

HACHI: Human-AI Co-design for Highly Interpretable Clinical Prediction Models

Jean Feng^{1,*}, Avni Kothari^{1,*}, Patrick Vossler¹, Andrew Bishara¹, Lucas Zier¹, Newton Addo¹, Aaron Kornblith¹, Yan Shuo Tan², Chandan Singh³

¹University of California, San Francisco · ²National University of Singapore · ³Microsoft Research · *Equal contribution



BACKGROUND & MOTIVATION

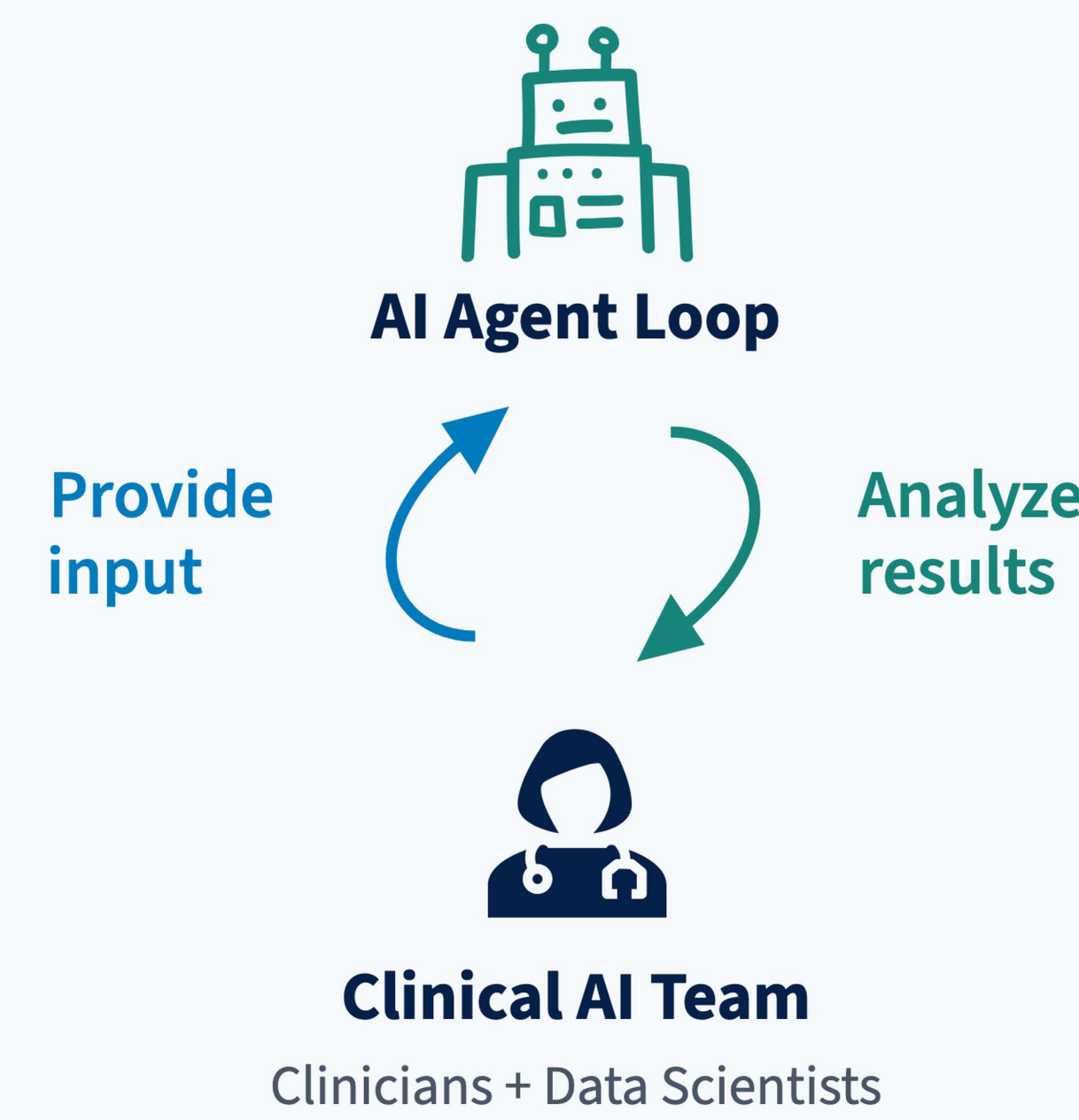
Clinical prediction models (CPMs) are simple scoring tools (e.g., PECARN for pediatric head trauma, SOFA for ICU mortality) that help clinicians assess patient risk. Most rely on structured data, but clinical notes often contain richer information — symptoms, reasoning, and context that structured fields miss. Building CPMs from notes requires lengthy collaboration between clinicians and data scientists.

