstorm useful features to include in a CPM, convert these features into simple tabular elements, and apply statistical tools to fit models [24–28]. To find concepts in clinical notes that are most predictive of the target of interest, these methods have an AI agent iteratively explore the infinite space of concepts based on its clinical knowledge, which can reveal novel predictive signals that are difficult to identify upfront [28, 29]. However, existing works only study the AI agent working in isolation. Without external guidance from domain experts, the AI agent may learn concepts that are overly simplistic or lack clinical credibility, leading to CPMs that are promising but lack the clinical sophistication necessary for translation into clinical practice.

This work studies how clinical AI teams can collaborate with AI agents to improve the CPM learning process. This cannot be achieved simply through one-time prompt tuning because of the so-called "gulf of envisionment" [30, 31]: In the initial stage of building a CPM, the clinical AI teams often lack alignment and clarity on the exact task specifications, so it is often unclear how to best instruct an AI agent to achieve the team's goals. To close this "gulf of envisionment," the team must clarify to the AI agent what specific data to analyze, how exactly the outcome should be defined, what concepts to explore, what statistical tools to use, and more. This can only be accomplished through an iterative, collaborative process, in which the human experts gradually understand the capabilities of the AI agent, gain clarity in how the task should be specified, and learn how to align the AI agent to learn the best possible CPM. Through this process, team members can also build trust in the AI agent's learning process and, ultimately, the resulting CPM.

To this end, we introduce a <u>H</u>uman+<u>A</u>gent <u>C</u>o-design framework for <u>H</u>ealthcare <u>I</u>nstruments (HACHI)[1]. In HACHI, the role of the AI agent is to convert unstructured clinical data into structured, reviewable concepts for learning the best-performing CPM, per instructions from the clinical AI team; the role of the clinical AI team is to analyze the results from each round using their high-level design and clinical reasoning to repeatedly refine the instructions to the AI agent. As such, the AI agent does much of the "heavy-lifting"—e.g., analyzing clinical notes, brainstorming candidate concepts, fitting statistical models, and selecting which concepts to include—while the clinical AI team provides overall, directional feedback on how the AI agent could improve its process (Fig. 1). HACHI allows the clinical AI team to fully collaborate with the AI agent because every step of HACHI is interpretable. We apply HACHI to two real-world clinical risk-prediction tasks—acute kidney injury and traumatic brain injury—to evaluate the efficacy of this human+AI agent co-design process.

## 2 Results

The goal of HACHI is to learn a CPM with $k$ concepts that are interpretable, clinically reasonable, and maximize predictive performance, given a dataset of clinical notes and corresponding labels. Concepts are defined as answers to yes/no questions, such as "Does a patient have...?". Because clinical AI teams often do not have a precise definition of this high-level goal at the beginning of a project, HACHI involves multiple rounds of fitting CPMs, each time with an AI-powered CPM learning procedure that

---

[1]HACHI references Hachikō (1923-1935), an Akita dog commemorated in Japan for his strong loyalty to his owner. The name highlights how the process of clinical AI teams learning to work with an AI agent is similar to the iterative process of owners learning to train their dogs.

is precisely specified but potentially imperfect. The clinical AI team gathers new insights from each round, which are then used to improve how the AI-powered CPM learning procedure is run. After multiple rounds, the learning procedure becomes better-aligned with the clinical AI team's goals, resulting in a CPM with the desired characteristics.

The AI-guided CPM learning procedure used in this work (Fig. 1 blue box) is a greedy hill-climbing procedure based on [28] that iteratively refines a $k$-concept CPM by cyclically traversing each concept position and replacing the incumbent concept when a better-performing concept is found. The AI agent iteratively brainstorms, evaluates, and selects CPM concepts per the following steps:

1. **Initialize**: The AI agent creates a set of keyphrases to represent each observation, per instructions given in `KeyphrasePrompt`. The CPM is initialized with $k$ concepts per an `InitializationPrompt`, and the data is partitioned into a training set and a validation set.
2. **Propose candidate concepts**: Using the training data, the AI agent fits a statistical model to identify the top keyphrases associated with the outcome. Combining this list with its prior knowledge, the AI agent proposes candidate concepts in the form of yes/no questions for inclusion in the CPM, per instructions given in `ProposalPrompt`.
3. **Evaluate candidate concepts**: For each observation, the AI agent reviews the clinical notes and extracts the value of each candidate concept per instructions given in `ExtractionPrompt`. For each candidate concept, the AI agent fits a CPM on the training data and evaluates its performance on the validation data. The candidate concept with the best-performing CPM is selected.
4. **Iterate**: Repeat steps 2-4 until convergence or until a maximum number of iterations is reached.

The clinical AI team provides feedback to the AI agent by refining the hyperparameters of this procedure. The team's key lever is through modification of free-text prompts given as instructions to the AI agent, which has the advantage of being readable by any human, even those with limited AI expertise; a default set of prompts is provided to initialize the process. Nevertheless, the team has access to other levers as well, such as modifying which observations are included in the data, developing and/or adding new statistical tools for the AI agent to use, and deciding how samples are weighted. To help teams come up with feedback for each round, a Protected Health Information (PHI)-compliant web interface is provided so that all team members can review the AI agent's results (Fig. 2). Through this review process, the clinical AI team can decide how many rounds of HACHI to run; the following case studies found 3-4 rounds to be sufficient.

Below, we describe how clinical AI teams leveraged HACHI to learn CPMs for two prediction tasks at UCSF: whether a child presenting to the Emergency Department after head trauma will be diagnosed with traumatic brain injury (TBI) and whether an adult preparing for general surgery will develop acute kidney injury (AKI) post-surgery.

## 2.1 Case Study 1: Traumatic Brain Injury (TBI)

TBI is a leading cause of morbidity and mortality in children presenting to the emergency department (ED) following head trauma [32, 33]. As such, it is critical to identify which children require computed tomography (CT) imaging to avoid missing TBI cases; on the other hand, ionizing radiation from CT scans can cause lethal malignancies [34, 35]. To guide CT imaging decisions in pediatric head trauma, the Pediatric Emergency Care Applied Research Network (PECARN) rule is currently the most widely used CPM for identifying patients at high risk of TBI [2]. PECARN incorporates clinical signs and symptoms such as altered mental status, loss of consciousness, severe mechanism of injury, and physical examination findings (e.g., scalp hematoma, skull fractures). While the PECARN rule is designed to achieve very high sensitivity to minimize missed TBIs, its specificity remains relatively low, potentially exposing patients to unnecessary harmful radiation [36, 37].

**Table 1**: **Description of HACHI rounds for TBI.** Purple indicates features that raised concerns. Abbreviations — LOC: Loss of Conciousness, GCS: Glasgow Coma Scale.

| Round | Key Feedback | Changes Made | Learned Concepts |
|---|---|---|---|
| 1 | Initial exploration with default prompts | None (baseline) | LOC (2.70), Brain bleed (1.74), Neurological event (1.58), Note mentions GCS (1.23), Seizure-free (1.10) |
| 2 | (1) Some concepts reflect note-writing style instead of patient characteristics; (2) Brain bleed concept suggest data leakage, investigation revealed patients with existing CT results from prior facilities | (1) Require concept questions to have prefix "Does the note mention the patient having..."; (2) Removed cases with existing TBI diagnosis or mentioned CT results | LOC (3.46), Normal GCS & age $\geq$2yr (1.50), Convulsions (1.04), Altered mental status (0.82), Intact cranial nerves ($-0.65$) |
| 3 | Some coefficients contradict clinical intuition | Modified greedy concept selection to require sign of estimated coefficient to match LLM's clinical prior | LOC (4.34), History of mild TBI (1.41), Occipital hematoma (1.21), Vision changes (1.07), Memory intact ($-0.53$) |
| 4 | Model shows poor generalizability across Oakland and Mission Bay campuses due to imbalanced representation | Introduced sample weights to weight the two campuses equally | LOC (1.58), Altered mental status (0.73), Headache (0.25), Head trauma (0.08), Normal gait ($-0.22$) |

To determine whether a better-performing CPM for TBI could be learned, a clinical AI team was assembled, including one pediatric ED clinician (A. Kornblith), an emergency clinical data analyst (N.A.), and three data scientists (J.F., C.S., A. Kothari). Given that PECARN includes only a small number of concepts, the team decided to learn a 5-concept CPM. The team collated a retrospective case-control dataset with 400 cases and 400 controls from all encounters with documented ICD codes for head trauma or traumatic brain injury (TBI) diagnoses at two academic pediatric EDs within the Benioff Children's Hospital (Oakland and Mission Bay Campuses) between March 1, 2014, and December 31, 2024. ED Triage Notes, ED Provider Notes, and Nursing Notes were included in the analysis. To minimize the chance of data leakage, only the History & Physical section of ED provider notes were included and, for the other note types, only those with timestamps prior to the encounter's CT scan were included if a CT scan was performed. Four rounds of HACHI were completed, with Table 1 showing the learned concepts, feedback, and modifications from each round. The full set of prompts used in each round are provided in the Supplementary Materials.

