

The Natural Language of Clinical Problem-Solving

Raja-Elie E. Abdulnour, M.D.

Editor-in-Chief, NEJM Journal Watch
Editor, Clinical Development and AI Innovation, NEJM Group
'23-'24 Diagnostic Excellence Scholar, National Academy of Medicine
Pulmonary and Critical Care Medicine, Brigham and Women's Hospital
Assistant Professor of Medicine, Harvard Medical School

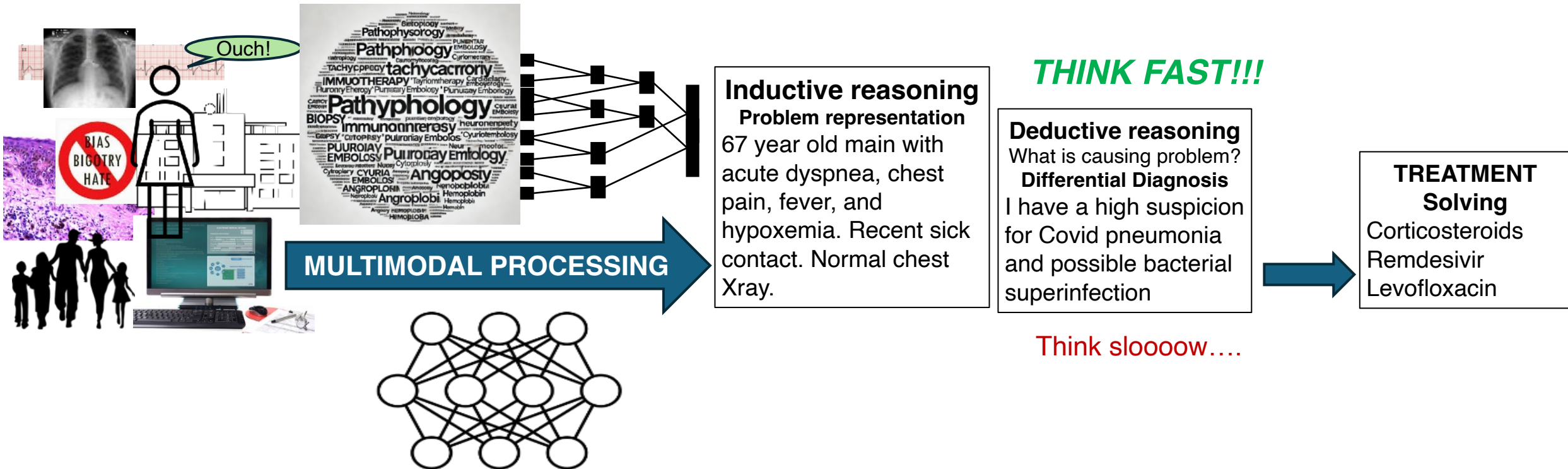


Disclosures



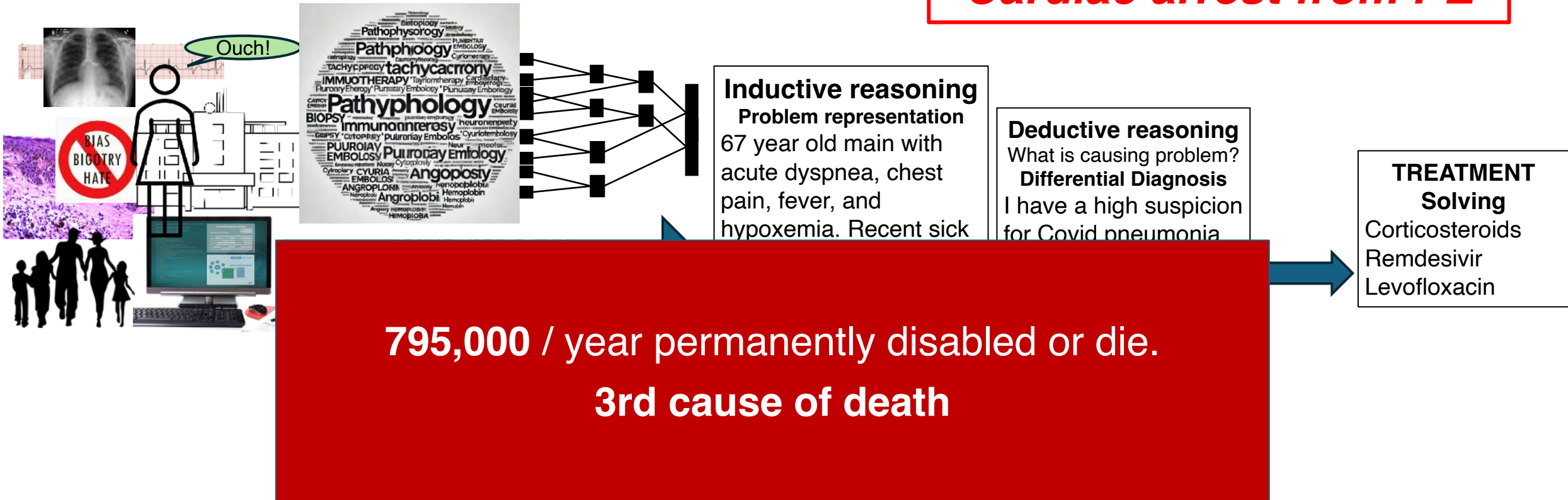
- MMS / NEJM Group
 - NEJM Healer
 - NEJM Healer was acquired by Lecturio (no financial COI)
- Lumeris, consultant
- I am NOT a computer scientist