The gap: Large language models (LLMs) can now extract clinical concepts from notes at scale, but a purely automated approach may miss data quality problems, spurious correlations, and fairness issues that only domain experts would catch.

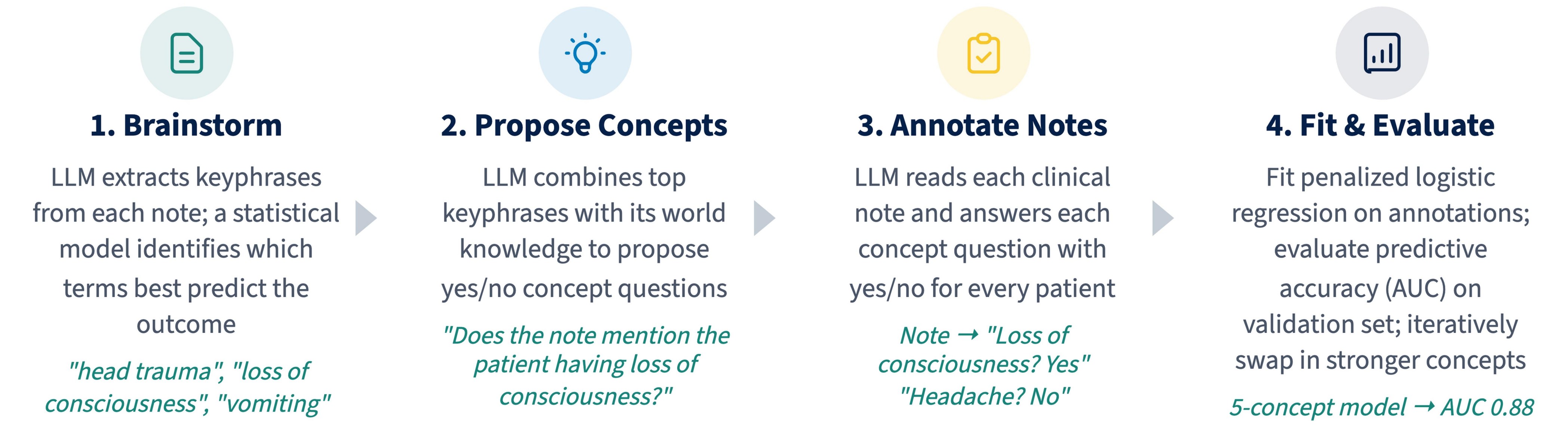
HACHI lets clinical teams iteratively co-design interpretable, note-based CPMs with an AI agent. Each CPM is a logistic regression whose features are yes/no questions ("concepts") that an LLM answers by reading the patient's note.