**Round 1 → 2: Data leakage and spurious correlations.** `Round 1` was run with the default template of prompts. To the surprise of the clinical AI team, this CPM already achieved a high AUC of 0.92. While this performance could be partly explained by the CPM's overlap in concepts with PECARN (e.g., whether the patient had experienced loss of consciousness (LOC) and having a neurological event), the team also noted two major issues that could be leading to inflated performance. First, the AI agent had learned the concept "whether a note mentions Glasgow Coma Scale (GCS)," which meant that the AI agent had learned to rely on spurious correlations between the note writing style and TBI diagnosis. Second, the AI agent identified that having a "brain bleed" substantially increased the risk of TBI, but brain bleeds (i.e., intracranial hemorrhages) are typically known only *after* a patient has CT scan results. After investigating this issue of data leakage, the team found that a significant proportion of patients were transferred from another ED with an existing diagnosis of TBI. So for `Round 2`, (1) the AI agent was instructed to define concept questions using the prefix "Does the note mention the patient having...", to ensure that concepts extract attributes of the patient and not note writing style, and (2) all patients who already had prior CT scan results were removed from the dataset, resulting in 304 remaining cases.

**Round 2 → 3: Directional alignment with clinical prior.** `Round` 2's CPM had a slightly lower AUC of 0.90, which was expected since Round 1's AUC was inflated. While all the learned concepts were clinically relevant, the new concern raised by the clinical AI team was that some of the learned coefficients did not match clinical intuition. For instance, having a normal GCS had a positive coefficient (i.e., increased risk), contrary to clinical expectations. After discussion with the clinical AI team, the data scientists suggested modifying the greedy concept selection process in `Round` 3: rather than only selecting the candidate concept based on AUC, selection additionally required the coefficient's sign to match the LLM's clinical prior.

**Round 3 → 4: Model generalizability and fairness.** With the added changes, `Round` 3's CPM simultaneously had clinically reasonable concepts, directionally aligned coefficients, and maintained the same AUC of 0.9. As such, the clinical AI team decided to run one last round, where the team assessed the generalizability of the model across the two Oakland and Mission Bay campuses. This revealed a surprising gap in performance: the AUCs at the Oakland and Mission Bay campuses were 0.93 and 0.71, respectively. Investigating further, the team noted that the ratio of patients from Oakland to Mission Bay in the dataset was 3:1 and that the top feature learned by the procedure, LOC, was highly discriminative for patients at the Oakland campus but not at Mission Bay. This was clinically plausible, as the Oakland campus functions as a safety-net, pediatric Level 1 trauma referral center that sees higher-acuity and more selectively referred children, whereas the Mission Bay campus serves as a quaternary care hospital that sees a broader population. So for `Round` 4, the team decided to update the learning procedure to allow for sample weights, so that the two campuses could be equally weighted. `Round` 4's CPM not only achieved a better overall AUC of 0.91, but also better campus-specific AUCs of 0.93 and 0.80 at Oakland and Mission Bay, respectively. The concepts in this final model are LOC, altered mental status, headache, head trauma, and having a normal gait. The concept "head trauma" was initially surprising, but analyzing the LLM's annotations revealed that in real-world scenarios, there are a small percentage of cases (5%) where there is a suspicion of unwitnessed head trauma but no reliable patient history is available (e.g., abandoned child, unclear mechanism of injury, no chief complaint on file).

**Comparator methods.** As comparison, we evaluated the performance of PECARN by extracting relevant attributes from the same clinical notes. In addition, we simulated the common approach to learning CPMs where a clinical expert outputs a list of potential risk factors and protective factors for inclusion in a CPM, but does not further refine this list. To simulate this one-time brainstorming of clinical factors, we collated features that were originally considered for inclusion to the PECARN model [38] as well as features brainstormed by OpenEvidence [39]. Nevertheless, the AUCs of these models were 0.75 and 0.88, respectively, which were substantially worse than that achieved by HACHI. Finally, because PECARN was designed to optimize sensitivity, we compare the specificities of the CPMs from the different methods holding sensitivity constant in Table A2 of the Appendix, where we again find that the final model from HACHI outperforms the comparator methods.

## 2.2 Case Study 2: Acute Kidney Injury (AKI)

AKI is a common but serious postoperative complication, occurring in 2–25% of patients. When it occurs, AKI increases associated mortality and costs by two- to fivefold [40–42]. Early identification of high-risk patients enables preventive measures such as optimizing hemodynamics, avoiding nephrotoxic medications, and closer monitoring. While there are multiple well-known AKI prediction models for cardiac surgeries, there are far fewer for predicting AKI risk for lower-risk general surgeries. The most well-known interpretable AKI prediction model for General Surgery is the Kheterpal model [43], whose predictors include patient characteristics, medications, and surgery type. Nevertheless, use of this model is limited in clinical practice, since its performance is highly dependent on the patient case mix.

**Table 2**: Summary of HACHI rounds for AKI. Blue indicates concepts learned in response to feedback from the clinical AI team. CKD: Chronic Kidney Disease.

| Round | Key Feedback | Changes Made | Learned Concepts |
|---|---|---|---|
| **1** | Initial exploration with default prompts focused on patient characteristics | None (baseline) | Leukocytosis (1.31), Abdominal distention (0.87), Cardiac dysfunction (0.74), Tachycardia (0.71), Swelling (0.65), Systemic infection (0.57), Renal impairment (0.53), Diabetes mellitus (0.47), No kidney disease (−0.36), Low hematocrit (−1.01) |
| **2** | (1) Prompts only considered patient characteristics, but surgery type is often more influential for AKI risk; (2) Learned concepts are too broad and should be more specific; (3) clinician brainstormed additional factors such as infection and need for blood transfusion | (1) Updated all prompts to include both patient and surgical risk factors; (2) Updated prompts to encourage more specific concepts, such as focusing on higher-risk factors; (3) Added clinician-suggested concepts as in-context learning examples in prompts | CKD (1.28), Exploratory laparotomy (1.03), Sepsis (0.98), Fluid retention (0.89), Respiratory condition (0.47), Need blood transfusion during surgery (0.43), High-risk surgery (0.40), Heart disease (0.37), Malignancy (0.31), Minimally invasive surgery (−1.01) |
| **3** | (1) Concepts too vaguely defined, which may lead to inconsistent extraction between clinicians; (2) Medication concepts missing | (1) Required concept questions to include precise definitions with examples; (2) Added medications to brainstorming prompts | CKD (1.05), Major/high-risk surgery e.g., major abdominal surgery or anticipated blood loss > 500mL (0.91), Tachycardia >100 bpm or palpitations (0.88), Active systemic infection e.g., sepsis, bacteremia, or requiring therapeutic antibiotics (0.72), Heart failure (0.64), Sleep apnea (0.51), Urgent/emergent surgery (0.50), Obesity (0.39), Hypertension (0.31), Minimally invasive surgery e.g., laparoscopic or robotic-assisted procedures (−0.74) |

To determine if a better-performing CPM for AKI could be learned, a clinical AI team was assembled with an anesthesiologist (A.B.), a clinical data analyst (J.Y.), and two data scientists (J.F. and A. Kothari). A retrospective case-control dataset was selected uniformly at random from all General Surgery patients between January 2016 and March 2024, with 800 cases and 800 controls. The outcome is defined as whether AKI develops within 7 days following surgery per the Kidney Disease Improving Global Outcomes (KDIGO) criteria (stage 1 or higher) [44]. To capture preoperative patient characteristics, the dataset consists of preoperative anesthesia notes. Given that the Kheterpal model [43] used 11 clinical features, the team decided that a 10-concept CPM would be appropriate. Three rounds of HACHI were completed, as shown in Table 2. The full set of prompts used in each round are provided in the Supplementary Materials.

**Round 1 → 2: Missing concept category.** Round 1 was run with the default template of prompts for learning patient characteristics associated with the prediction target. This CPM achieved an AUC of 0.60 and identified many factors known to be associated with AKI that had been previously identified in [43], such as renal impairment and diabetes mellitus. The method also identified features not in [43], including hematocrit levels, leukocytosis, and cardiac dysfunction. However, the clinical AI team noted that the default prompts had directed the AI agent to only consider patient characteristics, when surgery type is often even more influential to a patient's AKI risk. So for Round 2, all the LLM prompts were updated to instruct the AI agent to consider both patient and surgical risk factors. In addition, the clinician suggested potential risk factors that could be provided to the AI agent as in-context learning examples, such as a patient's functional capacity, surgical duration, and whether the patient was coming from the ICU.

**Round 2 → 3: Increased precision.** With the inclusion of surgical risk factors, `Round 2`'s CPM had a substantially higher AUC of 0.70. Some of these factors are already known (e.g., chronic kidney disease, fluid retention) and others are known but not included in commonly-used CPMs (e.g., minimally invasive surgery, active infection requiring antibiotics, requiring blood transfusion during surgery). While the performance increase was promising, the clinical AI team discussed potential hurdles when implementing this CPM in practice. The major concern was that some of the CPM concepts were too vaguely defined, so that concept extractions for the same patient may differ between clinicians, reducing reliability in practice. So for `Round 3`, the AI agent was instructed to define concept questions more precisely using the prefix "Does the note mention the patient having ..., e.g.,....?" and likewise for surgery. For instance, a valid question could be "Does the note mention the patient having good exercise tolerance, e.g., having at least 5 METS of functional capacity or ability to climb two flights of stairs?" The clinician also suspected that home medications may also be predictive of AKI, so the AI agent was also encouraged to consider inclusion of medication keyphrases and concepts. The CPM in `Round 3`, the final round, achieved the highest AUC of 0.73 on the internal validation set (Period 1), suggesting that the increased precision reduced variability of the extracted concepts and thereby increased performance. Critically, the concepts now had construct validity since they were clearly defined. To assess temporal generalizability, the clinical AI team assembled a second test dataset from the disjoint time period of April-December 2024 (Period 2). The AUCs in Period 2 similarly improved from Round 1 in HACHI, starting from 0.61 in Round 1 to 0.77 and 0.76 in Rounds 2 and 3, respectively.