Clinical Reasoning: How doctors think





Cardiac arrest from PE



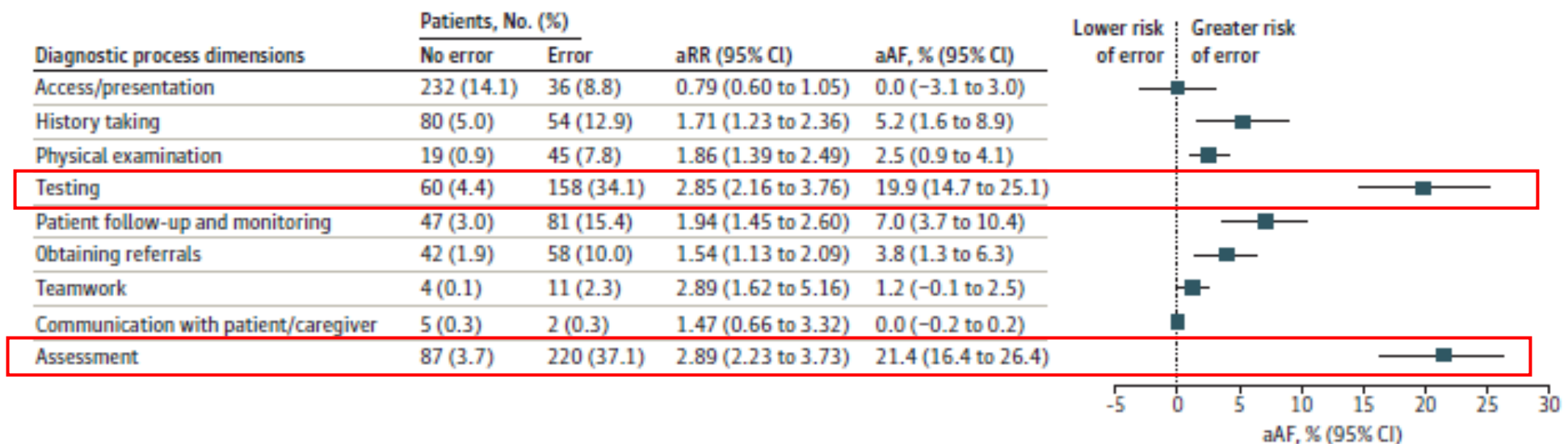
The Burden

Diagnostic Errors in Hospitalized Adults Who Died or Were Transferred to Intensive Care

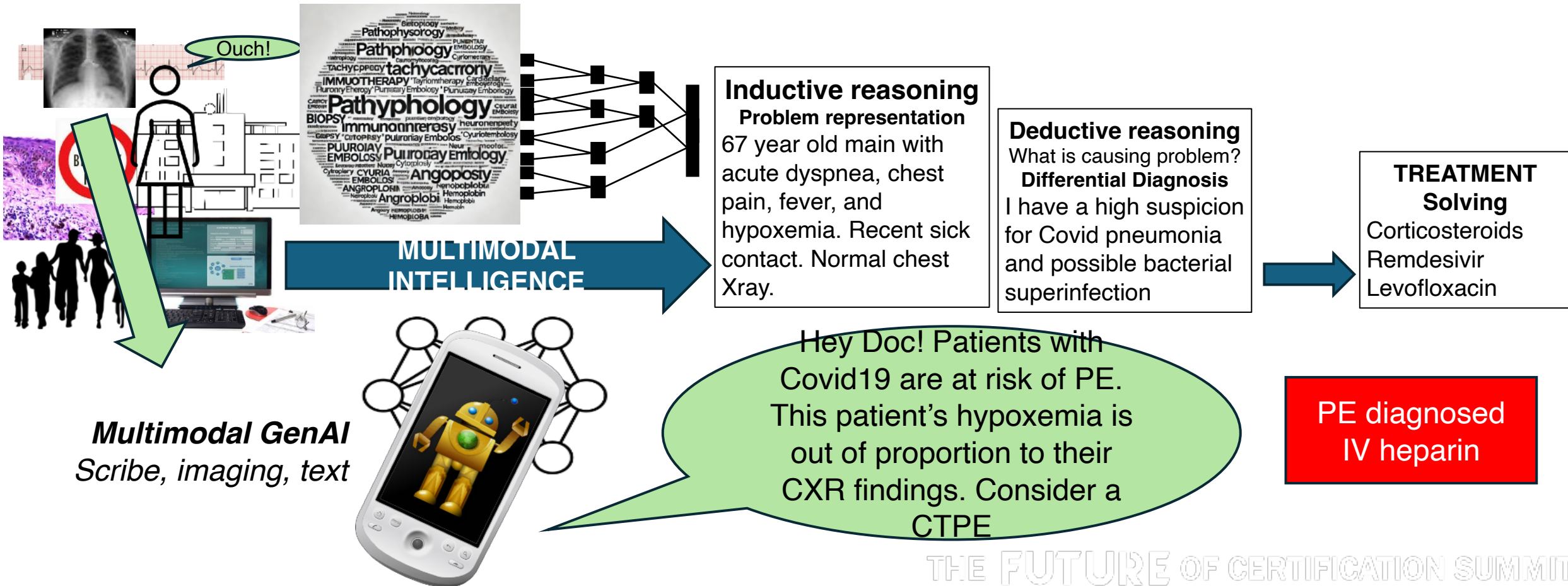
Andrew D. Auerbach, MD, MPH; Tiffany M. Lee, BA; Colin C. Hubbard, PhD; Sumant R. Ranji, MD; Katie Raffel, MD; Gilmer Valdes, PhD, DABR; John Boscardin, PhD; Anuj K. Dalal, MD; Alyssa Harris, MPH; Ellen Flynn, RN, MBA, JD; Jeffrey L. Schnipper, MD, MPH; for the UPSIDE Research Group