THE HACHI FRAMEWORK

Outer Loop: Human-AI Co-design



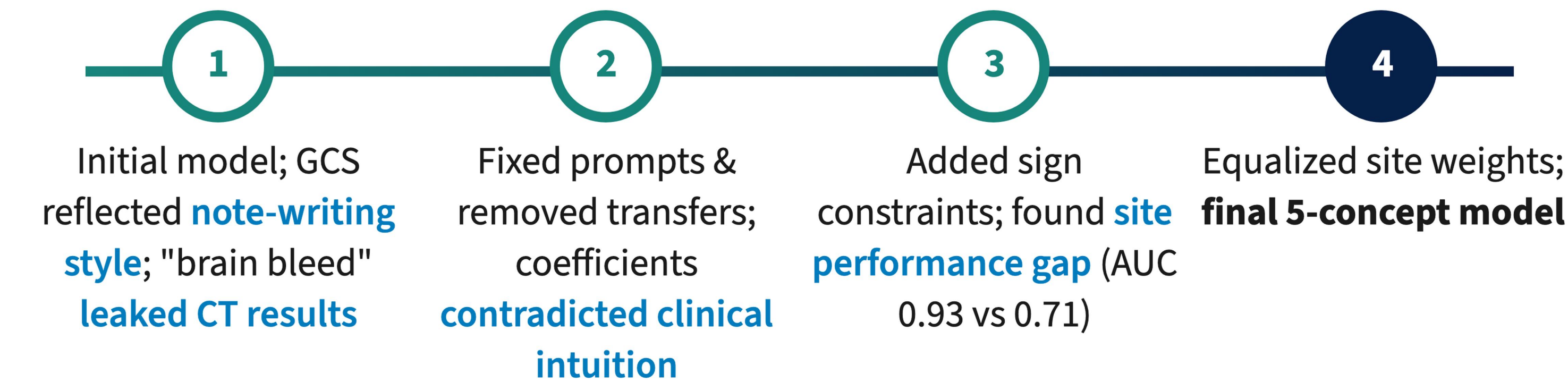
Inner Loop: AI Agent CPM Learning



CASE STUDY: TBI IN CHILDREN

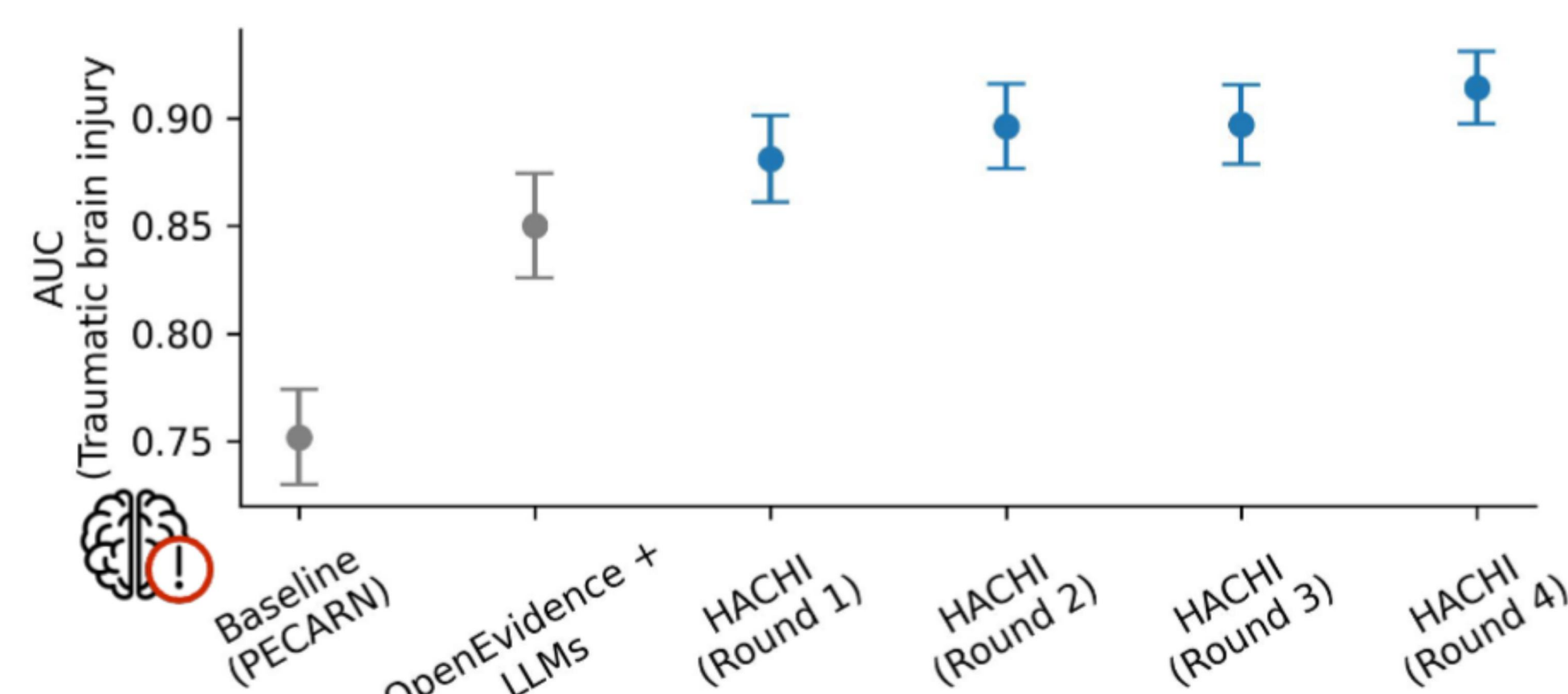