**Comparator methods.** As comparison, we evaluated the unweighted and weighted Kheterpal models [43] in two different ways: using LLM extractions from the preoperative clinical notes and using tabular data from the EHR. The models using LLM extractions achieved AUCs of 0.65-0.66 and those using tabular data achieved AUCs of 0.64. In addition, we simulated a clinical expert who brainstormed predictors for AKI risk without an iterative refinement process by asking OpenEvidence [39] to generate a list of potential predictors. Using the 30 concepts proposed through this approach, the resulting model achieved an AUC of 0.70. On the temporally disjoint test dataset, the models from [43] and the single-round brainstorming approach had similarly low AUCs of 0.67.

## 3  Discussion

This work demonstrates that meaningful collaboration between clinical AI teams and AI agents can produce high-performing CPMs that are simple to use, fully interpretable, and clinically credible. Across two clinical prediction tasks—traumatic brain injury in children and acute kidney injury in surgical patients—we found that the HACHI co-design framework consistently improved model performance, clinical relevance, and generalizability compared to existing clinical decision instruments and approaches that learn CPMs in a zero-shot fashion. Compared to the traditional process of training a CPM (which includes time spent on manual feature engineering, chart review, and iterative model refinement), HACHI constitutes a more efficient use of clinical AI teams: the teams in both case studies spent approximately 1-2 hours reviewing results and providing feedback per round, with 3-4 rounds sufficient to reach a satisfactory model.

These case studies highlight how human oversight is essential not only for validating AI outputs but also for shaping the model development process in ways that the AI agent could not achieve in isolation. In the TBI case study, the clinical AI team identified critical issues that would have gone undetected in a purely automated pipeline: data leakage from transferred patients with prior CT results, spurious associations between documentation style and outcomes, and performance disparities across clinical sites with different patient populations and care pathways. These insights led to modifications of the AI-guided CPM learning procedure not only in terms of the prompts given (e.g., constraining concept definitions to patient characteristics rather than note features and asking the

AI agent whether the concepts are expected to be risk factors versus protective factors) but also the data analyzed and how the model was evaluated (e.g., removing transferred patients and reweighting the training/validation data). Similarly, in the AKI case study, the AI agent initially focused exclusively on patient characteristics and learned concepts that were too vague for clinicians to implement in practice. The clinical AI team had to prompt the AI agent to additionally consider other concepts such as surgical factors and had to redefine how concept questions were formulated. Both case studies demonstrate how the clinical AI team must iteratively work with the AI agent to overcome the gulf of envisionment.

Through this collaborative process between the AI agent and clinical AI team, HACHI was able to uncover predictors commonly omitted from CPMs, while recovering commonly used predictors. In the TBI case study, the final model learned by HACHI is a much simpler form of the PECARN model, using only a subset of features such as LOC and altered mental status, alongside the more novel feature of "whether the patient has a normal gait." In the AKI case study, the final model identified multiple factors that have not been commonly included in interpretable AKI risk scores: minimally invasive surgery, tachycardia, and sleep apnea. The more novel factors uncovered in this study may be useful more generally, although additional validation is necessary.

HACHI is designed to be simple to use, in light of recent concerns that a global CPM is unlikely to work well for all patient populations, given the complex nature of medicine [45]. Instead, these works suggest training local CPMs for local populations. To address this concern, HACHI has minimal requirements. It only needs access to (i) clinical notes, which are generally easy to extract from the EHR, (ii) a PHI-compliant LLM (e.g., API endpoints in the cloud), and (iii) guidance from a clinical AI team. Furthermore, the deployment of HACHI learned CPMs is straightforward. Rather than requiring LLMs in clinical practice, clinicians can instead annotate the questions during a patient encounter. The code for running HACHI is open-source and ready-to-use (see Code Availability section). This allows teams to run HACHI for local patient populations, where it can discover a parsimonious set of concepts that are highly predictive locally.

HACHI builds on recent work using LLMs for concept bottleneck models [24–27, 46] and iterative concept refinement [27–29, 47]. The most closely related method is BC-LLM [28], which uses an LLM-in-the-loop rather than humans to iteratively propose and evaluate concepts. HACHI also relates to broader frameworks that combine LLMs with automatic verification loops [48–50], that incorporate human feedback to improve machine learning models [51–55] or that improve human-model interaction for difficult tasks [56–58]. Distinct from these prior works, HACHI integrates human domain expertise throughout the model development process, enabling clinical AI teams to identify data quality issues, enforce clinical consistency, and ensure practical usability—aspects that purely automated approaches cannot adequately address. While we use lasso-penalized logistic regression as the underlying statistical model here, HACHI can easily be customized to use models that better fit an application, e.g. sparse integer linear models [59], bayesian rule lists [60], concept decision trees [61–63], and other interpretable models [5, 64].

Several challenges remain before HACHI-derived CPMs can be deployed in clinical practice. Most critically, models developed retrospectively must undergo prospective validation to ensure they perform as expected when used in real-time clinical workflows. Nevertheless, the simplicity and interpretability of HACHI-derived CPMs may facilitate more rapid prospective validation compared to complex black-box models. Additionally, while LLM-based concept extraction is scalable, it still requires careful prompt engineering and quality control. In our studies, we found that the iterative review process naturally led to more precise concept definitions—the clinical AI team's feedback helped refine vague concepts like "protective factors" into precise questions with concrete examples, which substantially improved extraction consistency. Finally, as with any prediction model, careful attention to fairness, bias, and equity is essential, particularly when models are trained on data from specific institutions or populations that may not reflect the broader patient population.

Several limitations warrant consideration. Most importantly, the models in this study were developed and evaluated retrospectively using historical clinical data. While we assessed temporal generalizability using held-out time periods, retrospective evaluation may not fully capture how models will perform when deployed prospectively in real-time clinical workflows, so prospective validation is still needed. In addition, though these case studies involved multiple sites, the analyses and models were conducted at a single academic medical center with specific documentation practices and patient populations. The team members applying HACHI are also from this same institution, and it is to be determined if a different clinical AI team would have achieved similar performance and learned similar concepts. Finally, the quality of HACHI-derived models depends on the LLM's ability to accurately extract concepts from clinical notes and the types of clinical notes that are included. While prior work has found that GPT models are highly accurate at extracting concepts from medical texts [21, 65], LLMs differ in their extraction accuracy and, consequently, the CPM's performance may also differ.

Avenues for future research include extending HACHI's applicability to more complex clinical scenarios by incorporating multimodal data such as imaging, laboratory results, or time series data, as has been done in prior CBM works [66, 67]. In addition, other frameworks may be explored for incorporating human feedback into the co-design process, such as ways to reduce the burden on clinical experts while maintaining quality. Finally, although HACHI was tested in two clinical examples, the framework is applicable more broadly, particularly to fields beyond healthcare that also have large amounts of unstructured data. For this broader framework, perhaps a more suitable name is HACHII, "Human+Agent Co-design for Highly Interpretable Instruments."

## 4 Methods

### 4.1 AI agentic CPM learning procedure

The AI agent is an encapsulated large language model (LLM) managed entirely through text-based prompts and outputs. The agent conducts a greedy hill-climbing procedure to select $k$ concepts that are most predictive of the prediction target, subject to constraints specified by the clinical AI team. The procedure can be broken into the three main steps highlighted in Fig. 1a, beginning with the initialization in Step 1.

**Step 1: Initialize the CPM.** To brainstorm possible concepts, HACHI iterates through the clinical notes in the dataset, appending each note to a `KeyphrasePrompt`. This prompt is used to instruct the LLM to extract a list of 'keyphrases' or 'keywords' that represent the content of the note. The data is then split into a training and validation partition. Applying a statistical tool (e.g. ridge-penalized logistic regression), HACHI fits a bag-of-words model on the training partition to identify keyphrases that are most associated with the outcome of interest. Next, to initialize the concepts, HACHI appends the top keyphrases to a `ConceptInitializationProposalPrompt`. This prompt instructs the LLM to make use of these keyphrases, together with its world knowledge, to propose $k$ candidate "concept questions". Concept questions are defined as human-interpretable yes/no questions, e.g. "Does this note mention the patient having a history of smoking?" Concept questions can encompass multiple keyphrases (i.e., the concepts may be hierarchical).

To convert the extracted concepts into a CPM, the $k$ candidate "concept questions" are first extracted from the LLM's response. For each question, HACHI iterates through the clinical notes, appending both the question and the note to a `ConceptAnnotationPrompt`, which instructs the LLM to give a yes/no response to the question posed regarding the note. Finally, HACHI fits a $k$-concept CPM on the training partition.

Now that the CPM has been initialized, HACHI cyclically traverses the concept positions $j = 1, 2, \ldots, k$, evaluating each current concept for potential replacement by a more predictive concept. For each concept position, it conducts steps 2 and 3 as follows:

**Step 2: Propose candidate concepts.** The AI agent fits a bag-of-words model on the training partition to identify keyphrases that are most associated with the outcome of interest $Y$, but this time adjusting for existing concepts that are not up for replacement. HACHI appends the top keyphrases to a `ConceptReplacementProposalPrompt`, which instructs the LLM to make use of these keyphrases, together with its world knowledge, to propose $m$ candidate replacement "concept questions". The $m$ candidate "concept questions" are extracted from the LLM's response.