- 550 (23.0%) out of 2428 experienced diagnostic error
- 22.7% of inpatient patients who died or transferred to the ICU experienced diagnostic error

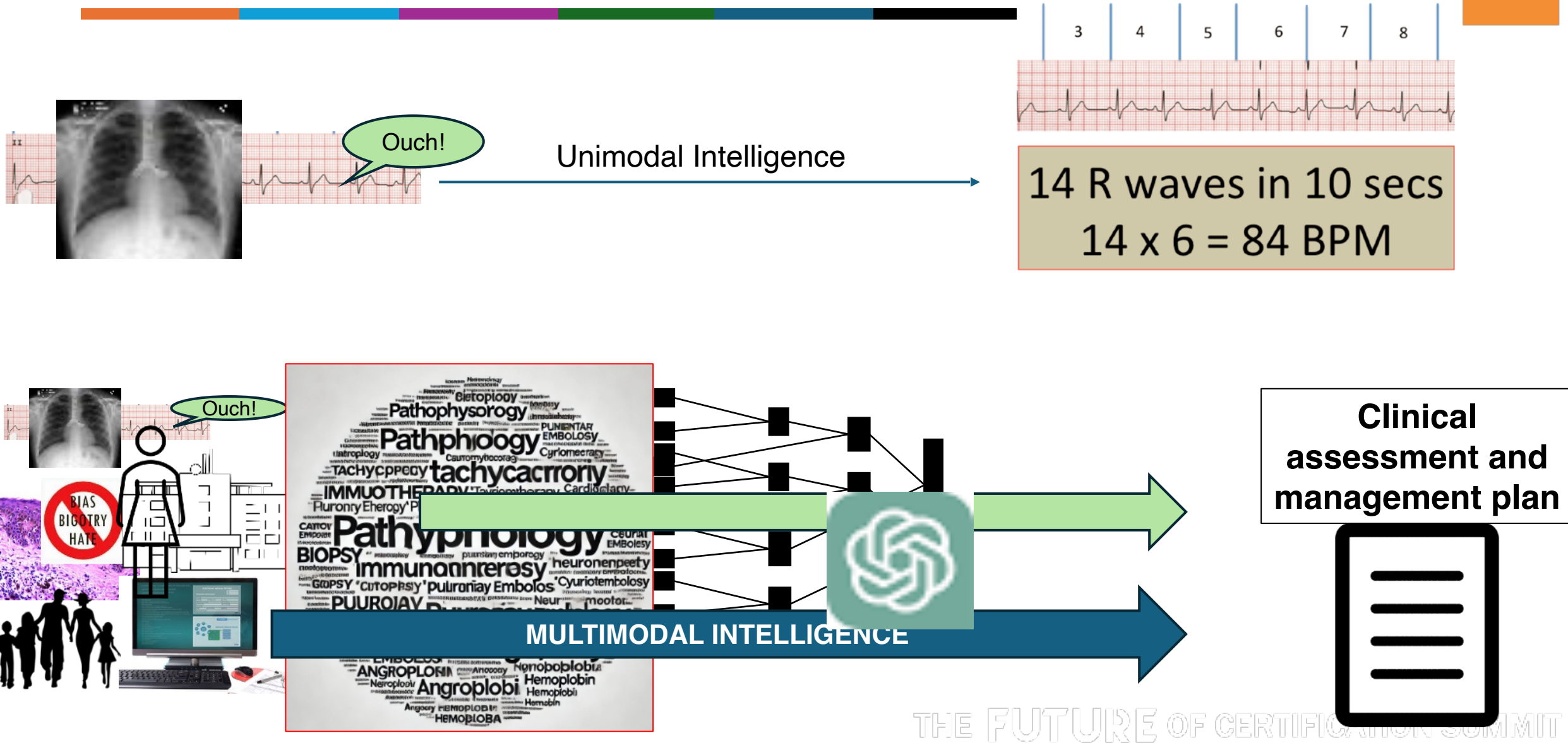
Figure 2. DEER Process Fault Dimensions: Prevalence, Adjusted Associations With Diagnostic Errors, and Adjusted Attributable Fractions (aAFs) (N = 2428)