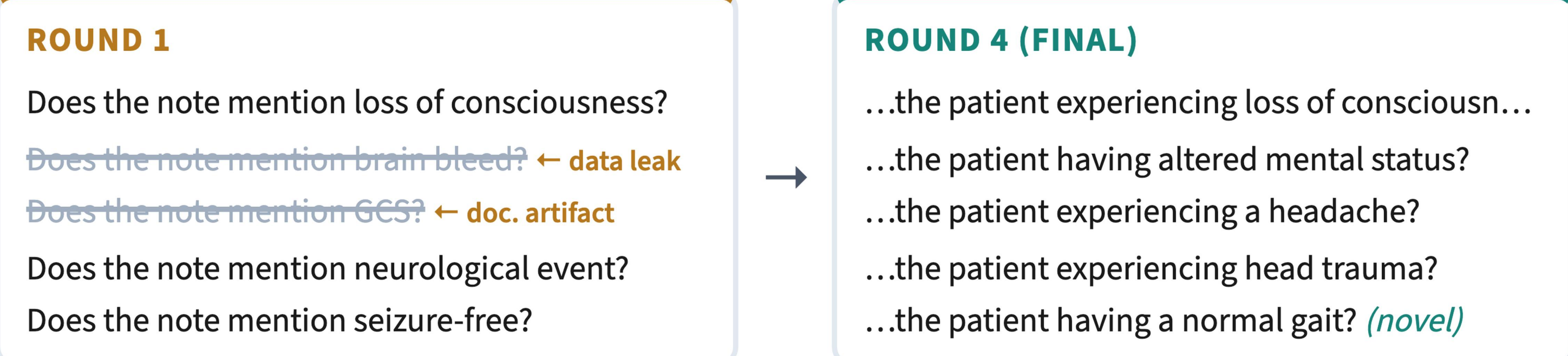
Predicting traumatic brain injury (TBI) in children presenting to the emergency department (ED) after head trauma. Team: pediatric ED clinician, data analyst, and data scientists.