**Step 3: Evaluate candidate concepts.** For each question, the AI agent iterates through the clinical notes, appending both the question and the note to a `ConceptAnnotationPrompt`, which instructs the LLM to give a yes/no response to the question posed regarding the note. The resulting annotations are then used by the AI agent to fit a $k$-concept CPM (again using lasso-penalized logistic regression) on the training partition for each candidate concept by combining the candidate with the existing $k-1$ concepts. Applying a statistical validation metric (e.g. AUC), HACHI evaluates the resulting CPM's performance on the validation partition. If the best-performing candidate concept outperforms the existing concept being considered for replacement, HACHI replaces it with this best-performing candidate.

**Implementation details.** In our experiments, we performed 10 iterations of the AI agent loop consisting of Steps 2 and 3. In ablation studies, we did not find significant improvement in the validation AUC with more iterations.

The AI-guided CPM learning procedure is stochastic, due to both sample splitting and inclusion of the AI agent. Furthermore, there is inherent uncertainty regarding the true concepts that are relevant for predicting the outcome of interest, due to finite sample sizes and the infinite number of possible concepts. To address both concerns, we run the AI-guided CPM learning procedure for multiple seeds in parallel. Differences between the learned CPMs can reveal uncertainty in the truly relevant concepts and stability of the learning process, enabling more efficient Human-AI interaction.

## 4.2 Human-AI interaction framework

**Interface for reviewing results from the AI agent.** To help the human expert determine how the prompts should be modified each round, we created an interactive PHI-compliant single-page locally-hosted webpage to show: (i) the learned factors, (ii) the LLM concept annotations at a patient/note level, (iii) which patients received an incorrect prediction, and (iv) the performance of the CPM on a held-out test set; see Fig. 2 for more details. A meeting between the data scientists and clinical experts is also held to review the learned CPMs and brainstorm ways to improve it. Data scientists then translate feedback from this meeting into actual LLM prompts and/or code updates.

**Human-in-the-loop control and prompt steering.** The clinical AI team exerts control over the AI agent primarily by modifying the four prompts: `KeyphrasePrompt`, `ConceptInitializationProposalPrompt`, `ConceptReplacementProposalPrompt`, and `ConceptAnnotationPrompt`. However, the team may also opt for other modifications as well, such as the composition of the dataset, modifications to the statistical tools used by the AI agent (e.g. the extension of the tool to allow for sample weights), and modifications to the code for how candidate concepts are selected in the optimization process. More complex modifications may require code-based changes to the AI-guided CPM learning procedure, rather than simply modifying the inputs to the procedure.

## 4.3 Dataset extended details

**Data collection for traumatic brain injury (TBI).** To ensure a comprehensive capture of relevant cases, we included any encounter with documented head trauma or a traumatic brain injury (TBI) diagnosis, as coded by the following ICD-10 (International Classification of Diseases, Tenth
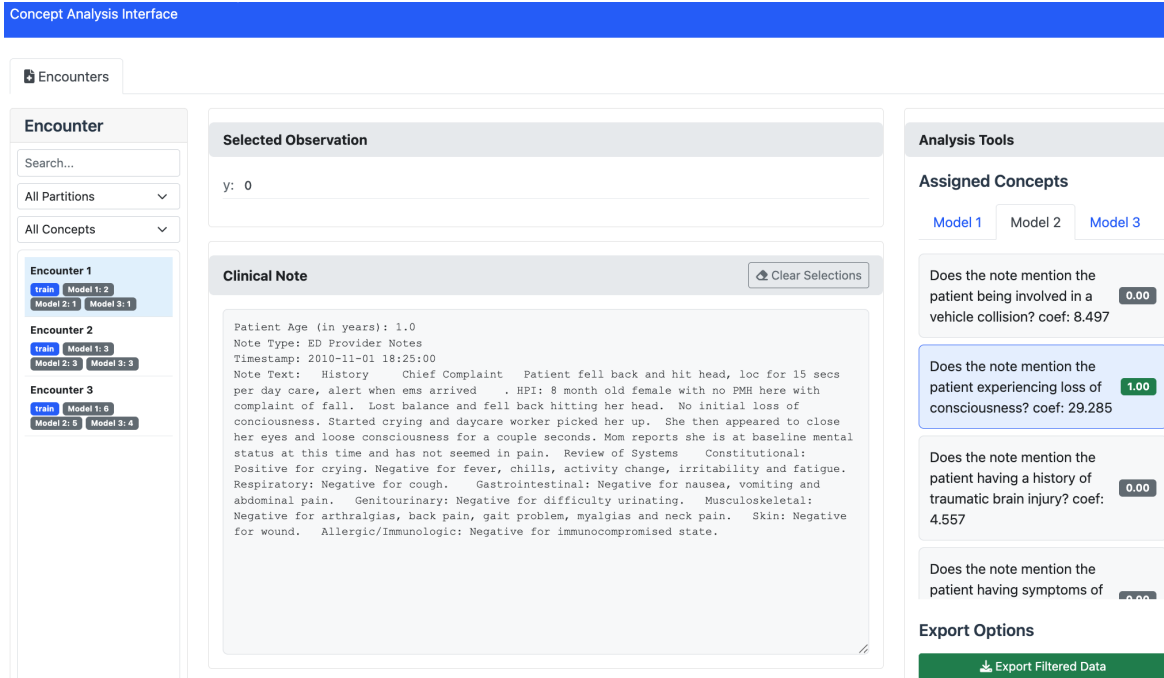
**Fig. 2**: **PHI-compliant web interface for auditing the AI-agent CPM learning procedure**. It shows (i) the learned factors, (ii) the LLM concept annotations at a patient/note level, (iii) which patients received an incorrect prediction, and (iv) the performance of the CPM on a held-out test set.

Revision [ICD-10]) codes: S06, S00.03, S00.83, S00.93, S02.0, S02.1, S02.80, S02.81, S02.91, S07.1, S07.8, S07.9, S09.8, and S09.90 for head trauma, and S06.1, S06.2, S06.3, S06.5, S06.6, S06.89, S06.9, S04.02, S04.03, S04.04, S07.1, and T74.4 for TBI. Encounters with at least one qualifying hospital-wide diagnostic code within the same encounter or occurring up to one calendar day before or after the ED visit were included, to capture both ED and Inpatient Diagnoses. Concatenated, the notes per encounter contain an average of 1,454 tokens. Table 3 (top) shows the prevalence of various concepts in the final dataset.

**Data collection for acute kidney injury (AKI).** The retrospective case-control cohort was assembled from all General Surgery cases between January 2016 to March 2024 at UCSF. Inclusion criteria included adult patients undergoing inpatient procedures (including emergency surgeries) with at least one serum creatinine (sCr) value in the 90 days preceding surgery and at least one serum sCr in the 7 days following surgery. The main outcome was AKI according to the KDIGO criteria in the 7 days following surgery: [48 h maximum postoperative sCr] - [last preoperative sCr] $\geq$ 0.3 mg/dL or [7 day maximum postoperative sCr]/[last preoperative sCr] $\geq$ 1.5 [44]. We analyze the preoperative note written by the anesthesiologist, which has an average of 1991 tokens. Table 3 (bottom) shows the prevalence of various concepts in the final dataset. Among the AKI cases, the prevalence of AKI Stage 1 was 0.74 (0.66, 0.80), AKI Stage 2 was 0.12 (0.08, 0.14), and AKI Stage 3 was 0.16 (0.12, 0.20).

**Table 3**: Final concepts in CPM for AKI and TBI learned by HACHI. Coefficients and prevalences shown.

| TBI Concept | Coefficient | Prevalence |
| --- | --- | --- |
| ... the patient having normal gait? | -0.22 | 0.56 (0.51, 0.61) |
| ... the patient experiencing head trauma? | 0.08 | 0.95 (0.93, 0.97) |
| ... the patient experiencing a headache? | 0.25 | 0.42 (0.37, 0.47) |
| ... the patient having altered mental status? | 0.73 | 0.13 (0.10, 0.16) |
| ... the patient experiencing unconsciousness? | 1.58 | 0.39 (0.34, 0.44) |

| AKI Concept | Coefficient | Prevalence |
| --- | --- | --- |
| ... the surgery being minimally invasive, e.g., laparoscopic or robotic-assisted procedures? | -0.74 | 0.41 (0.37, 0.44) |
| ... the patient having hypertension, e.g., high blood pressure or antihypertensive medication use? | 0.31 | 0.54 (0.51, 0.58) |
| ... the patient having obesity, e.g., high BMI or overweight status? | 0.39 | 0.32 (0.29, 0.36) |
| ... the surgery being urgent or emergent, e.g., requiring immediate attention or intervention? | 0.50 | 0.21 (0.18, 0.23) |
| ... the patient having a history of sleep apnea, e.g., witnessed apnea or CPAP use? | 0.51 | 0.20 (0.17, 0.23) |
| ... the patient having heart failure, e.g., reduced ejection fraction or history of heart failure exacerbations? | 0.64 | 0.08 (0.06, 0.10) |
| ... the patient having an active systemic infection, e.g., sepsis, bacteremia, or requiring therapeutic antibiotics? | 0.72 | 0.07 (0.05, 0.09) |
| ... the patient having tachycardia, e.g., heart rate > 100 bpm or palpitations? | 0.88 | 0.16 (0.14, 0.19) |
| ... the surgery being major or high-risk, e.g., major abdominal surgery or anticipated blood loss > 500mL? | 0.91 | 0.40 (0.36, 0.43) |
| ... the patient having chronic kidney disease, e.g., a history of CKD or elevated creatinine levels? | 1.05 | 0.13 (0.11, 0.16) |

# 5 Declaration Statements

## 5.1 Data Availability

Data that support the findings of this study are not publicly available due to the use of protected health information. Deidentified data can be made available from the corresponding author upon reasonable request.