Artificial Clinical Intelligence



Generative AI: A paradigm shift



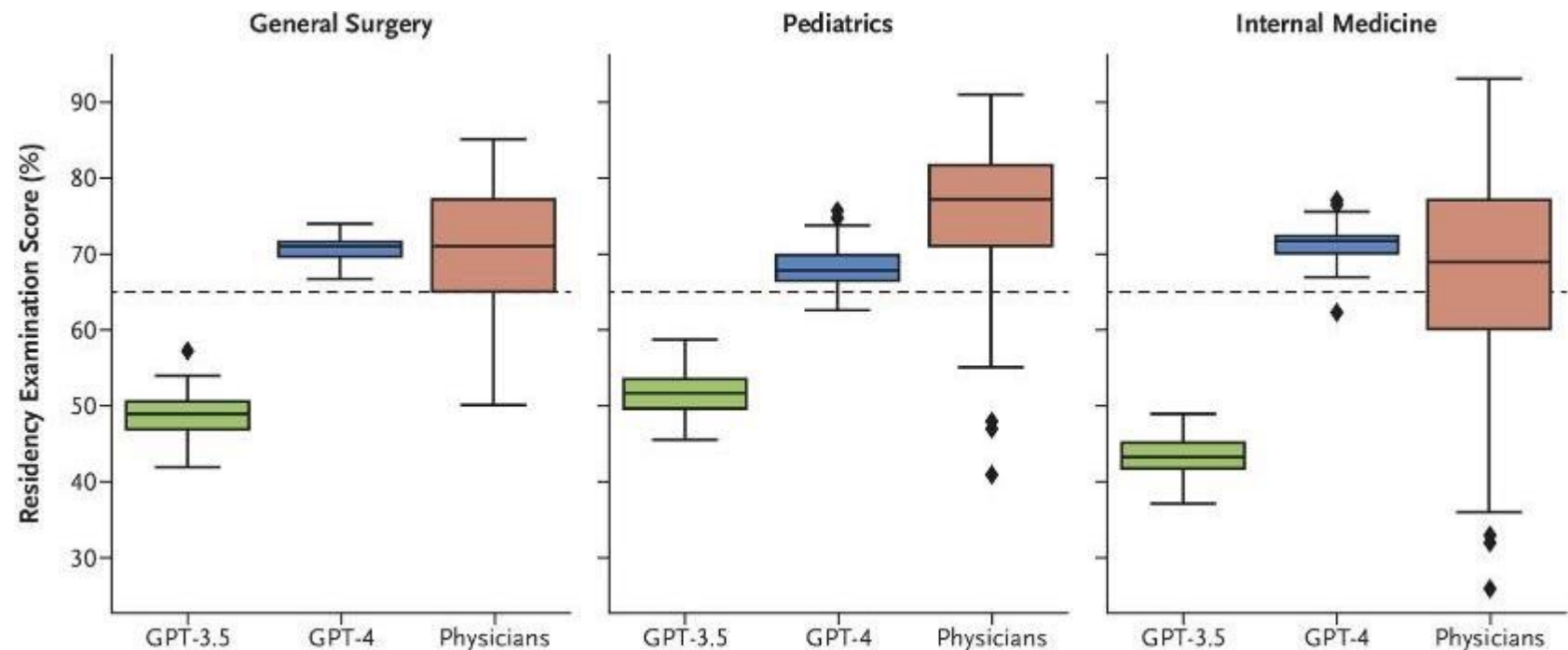
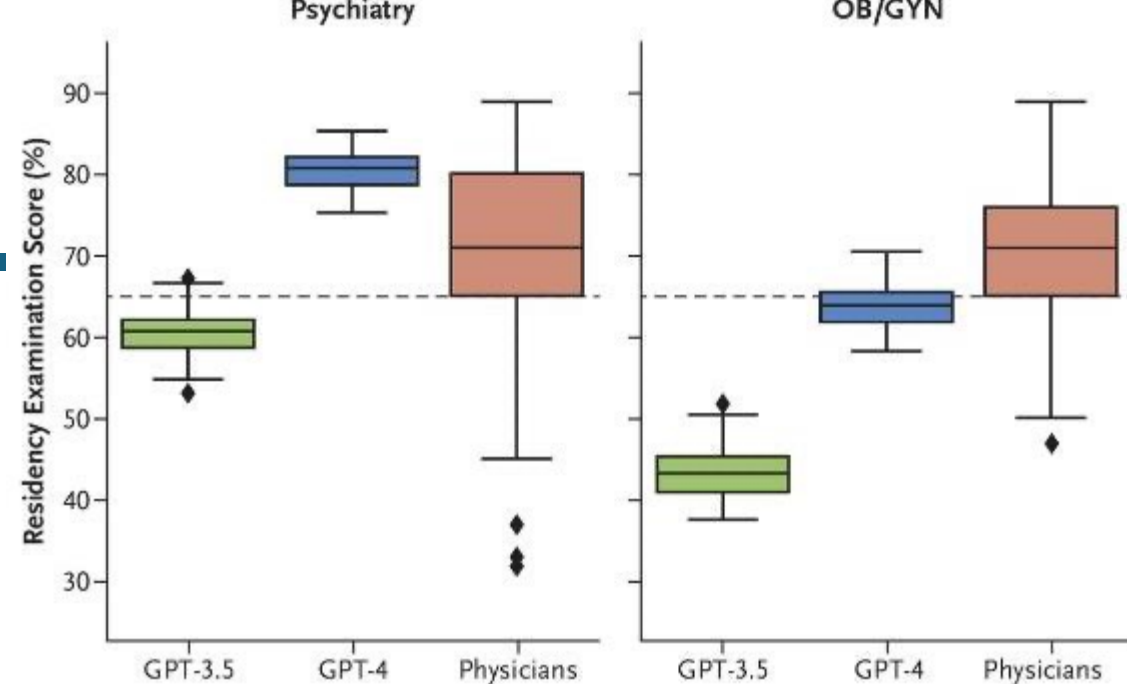
GPT vs Physicians on Board Exams