CO-DESIGN ROUNDS



Final 5-Concept Model — AUC 0.91

Simpler than PECARN (current standard, AUC 0.75) · Outperformed OpenEvidence + LLMs (AUC 0.85)



CASE STUDY: ACUTE KIDNEY INJURY

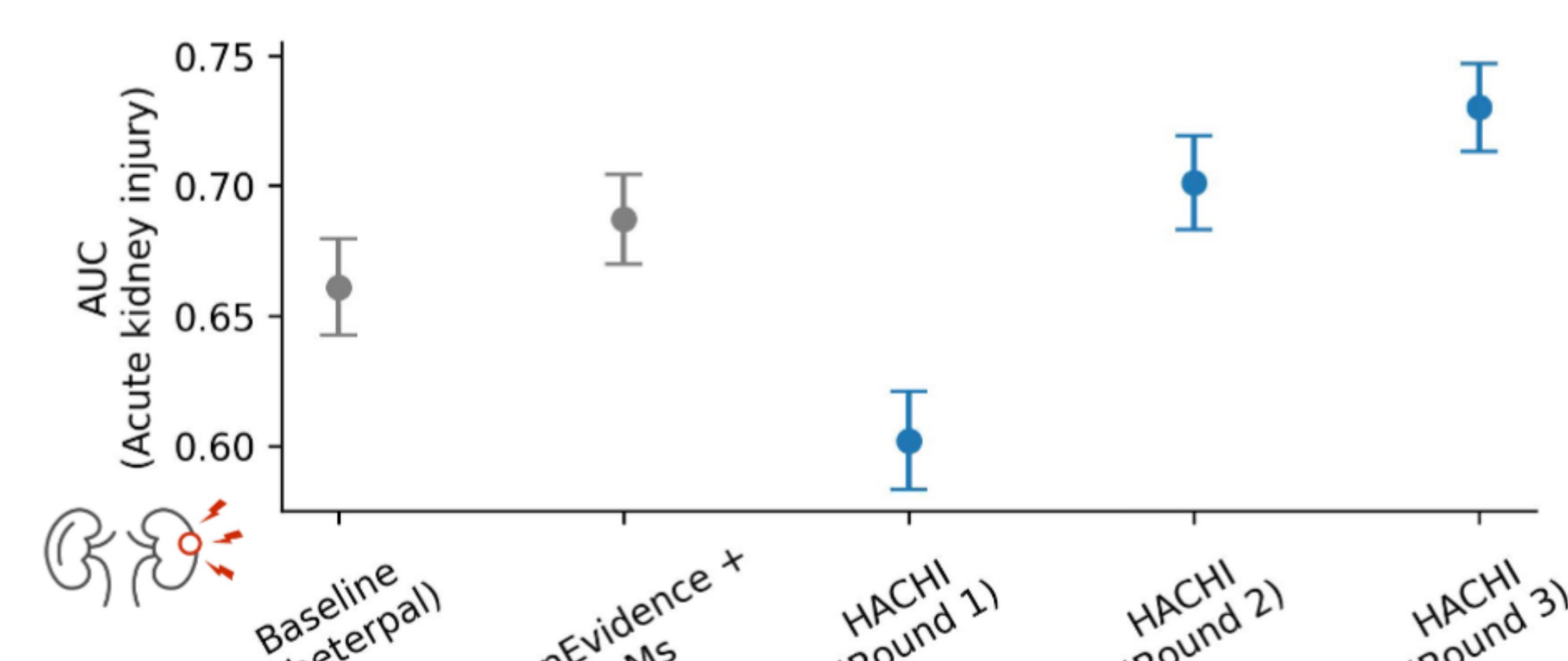
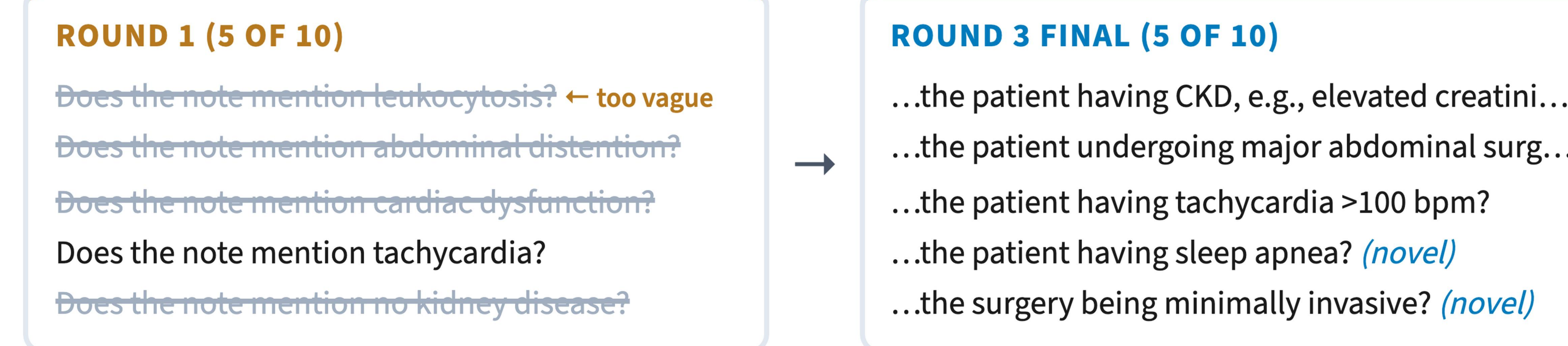
Predicting acute kidney injury (AKI) within 7 days of general surgery. Team: anesthesiologist, data analyst, and data scientists.

CO-DESIGN ROUNDS



Final 10-Concept Model — AUC 0.73

Outperformed Kheterpal index (existing risk tool, AUC 0.64–0.66) · Temporal validation AUC 0.76



PHI-COMPLIANT REVIEW INTERFACE

Clinicians review the AI agent's outputs through a locally hosted interface that surfaces **incorrect predictions** first.

Concept Analysis Interface		
CLINICAL NOTE	CONCEPT ANNOTATIONS	COEFF
Pt age 1.0y. Chief Complaint: fell back, hit head, LOC 15 sec. Started crying, daycare worker picked her up. She appeared to close her eyes and lose consciousness for a couple seconds...	0 ...patient involved in a vehicle collision?	8.50
	1 ...patient experiencing loss of consciousness?	29.3
	0 ...patient having a history of traumatic brain injury?	4.56
	0 ...patient having altered mental status?	1.22
	1 ...patient having a normal gait?	-0.22

This review process helped teams spot **data leakage**, **spurious correlations**, and **documentation artifacts** across the case studies.

TAKEAWAYS

- Each round of human feedback improved model quality — catching data leakage, spurious correlations, and fairness gaps that automated pipelines missed.
- 3–4 rounds of 1–2 hours of clinician time were sufficient to produce final models.
- The framework generalized across two different clinical settings, teams, and data sources.

Open-source code and interface on GitHub →