## 5.2 Code Availability

An open-source python package for the HACHI framework is publicly available at http://github.com/jjfenglab/HACHI, which includes code for reproducing this work, web interfaces for reviewing outputs from the AI agent, and tutorials for running the AI-guided CPM learning procedure.

## 5.3 Acknowledgements

API gateway. The authors would also like to thank Joanne Yim for helping with extracting data for the AKI case study.

## 5.4 Author Contributions

J.F., A. Kothari, L.Z., P.V., Y.T., and C.S. developed the methodology and software. J.F., A. Kothari, A.B., N.A., and A. Kornblith curated data. J.F., A. Kothari, P.V., A.B., N.A., A. Kornblith, Y.T., and C.S. analyzed results and validated the methodology. J.F., A. Kothari, Y.T., and C.S. wrote the main manuscript text. All authors contributed to the conceptualization of the work and reviewed/edited the manuscript.

## 5.5 Competing Interests

The authors have no competing interests as defined by Nature Portfolio, or other interests that might be perceived to influence the results and/or discussion reported in this paper.

## References

[1] Ranzani, O.T., Singer, M., Salluh, J.I.F., Shankar-Hari, M., Pilcher, D., Berger-Estilita, J., Coopersmith, C.M., Juffermans, N.P., Laffey, J., Reinikainen, M., Neto, A.S., Tavares, M., Timsit, J.-F., Arias Lopez, M.D.P., Arulkumaran, N., Aryal, D., Azoulay, E., Celi, L.A., Chaudhuri, D., De Lange, D., De Waele, J., Dos Santos, C.C., Du, B., Einav, S., Engelbrecht, T., Fazla, F., Ferrer, R., Finazzi, S., Fujii, T., Gershengorn, H.B., Greene, J.D., Haniffa, R., Hao, S., Hasan, M.S., Hollenberg, S., Ippolito, M., Jung, C., Kirov, M., Kobari, S., Lakbar, I., Lipman, J., Liu, V., Liu, X., Lobo, S.M., Magatti, D., Martin, G.S., Metnitz, B., Metnitz, P., Myatra, S.N., Oczkowski, S., Paiva, J.-A., Paruk, F., Pekkarinen, P.T., Piquilloud, L., Pölkki, A., Prescott, H.C., Blaser, A.R., Rezende, E., Robba, C., Rochwerg, B., Ruckly, S., Samei, R., Schenck, E.J., Secombe, P., Sendagire, C., Siaw-Frimpong, M., Simpkin, A.J., Soares, M., Summers, C., Szczeklik, W., Takala, J., Tanaka, S., Tricella, G., Vincent, J.-L., Wendon, J., Zampieri, F.G., Rhodes, A., Moreno, R.: Development and validation of the sequential organ failure assessment (SOFA)-2 score. JAMA **334**(23), 2090–2103 (2025)

[2] Kuppermann, N., Holmes, J.F., Dayan, P.S., Hoyle, J.D., Atabaki, S.M., Holubkov, R., Nadel, F.M., Monroe, D., Stanley, R.M., Borgialli, D.A., *et al.*: Identification of children at very low risk of clinically-important brain injuries after head trauma: a prospective cohort study. The Lancet **374**(9696), 1160–1170 (2009)

[3] Teasdale, G., Jennett, B.: Assessment of coma and impaired consciousness. a practical scale. Lancet **2**(7872), 81–84 (1974)

[4] Laupacis, A., Sekar, N., Stiell, I.G.: Clinical prediction rules. a review and suggested modifications of methodological standards. JAMA **277**(6), 488–494 (1997)

[5] Rudin, C.: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nature Machine Intelligence **1**(5), 206–215 (2019)

[6] Koh, P.W., Nguyen, T., Tang, Y.S., Mussmann, S., Pierson, E., Kim, B., Liang, P.: Concept bottleneck models. ICML **abs/2007.04612** (2020)

[7] Wells, P.S., Anderson, D.R., Rodger, M., Ginsberg, J.S., Kearon, C., Gent, M., Turpie, A.G., Bormanis, J., Weitz, J., Chamberlain, M., Bowie, D., Barnes, D., Hirsh, J.: Derivation of a simple

clinical model to categorize patients probability of pulmonary embolism: increasing the models utility with the SimpliRED D-dimer. Thromb. Haemost. **83**(3), 416–420 (2000)

[8] Chung, F., Abdullah, H.R., Liao, P.: STOP-bang questionnaire: A practical approach to screen for obstructive sleep apnea. Chest **149**(3), 631–638 (2016)

[9] McLendon, K., Goyal, A., Attia, M.: Deep venous thrombosis risk factors. In: StatPearls. StatPearls Publishing, Treasure Island (FL) (2025)

[10] Jain, S., Margetis, K., Iverson, L.M.: Glasgow coma scale. In: StatPearls. StatPearls Publishing, Treasure Island (FL) (2025)

[11] Damen, J.A.A.G., Hooft, L., Schuit, E., Debray, T.P.A., Collins, G.S., Tzoulaki, I., Lassale, C.M., Siontis, G.C.M., Chiocchia, V., Roberts, C., Schlüssel, M.M., Gerry, S., Black, J.A., Heus, P., Schouw, Y.T., Peelen, L.M., Moons, K.G.M.: Prediction models for cardiovascular disease risk in the general population: systematic review. BMJ **353**, 2416 (2016)

[12] Wynants, L., Van Calster, B., Collins, G.S., Riley, R.D., Heinze, G., Schuit, E., Bonten, M.M.J., Dahly, D.L., Damen, J.A.A., Debray, T.P.A., Jong, V.M.T., De Vos, M., Dhiman, P., Haller, M.C., Harhay, M.O., Henckaerts, L., Heus, P., Kammer, M., Kreuzberger, N., Lohmann, A., Luijken, K., Ma, J., Martin, G.P., McLernon, D.J., Andaur Navarro, C.L., Reitsma, J.B., Sergeant, J.C., Shi, C., Skoetz, N., Smits, L.J.M., Snell, K.I.E., Sperrin, M., Spijker, R., Steyerberg, E.W., Takada, T., Tzoulaki, I., Kuijk, S.M.J., Bussel, B., Horst, I.C.C., Royen, F.S., Verbakel, J.Y., Wallisch, C., Wilkinson, J., Wolff, R., Hooft, L., Moons, K.G.M., Smeden, M.: Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. BMJ **369**, 1328 (2020)

[13] Feng, Y., Wang, A.Y., Jun, M., Pu, L., Weisbord, S.D., Bellomo, R., Hong, D., Gallagher, M.: Characterization of risk prediction models for acute kidney injury: A systematic review and meta-analysis. JAMA Netw. Open **6**(5), 2313359 (2023)

[14] Os, H.J.A., Kanning, J.P., Wermer, M.J.H., Chavannes, N.H., Numans, M.E., Ruigrok, Y.M., Zwet, E.W., Putter, H., Steyerberg, E.W., Groenwold, R.H.H.: Developing clinical prediction models using primary care electronic health record data: The impact of data preparation choices on model performance. Front. Epidemiol. **2**, 871630 (2022)

[15] Ramaswamy, V.V., Kim, S.S.Y., Fong, R., Russakovsky, O.: Overlooked factors in concept-based explanations: Dataset choice, concept learnability, and human capability. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10932–10941 (2023)

[16] Obra, J.K., Singh, C., Watkins, K., Feng, J., Obermeyer, Z., Kornblith, A.: Potential for algorithmic bias in clinical decision instrument development. NPJ Digit. Med., 1–7 (2025)

[17] Mattei, T.A., Teasdale, G.M.: The story of the development and adoption of the glasgow coma scale: Part I, the early years. World Neurosurg. **134**, 311–322 (2020)

[18] Jiang, L.Y., Liu, X.C., Nejatian, N.P., Nasir-Moin, M., Wang, D., Abidin, A., Eaton, K., Riina, H.A., Laufer, I., Punjabi, P., Miceli, M., Kim, N.C., Orillac, C., Schnurman, Z., Livia, C., Weiss, H., Kurland, D., Neifert, S., Dastagirzada, Y., Kondziolka, D., Cheung, A.T.M., Yang, G., Cao, M., Flores, M., Costa, A.B., Aphinyanaphongs, Y., Cho, K., Oermann, E.K.: Health system-scale language models are all-purpose prediction engines. Nature (2023)

[19] Seinen, T.M., Fridgeirsson, E.A., Ioannou, S., Jeannetot, D., John, L.H., Kors, J.A., Markus,

A.F., Pera, V., Rekkas, A., Williams, R.D., Yang, C., Mulligen, E.M., Rijnbeek, P.R.: Use of unstructured text in prognostic clinical prediction models: a systematic review. J. Am. Med. Inform. Assoc. **29**(7), 1292–1302 (2022)