retrospective analysis of physicians' performance on the **2022 Israeli medical board certification examinations** across five core medical specialties.

Compared 849 physicians with GPT

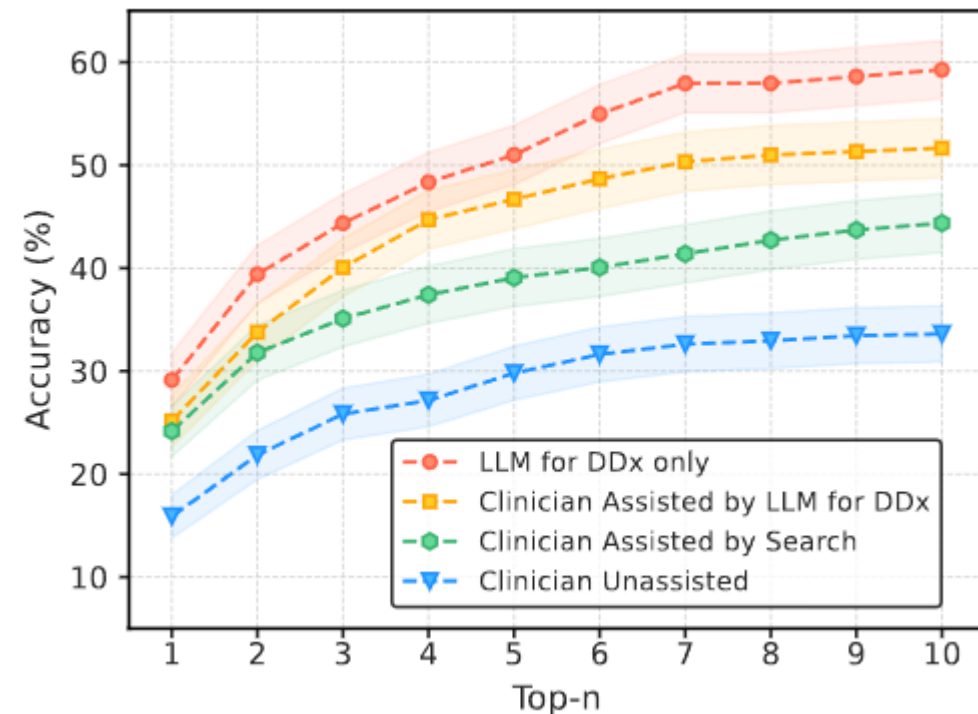
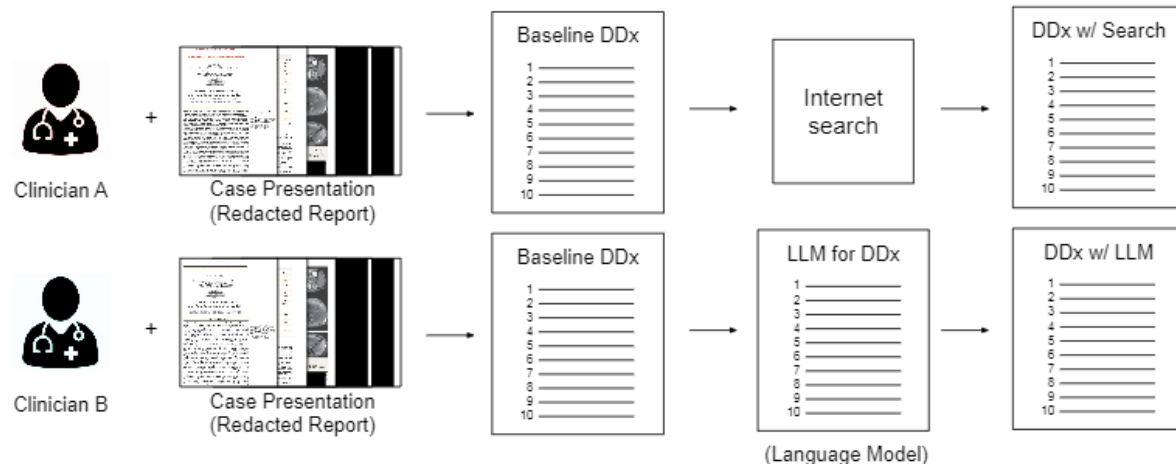
Accounted for model stochasticity by GPT model on 120 attempts

Katz et al. NEJM AI, 2024



LLMs Can Solve Case Reports

- 358 NEJM CPCs, including 56 not included in training data, using a fine-tuned Palm2 compared to human clinicians.



LLMs Express Clinical Reasoning

- Residents, attending, and GPT-4 solving NEJM Healer cases – 236 sections in total
- Assessed expression of reasoning process with r-IDEA

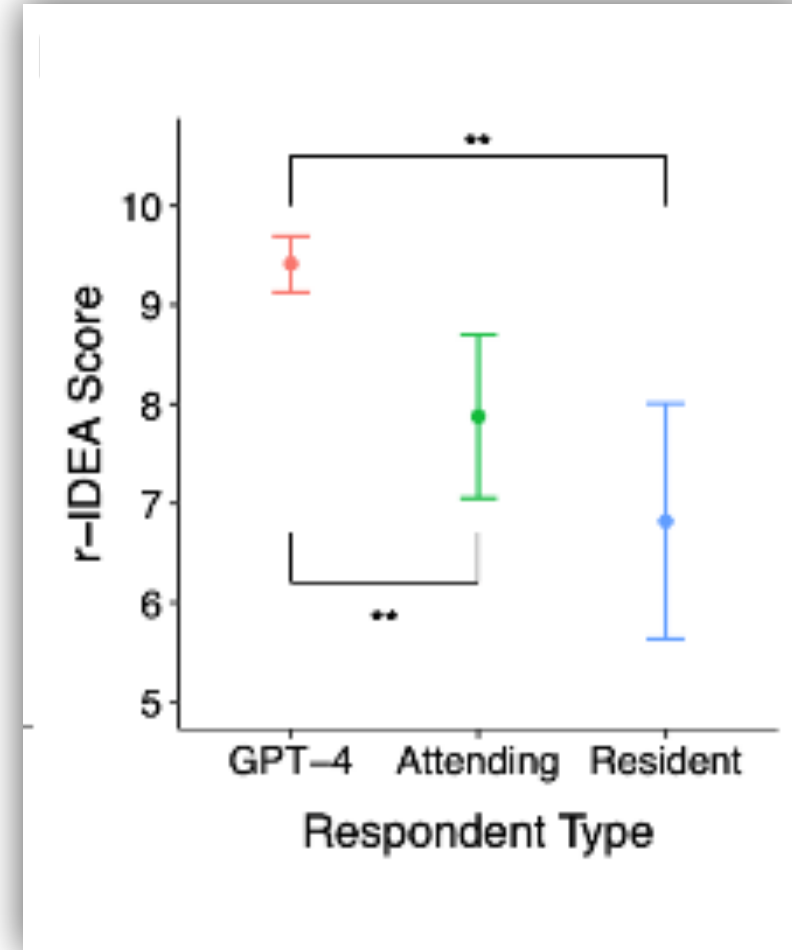
JGIM Schaye et al.: Development of a Clinical Reasoning Assessment Tool 509

Domain	Description	Assessment	Points
I - Interpretive Summary	Provides a concise summary statement that uses semantic vocabulary to highlight the most important elements from history, exam, and testing and to interpret and represent the patient's main problem(s). The presence or absence of the following features is assessed: a) Key risk factors; b) Chief complaint; c) Illness time course; and d) Use of semantic qualifiers (e.g. monoarticular vs polyarticular) or unified medical concepts (e.g. volume overload, cardiovascular risk factors). NB: Some problems have an implied time course (e.g. syncope, seizure).	No features present	0
		1 feature present	1
		2 features present	2
		3 features present	3
		4 features present	4
D - Differential Diagnosis	Offers more than one relevant diagnostic possibility, committing to what is most likely and considering what is less likely or unlikely yet important to consider for the main chief complaint. If the chief complaint is a diagnosis or syndrome (e.g. acute on chronic systolic heart failure) then differential to rate may be around the differential for that exacerbation (e.g. medication non-compliance vs. arrhythmia).	No differential	0
		Differential is implicitly stated, given as a diagnostic category (e.g. "cardiac"), OR implicitly prioritized	1
		Differential is explicitly stated AND explicitly prioritized	2
E - Explanation of Lead Diagnosis	Explains the reasoning behind the lead diagnosis, including the epidemiology and key features and how these compare with the patient's presentation. If objective data points are not clearly linked to the lead diagnosis or alternative diagnosis, then only designate points to lead OR alternative diagnosis and NOT both.	No explanation	0
		1 objective data point in explanation of lead diagnosis	1
		≥2 objective data points in explanation of lead diagnosis	2
A - Alternative Diagnosis Explained	Explains the reasoning behind alternative diagnoses, including the epidemiology and key features and how these compare with the patient's presentation and alternative diagnosis. If objective data points are not clearly linked to the lead diagnosis or alternative diagnosis, then only designate points to lead OR alternative diagnosis and NOT both.	No explanation for any alternative diagnosis	0
		1 objective data point in explanation of at least one alternative diagnosis	1
		≥2 objective data points in explanation of at least one alternative diagnosis	2
Revised-IDEA Score	Overall evaluation of demonstration of clinical reasoning in the assessment section of admission notes.	Sum of I + D + E + A points (score ≥6 indicates high-quality clinical reasoning documentation)	0-10

Figure 1 The Revised-IDEA assessment tool for clinical reasoning documentation.