[20] Seinen, T.M., Kors, J.A., Mulligen, E.M., Rijnbeek, P.R.: Using structured codes and free-text notes to measure information complementarity in electronic health records: Feasibility and validation study. J. Med. Internet Res. **27**(1), 66910 (2025)

[21] Agrawal, M., Hegselmann, S., Lang, H., Kim, Y., Sontag, D.: Large language models are few-shot clinical information extractors. Empirical Methods in Natural Language Processing, 1998–2022 (2022)

[22] Yang, X., Chen, A., PourNejatian, N., Shin, H.C., Smith, K.E., Parisien, C., Compas, C., Martin, C., Costa, A.B., Flores, M.G., Zhang, Y., Magoc, T., Harle, C.A., Lipori, G., Mitchell, D.A., Hogan, W.R., Shenkman, E.A., Bian, J., Wu, Y.: A large language model for electronic health records. NPJ Digit. Med. **5**(1), 194 (2022)

[23] Guevara, M., Chen, S., Thomas, S., Chaunzwa, T.L., Franco, I., Kann, B.H., Moningi, S., Qian, J.M., Goldstein, M., Harper, S., Aerts, H.J.W.L., Catalano, P.J., Savova, G.K., Mak, R.H., Bitterman, D.S.: Large language models to identify social determinants of health in electronic health records. npj Digital Medicine **7**(1), 1–14 (2024)

[24] Oikarinen, T., Das, S., Nguyen, L.M., Weng, T.-W.: Label-free concept bottleneck models. International Conference on Learning Representations (2023) [cs.LG]

[25] McInerney, D.J., Young, G., Meent, J.-W., Wallace, B.C.: CHiLL: Zero-shot custom interpretable feature extraction from clinical notes with large language models. In: The 2023 Conference on Empirical Methods in Natural Language Processing (2023)

[26] Benara, V., Singh, C., Morris, J., Antonello, R., Stoica, I., Huth, A., Gao, J.: Crafting interpretable embeddings for language neuroscience by asking llms questions. Advances in Neural Information Processing Systems **37**, 124137–124162 (2025)

[27] Kim, J., Wang, Z., Qiu, Q.: Constructing concept-based models to mitigate spurious correlations with minimal human effort. In: European Conference on Computer Vision, pp. 137–153 (2024). Springer

[28] Feng, J., Kothari, A., Zier, L., Singh, C., Tan, Y.S.: Bayesian concept bottleneck models with LLM priors. In: The Thirty-ninth Annual Conference on Neural Information Processing Systems (2025). https://openreview.net/forum?id=oXSkzIXgbk

[29] Ludan, J.M., Lyu, Q., Yang, Y., Dugan, L., Yatskar, M., Callison-Burch, C.: Interpretable-by-design text classification with iteratively generated concept bottleneck. arXiv preprint arXiv:2310.19660 (2023)

[30] Subramonyam, H., Pea, R., Pondoc, C., Agrawala, M., Seifert, C.: Bridging the gulf of envisioning: Cognitive challenges in prompt based interactions with LLMs. In: Proceedings of the CHI Conference on Human Factors in Computing Systems, vol. 31, pp. 1–19. ACM, New York, NY, USA (2024)

[31] Kothari, A., Vossler, P., Digitale, J., Forouzannia, M., Rosenberg, E., Lee, M., Bryant, J., Molina, M., Marks, J., Zier, L., Feng, J.: When the domain expert has no time and the LLM developer

has no clinical expertise: Real-world lessons from LLM co-design in a safety-net hospital. Proc. Conf. AAAI Artif. Intell. (2026) [cs.CY]

[32] Langlois, J.A., Rutland-Brown, W., Thomas, K.E.: Traumatic brain injury in the united states: emergency department visits, hospitalizations, and deaths (2006)

[33] Coronado, V.G., Xu, L., Basavaraju, S.V., McGuire, L.C., Wald, M.M., Faul, M.D., Guzman, B.R., Hemphill, J.D., Disease Control, C., (CDC), P., et al.: Surveillance for traumatic brain injury-related deaths: United states, 1997-2007 (2011)

[34] Brenner, D.J.: Estimating cancer risks from pediatric ct: going from the qualitative to the quantitative. Pediatric radiology **32**(4), 228–231 (2002)

[35] Brenner, D.J., Hall, E.J.: Computed tomography—an increasing source of radiation exposure. New England journal of medicine **357**(22), 2277–2284 (2007)

[36] Easter, J.S., Bakes, K., Dhaliwal, J., Miller, M., Caruso, E., Haukoos, J.S.: Comparison of pecarn, catch, and chalice rules for children with minor head injury: a prospective cohort study. Annals of emergency medicine **64**(2), 145–152 (2014)

[37] Holmes, J.F., Yen, K., Ugalde, I.T., Ishimine, P., Chaudhari, P.P., Atigapramoj, N., Badawy, M., McCarten-Gibbs, K.A., Nielsen, D., Sage, A.C., *et al.*: Pecarn prediction rules for ct imaging of children presenting to the emergency department with blunt abdominal or minor head trauma: a multicentre prospective validation study. The Lancet Child & Adolescent Health **8**(5), 339–347 (2024)

[38] Yen, K., Kuppermann, N., Lillis, K., Monroe, D., Borgialli, D., Kerrey, B.T., Sokolove, P.E., Ellison, A.M., Cook, L.J., Holmes, J.F., Intra-abdominal Injury Study Group for the Pediatric Emergency Care Applied Research Network (PECARN): Interobserver agreement in the clinical assessment of children with blunt abdominal trauma. Acad. Emerg. Med. **20**(5), 426–432 (2013)

[39] Hurt, R.T., Stephenson, C.R., Gilman, E.A., Aakre, C.A., Croghan, I.T., Mundi, M.S., Ghosh, K., Edakkanambeth Varayil, J.: The use of an artificial intelligence platform OpenEvidence to augment clinical decision-making for primary care physicians. J. Prim. Care Community Health **16**, 21501319251332215 (2025)

[40] Wijeysundera, D.N., Karkouti, K., Beattie, W.S., Rao, V., Ivanov, J.: Improving the identification of patients at risk of postoperative renal failure after cardiac surgery. Anesthesiology **104**(1), 65–72 (2006)

[41] Thakar, C.V., Arrigain, S., Worley, S., Yared, J.-P., Paganini, E.P.: A clinical score to predict acute renal failure after cardiac surgery. Journal of the American Society of Nephrology **16**(1), 162–168 (2005)

[42] Chertow, G.M., Lazarus, J.M., Christiansen, C.L., Cook, E.F., Hammermeister, K.E., Grover, F., Daley, J.: Preoperative renal risk stratification. Circulation **95**(4), 878–884 (1997)

[43] Kheterpal, S., Tremper, K.K., Heung, M., Rosenberg, A.L., Englesbe, M., Shanks, A.M., Campbell, D.A. Jr: Development and validation of an acute kidney injury risk index for patients undergoing general surgery: results from a national data set: Results from a national data set. Anesthesiology **110**(3), 505–515 (2009)

[44] Kellum, J.A., Lameire, N., Aspelin, P., Barsoum, R.S., Burdmann, E.A., Goldstein, S.L., Herzog, C.A., Joannidis, M., Kribben, A., Levey, A.S., MacLeod, A.M., Mehta, R.L., Murray, P.T., Naicker, S., Opal, S.M., Schaefer, F., Schetz, M., Uchino, S.: Kidney disease: Improving global outcomes (KDIGO) acute kidney injury work group. KDIGO clinical practice guideline for acute kidney injury. Kidney Int. Suppl. (2011) **2**(1), 1 (2012)

[45] Futoma, J., Simons, M., Panch, T., Doshi-Velez, F., Celi, L.A.: The myth of generalisability in clinical research and machine learning in health care. Lancet Digit Health **2**(9), 489–492 (2020)

[46] Patel, A., Rao, D., Kothary, A., McKeown, K., Callison-Burch, C.: Learning interpretable style embeddings via prompting LLMs. Proc. Conf. Empir. Methods Nat. Lang. Process. (2023) [cs.CL]

[47] Sun, Y., Huang, Q., Tang, Y., Tung, A.K.H., Yu, J.: A general framework for producing interpretable semantic text embeddings. In: The Thirteenth International Conference on Learning Representations (2025). https://openreview.net/forum?id=23uY3FpQxc

[48] Romera-Paredes, B., Barekatain, M., Novikov, A., Balog, M., Kumar, M.P., Dupont, E., Ruiz, F.J., Ellenberg, J.S., Wang, P., Fawzi, O., et al.: Mathematical discoveries from program search with large language models. Nature **625**(7995), 468–475 (2024)

[49] Novikov, A., Vũ, N., Eisenberger, M., Dupont, E., Huang, P.-S., Wagner, A.Z., Shirobokov, S., Kozlovskii, B., Ruiz, F.J., Mehrabian, A., et al.: Alphaevolve: A coding agent for scientific and algorithmic discovery. arXiv preprint arXiv:2506.13131 (2025)

[50] Singh, C., Morris, J.X., Aneja, J., Rush, A.M., Gao, J.: Explaining patterns in data with language models via interpretable autoprompting. arXiv preprint arXiv:2210.01848 (2022)