LLMs Express Clinical Reasoning

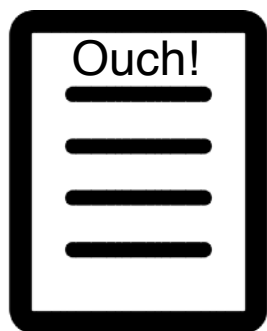
- Residents, attending, and GPT-4 solving NEJM Healer cases – 236 sections in total
- Assessed expression of reasoning process with r-IDEA
- GPT-4 had significantly higher r-IDEA scores (9.41 vs 7.83 for attendings and 6.82 for residents)
- No difference in efficiency, accuracy, quality, cannot miss
- Increase of incorrect reasoning (12% vs 3%), though all minor examples



Patient Triage In The ED

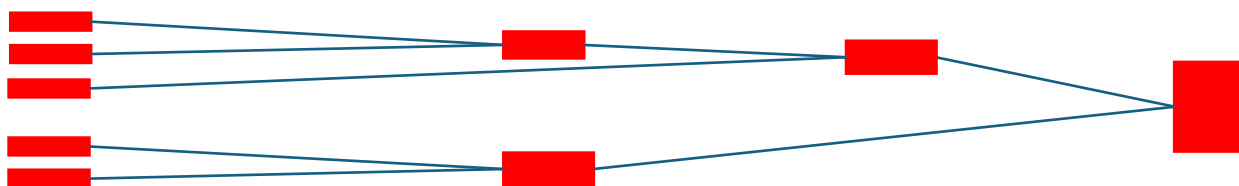


“You are an Emergency Department physician. Below are the symptoms of a patient presenting to the Emergency Department.



1'000
patients

Clinical
history
and
physical
exam only



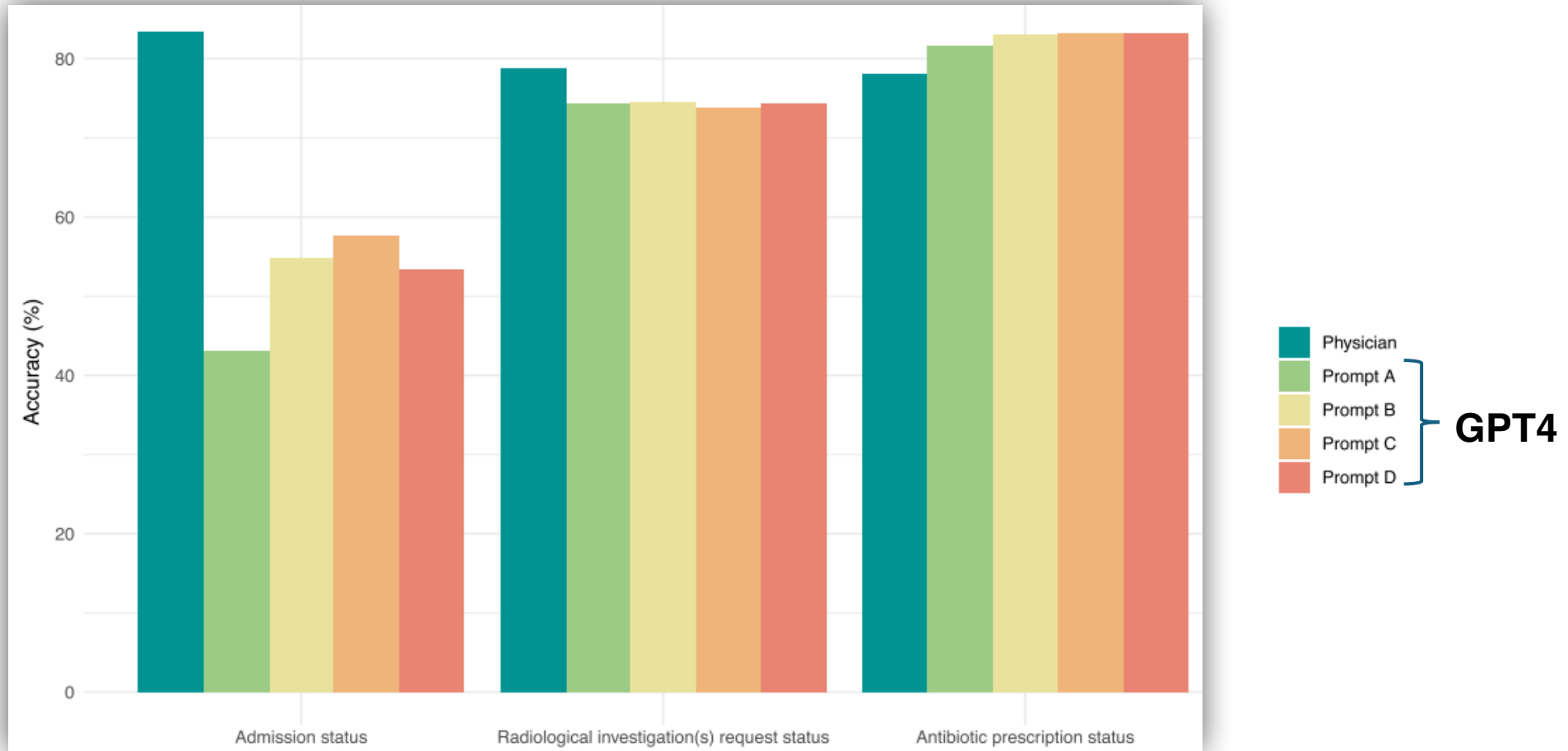
Please return whether the patient should be admitted to hospital.