[51] Monarch, R.M.: Human-in-the-Loop Machine Learning: Active Learning and Annotation for Human-centered AI. Simon and Schuster, ??? (2021)

[52] Mosqueira-Rey, E., Hernández-Pereira, E., Alonso-Ríos, D., Bobes-Bascarán, J., Fernández-Leal, Á.: Human-in-the-loop machine learning: a state of the art. Artificial Intelligence Review **56**(4), 3005–3054 (2023)

[53] Lage, I., Doshi-Velez, F.: Learning interpretable concept-based models with human feedback. arXiv preprint arXiv:2012.02898 (2020)

[54] Brewster, R.C., Tse, G., Fan, A.L., Elborki, M., Newell, M., Gonzalez, P., Hoq, A., Chang, C., Chowdhury, M., Geeti, A., et al.: Evaluating human-in-the-loop strategies for artificial intelligence-enabled translation of patient discharge instructions: a multidisciplinary analysis. NPJ digital medicine **8**(1), 629 (2025)

[55] Sivaraman, V., Vaishampayan, A., Li, X., Buck, B.R., Ma, Z., Boyce, R.D., Perer, A.: Tempo: Helping data scientists and domain experts collaboratively specify predictive modeling tasks. In: Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems, pp. 1–18. ACM, New York, NY, USA (2025)

[56] Gao, J., Gebreegziabher, S.A., Choo, K.T.W., Li, T.J.-J., Perrault, S.T., Malone, T.W.: A taxonomy for human-llm interaction modes: An initial exploration. In: Extended Abstracts of the CHI Conference on Human Factors in Computing Systems, pp. 1–11 (2024)

[57] Wu, S., Galley, M., Peng, B., Cheng, H., Li, G., Dou, Y., Cai, W., Zou, J., Leskovec, J., Gao,

J.: Collabllm: From passive responders to active collaborators. In: Forty-second International Conference on Machine Learning

[58] Chauhan, K., Tiwari, R., Freyberg, J., Shenoy, P., Dvijotham, K.: Interactive concept bottleneck models. In: Proceedings of the Aaai Conference on Artificial Intelligence, vol. 37, pp. 5948–5955 (2023)

[59] Ustun, B., Rudin, C.: Supersparse linear integer models for optimized medical scoring systems. Machine Learning **102**(3), 349–391 (2016)

[60] Yang, H., Rudin, C., Seltzer, M.: Scalable bayesian rule lists. In: International Conference on Machine Learning, pp. 3921–3930 (2017). PMLR

[61] Grari, V., Arni, T., Laugel, T., Lamprier, S., Zou, J., Detyniecki, M.: Act: Agentic classification tree. arXiv preprint arXiv:2509.26433 (2025)

[62] Singh, C., Morris, J., Rush, A.M., Gao, J., Deng, Y.: Tree prompting: Efficient task adaptation without fine-tuning. In: Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 6253–6267 (2023)

[63] Ragkousis, A., Parbhoo, S.: Tree-based leakage inspection and control in concept bottleneck models. arXiv preprint arXiv:2410.06352 (2024)

[64] Singh, C., Nasseri, K., Tan, Y.S., Tang, T., Yu, B.: imodels: a python package for fitting interpretable models. Journal of Open Source Software **6**(61), 3192 (2021) https://doi.org/10.21105/joss.03192

[65] Kornblith, A.E., Singh, C., Innes, J.C., Chang, T.P., Adelgais, K.M., Holsti, M., Kim, J., McClain, B., Nishijima, D.K., Rodgers, S., *et al.*: Analyzing patient perspectives with large language models: a cross-sectional study of sentiment and thematic classification on exception from informed consent. Scientific reports **15**(1), 6179 (2025)

[66] Shi, T., Yan, G., Oikarinen, T., Weng, T.-W.: Multimodal concept bottleneck models. In: Mechanistic Interpretability Workshop at NeurIPS 2025

[67] Pang, W., Ke, X., Tsutsui, S., Wen, B.: Integrating clinical knowledge into concept bottleneck models. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 243–253 (2024). Springer

# A Extended Data

**Table A1**: **Complete list of concepts for TBI prediction across HACHI iterations.** In round 1, we additionally consider running HACHI in an outcome-agnostic way, where the LLM does not know which outcome it is predicting and must rely on the prediction model to identify useful features. Nevertheless, the model successfully identifies useful features, even recovering a few of the original PECARN concepts.

| Round | Coef. | Concepts |
|---|---|---|
| **Round 1** | 2.70 | Does the note mention the patient experiencing a loss of consciousness? |
| | 1.74 | Does the note mention a brain bleed? |
| | 1.58 | Does the note mention the patient having a neurological event? |
| | 1.23 | Does the note mention a Glasgow Coma Scale score? |
| | 1.10 | Does the note mention the patient being seizure-free? |
| **Round 2** | 3.46 | Does the note mention the patient experiencing loss of consciousness? |
| | 1.50 | Does the note mention the patient having a normal Glasgow Coma Scale score and being at least 2 years old? |
| | 1.04 | Does the note mention the patient experiencing convulsions? |
| | 0.82 | Does the note mention the patient having altered mental status? |
| | -0.65 | Does the note mention the patient having intact cranial nerves? |
| **Round 3** | 4.34 | Does the note mention the patient experiencing loss of consciousness? |
| | 1.41 | Does the note mention the patient having a history of mild TBI? |
| | 1.21 | Does the note mention the patient having an occipital hematoma? |
| | 1.07 | Does the note mention the patient having vision changes? |
| | -0.53 | Does the note mention the patient having memory intact? |
| **Round 4** | 1.58 | Does the note mention the patient experiencing unconsciousness? |
| | 0.73 | Does the note mention the patient having altered mental status? |
| | 0.25 | Does the note mention the patient experiencing a headache? |
| | 0.08 | Does the note mention the patient experiencing head trauma? |
| | -0.22 | Does the note mention the patient having normal gait? |

**Table A2**: **Comparing the sensitivity and specificity of TBI models from HACHI and comparator methods.** To make the models comparable, the threshold is chosen so that the sensitivity is as close to 0.90 as possible. 95% CI is shown in parentheses.

| Model | Baseline (PECARN) | | OpenEvidence+ LLMs | | HACHI (Round 1) | | HACHI (Round 2) | | HACHI (Round 3) | | HACHI (Round 4) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Sensitivity | 0.942 | (0.942, 1.000) | 0.913 | (0.900, 1.000) | 0.906 | (0.900, 0.993) | 0.913 | (0.900, 0.938) | 0.906 | (0.900, 0.950) | 0.906 | (0.900, 0.938) |
| Specificity | 0.556 | (0.000, 0.613) | 0.642 | (0.000, 0.884) | 0.481 | (0.154, 0.906) | 0.615 | (0.510, 0.875) | 0.749 | (0.332, 0.891) | 0.807 | (0.545, 0.883) |

**Table A3**: Complete list of concepts for AKI prediction across HACHI iterations.

| Round | Coef. | Concepts |
|---|---|---|
| **1** | 1.31 | Does the note mention this patient having leukocytosis? |
| | 0.87 | Does the note mention this patient having abdominal distention? |
| | 0.74 | Does the note mention this patient having cardiac dysfunction? |
| | 0.71 | Does the note mention this patient having tachycardia? |
| | 0.65 | Does the note mention this patient experiencing swelling? |
| | 0.57 | Does the note mention this patient having a systemic infection? |
| | 0.53 | Does the note mention this patient having renal impairment? |
| | 0.47 | Does the note mention this patient having diabetes mellitus? |
| | -0.36 | Does the note mention this patient having no kidney disease? |
| | -1.01 | Does the note mention this patient having low hematocrit? |
| **2** | 1.28 | Does the note mention the patient having chronic kidney disease? |
| | 1.03 | Does the note mention the patient undergoing an exploratory laparotomy? |
| | 0.98 | Does the note mention the patient having sepsis? |
| | 0.89 | Does the note mention the patient having fluid retention? |
| | 0.47 | Does the note mention the patient having a respiratory condition? |
| | 0.43 | Does the note mention the patient requiring a blood transfusion during surgery? |
| | 0.40 | Does the note mention the surgery being high-risk? |
| | 0.37 | Does the note mention the patient having heart disease? |
| | 0.31 | Does the note mention the patient having a malignancy? |
| | -1.01 | Does the note mention the patient having minimally invasive surgery? |
| **3** | 1.05 | Does the note mention the patient having chronic kidney disease, e.g., a history of CKD or elevated creatinine levels? |
| | 0.91 | Does the note mention the surgery being major or high-risk, e.g., major abdominal surgery or anticipated blood loss $> 500$mL? |
| | 0.88 | Does the note mention the patient having tachycardia, e.g., heart rate $> 100$ bpm or palpitations? |
| | 0.72 | Does the note mention the patient having an active systemic infection, e.g., sepsis, bacteremia, or requiring therapeutic antibiotics? |
| | 0.64 | Does the note mention the patient having heart failure, e.g., reduced ejection fraction or history of heart failure exacerbations? |
| | 0.51 | Does the note mention the patient having a history of sleep apnea, e.g., witnessed apnea or CPAP use? |
| | 0.50 | Does the note mention the surgery being urgent or emergent, e.g., requiring immediate attention or intervention? |
| | 0.39 | Does the note mention the patient having obesity, e.g., high BMI or overweight status? |
| | 0.31 | Does the note mention the patient having hypertension, e.g., high blood pressure or antihypertensive medication use? |
| | -0.74 | Does the note mention the surgery being minimally invasive, e.g., laparoscopic or robotic-assisted procedures? |

# B  Example prompts for TBI in Round 1

Here we show examples of prompts used for TBI in Round 1; see all prompts on Github.

## B.1  KeyphrasePrompt

```
You are extracting clinical descriptors of a pediatric patient from emergency room medical notes that may influence
↪   their risk of being diagnosed with traumatic brain injury (TBI).  Focus entirely on factors explicitly mentioned
↪   in the notes. Do not infer or add details unless explicitly stated.

Extract factors that may:
- Increase TBI risk, be specific and detailed (e.g., substance use, prior head injuries)
- Decrease TBI risk, be specific and detailed (e.g., protective equipment use, medication compliance)
- Indicate TBI severity or complications, be specific and detailed (e.g., neurological symptoms, cognitive changes)

For each factor, provide:
1. The primary descriptor (exact or paraphrased from the note)
2. Synonyms or alternative phrasings
3. Broader generalizations that capture similar concepts

Requirements:
- Each descriptor should be $\leq$3 words
- The descriptor should state distinguishing features.
- List specific terms, synonyms, and increasingly general concepts. However, do not state overly general concepts like
↪   "protective factor", "risk factor", "low risk", or "high risk" because we are using these terms to
↪   brainstorm/learn specific factors that clinicians can use to predict TBI risk.
```

- Separate items with commas


Note(s):
"{note}"

Do not infer or fabricate details. Base all responses strictly on the medical note text.

Output Format:
```
{
    "reasoning": "Describe any reasoning here",
    "keyphrases": [
        "Primary descriptor, Synonym1, Synonym2, Broader category",
        "Primary descriptor, Alternative phrasing, More general concept",
        ...
    ]
}
```

Examples:
```
{
    "reasoning": "Describe any reasoning here",
    "keyphrases": [
        "Crowded residence, Overcrowded living, Constrained environment, Environmental risk",
        "Cervical strain, Neck pain, Muscle spasm",
        "Physical assault, Blunt trauma, Injury infliction",
        ...
    ]
}
```

## B.2 ConceptInitializationProposalPrompt

We are fitting a concept bottleneck model to predict traumatic brain injury (TBI) for pediatric patients admitted to
↪   emergency room. The goal is to come up with a few meta-concepts that clinicians can use to easily risk-stratify
↪   patients for whether they are at high risk of being diagnosed with TBI. As training data, we are using clinical
↪   notes to retrospectively extract patient characteristics.

The top words most associated with TBI were estimated using a simple logistic regression model with tf-idf features.
↪   Based on the top words, come up with {max_meta_concepts} meta-concepts. A meta-concept has to be defined in terms
↪   of a yes/no question, e.g. "Does the note mention this patient having a fall?". Pick meta-concepts that are as
↪   specific as possible, without being too broad.

Suggestions for generating candidate meta-concepts: Do not propose meta-concepts that are simply a union of two
↪   different concepts (e.g. "Does the note mention this patient having an infection or being elderly?" is not
↪   allowed), questions with answers that are almost always a yes (e.g. the answer to "Does the note mention this
↪   patient being hospitalized?" is almost always yes), or questions where the yes/no options are not clearly defined
↪   (e.g. "Does the note mention this patient experiencing difficulty?" is not clearly defined because difficulty may
↪   mean financial difficulty, physical difficulties, etc).

Words with top coefficients:

{top_features_df}

Propose at least {max_meta_concepts} candidates in the following JSON format:
```
{
  "reasoning": "Step-by-step explanation of how you systematically worked through the model concepts (from most to
  ↪   least predictive).",
  "concepts": [
    {
      "concept": "<QUESTION 1>",
      "words": ["<SYNONYM/ANTONYM 1>", "<SYNONYM/ANTONYM 2>", "<SYNONYM/ANTONYM 3>"]
    },
    {
      "concept": "<QUESTION 2>",
      "words": ["<SYNONYM/ANTONYM 1>", "<SYNONYM/ANTONYM 2>", "<SYNONYM/ANTONYM 3>"]
    },
    ...
    {
      "concept": "<QUESTION {max_meta_concepts}>",
      "words": ["<SYNONYM/ANTONYM 1>", "<SYNONYM/ANTONYM 2>", "<SYNONYM/ANTONYM 3>"]
    }
  ]
}
```

Example answer:

```
{
  "reasoning": "Step-by-step reasoning through the model concepts",
  "concepts": [
    {
      "concept": "Does the note mention the patient experiencing head trauma?",
      "words": ["head injury", "skull trauma", "cranial damage"]
    },
    {
      "concept": \Does the note mention loss of consciousness?"
      "words": ["LOC", "fainting, unresponsive"]
    },
    ...
    ]
}
```

## B.3 ConceptReplacementProposalPrompt

The goal is to come up with a concept bottleneck model (CBM) that extracts {num_concepts} meta-concepts from clinical
↪   notes from the electronic health record of pediatric patients admitted to the emergency room to predict risk of
↪   having a traumatic brain injury (TBI). A meta-concept is a binary feature extractor defined by a yes/no question.
↪   We have {num_concepts_fixed} meta-concepts so far:
{meta_concepts}

To come up with the {num_concepts}th meta-concept, I have done the following: I first fit a CBM on the
↪   {num_concepts_fixed} existing meta-concepts. Then to figure out how to improve this {num_concepts_fixed}-concept
↪   CBM, I first extracted a list of concepts that are present in each note, and then fit a linear regression model on
↪   the extracted concepts to predict the residuals of the {num_concepts_fixed}-concept CBM. These are the top
↪   extracted concepts in the resulting residual model, in descending order of importance:

{top_features_df}

To interpret this residual model, a general rule of thumb is that an extracted concept with a large positive
↪   coefficient means that a high risk of TBI is positively associated with the concept being mentioned in the note,
↪   and an extracted concept with a negative coefficient means that a higher risk of TBI is negatively associated with
↪   the concept being mentioned in the note.

Given the residual model, create cohesive candidates for the {num_concepts}th meta-concept. Be systematic and consider
↪   all the listed concepts in the residual model. Start from the most to the least predictive concept. For each
↪   concept, check if it matches an existing meta-concept or create a new candidate meta-concept. Work down the list,
↪   iterating through each concept. Clearly state each candidate meta-concept as a yes/no question.

Suggestions for generating candidate meta-concepts: Do not propose meta-concepts that are simply a union of two
↪   different concepts (e.g. "Does the note mention this patient having an infection or being elderly?" is not
↪   allowed), questions with answers that are almost always a yes (e.g. the answer to "Does the note mention this
↪   patient being hospitalized?" is almost always yes), or questions where the yes/no options are not clearly defined
↪   (e.g. "Does the note mention this patient experiencing difficulty?" is not clearly defined because difficulty may
↪   mean financial difficulty, physical difficulties, etc). Do not propose meta-concepts where you would expect over
↪   95% agreement or disagreement with the {num_concepts_fixed} existing meta-concepts (e.g. "Does the note mention
↪   the patient having homelessness?" overlaps too much with "Does the note mention the patient having housing
↪   insecurity?").

Finally, summarize all the generated candidates for the {num_concepts}-th meta-concept in a JSON. Merge any candidate
↪   meta-concepts that are essentially the same (where you would expect over 95% agreement) or essentially opposites
↪   (you would expect over 95% disagreement). In the JSON, include a list of comma-separated list of phrases that mean
↪   the same or opposite of each candidate meta-concept. Propose at least ten candidates. The final JSON should have
↪   the following format:

Propose at least six candidates in the following JSON format:
{
  "reasoning": "Step-by-step explanation of how you systematically worked through the model concepts (from most to
  ↪   least predictive) and combined them with the feedback suggestions.",
  "concepts": [
    {
      "concept": "<QUESTION 1>",
      "words": ["<SYNONYM/ANTONYM 1>", "<SYNONYM/ANTONYM 2>", "<SYNONYM/ANTONYM 3>"]
    },
    {
      "concept": "<QUESTION 2>",
      "words": ["<SYNONYM/ANTONYM 1>", "<SYNONYM/ANTONYM 2>", "<SYNONYM/ANTONYM 3>"]
    },
    ...
    {
      "concept": "<QUESTION 6>",
```

```
    "words": ["<SYNONYM/ANTONYM 1>", "<SYNONYM/ANTONYM 2>", "<SYNONYM/ANTONYM 3>"]
    }
  ]
}

Example answer:
{
  "reasoning": "Step-by-step reasoning through the residual model concepts and feedback",
  "concepts": [
    {
      "concept": "Does the note mention the patient experiencing head trauma?",
      "words": ["head injury", "skull trauma", "cranial damage"]
    },
    {
      "concept": \Does the note loss of consciousness?"
      "words": ["LOC", "fainting, unresponsive"]
    },
    ...
    ]
}
```

## B.4 ConceptAnnotationPrompt

```
You will be given clinical notes. I will give you a series of questions. Your task is answer each question with a
↪    probability from 0 to 1. Summarize the response with a JSON that includes your answer to all of the questions.
↪    Questions:
{prompt_questions}

clinical notes:
{sentence}

Example answer: {
    "0": 0,
    "1": 1,
    "2": 0.5
}
Answer all the questions and do not answer with anything else besides valid JSON. Do not add comments to the JSON.
```