Please return whether the patient requires radiological investigation (e.g X-ray, ultrasound scan, CT scan or MRI scan)

Please return whether the patient requires antibiotics

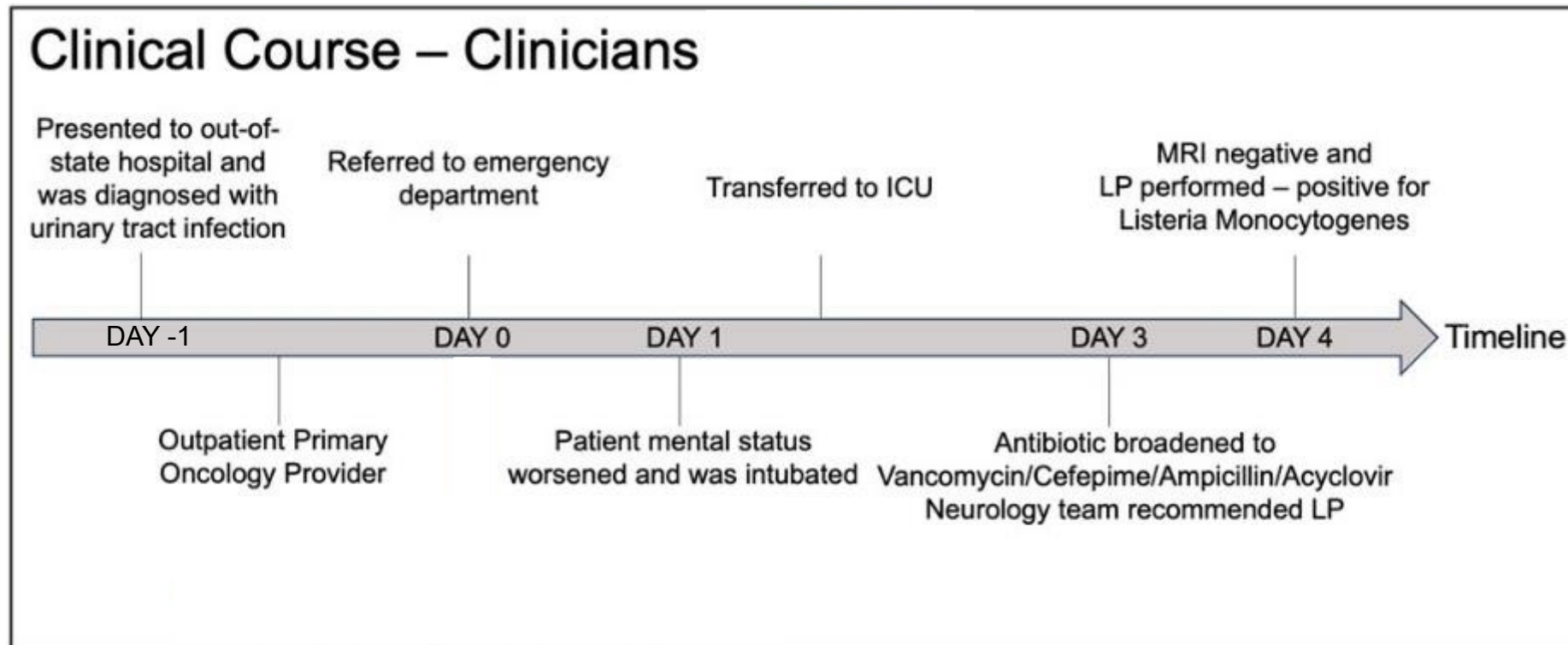
Patient Triage In The ED



Williams, C. Y. K., Miao, B. Y., Kornblith, A. E., & Butte, A. J. (2024).
Nature Communications 2024 15:1, 15(1)

A Day In The ICU

61-year-old female with a history of hypertension, breast cancer, and recurrent ovarian cancer, presenting with a recent history of nausea, poor oral intake, and altered mental status following a ferry ride



“I want you to be the consultant on the team. I will provide you the history of present illness with the vitals, physical examination findings, labs, and imaging data. Following that, I want you to provide me with 1) a brief summary of the patient, 2) 10 Differential diagnosis (ranked by percentage likelihood), 3) further diagnostic tests you would obtain and 4) an initial management plan.”

Evaluation and mitigation of the limitations of large language models in clinical decision-making

Received: 26 January 2024

Accepted: 29 May 2024

Published online: 04 July 2024

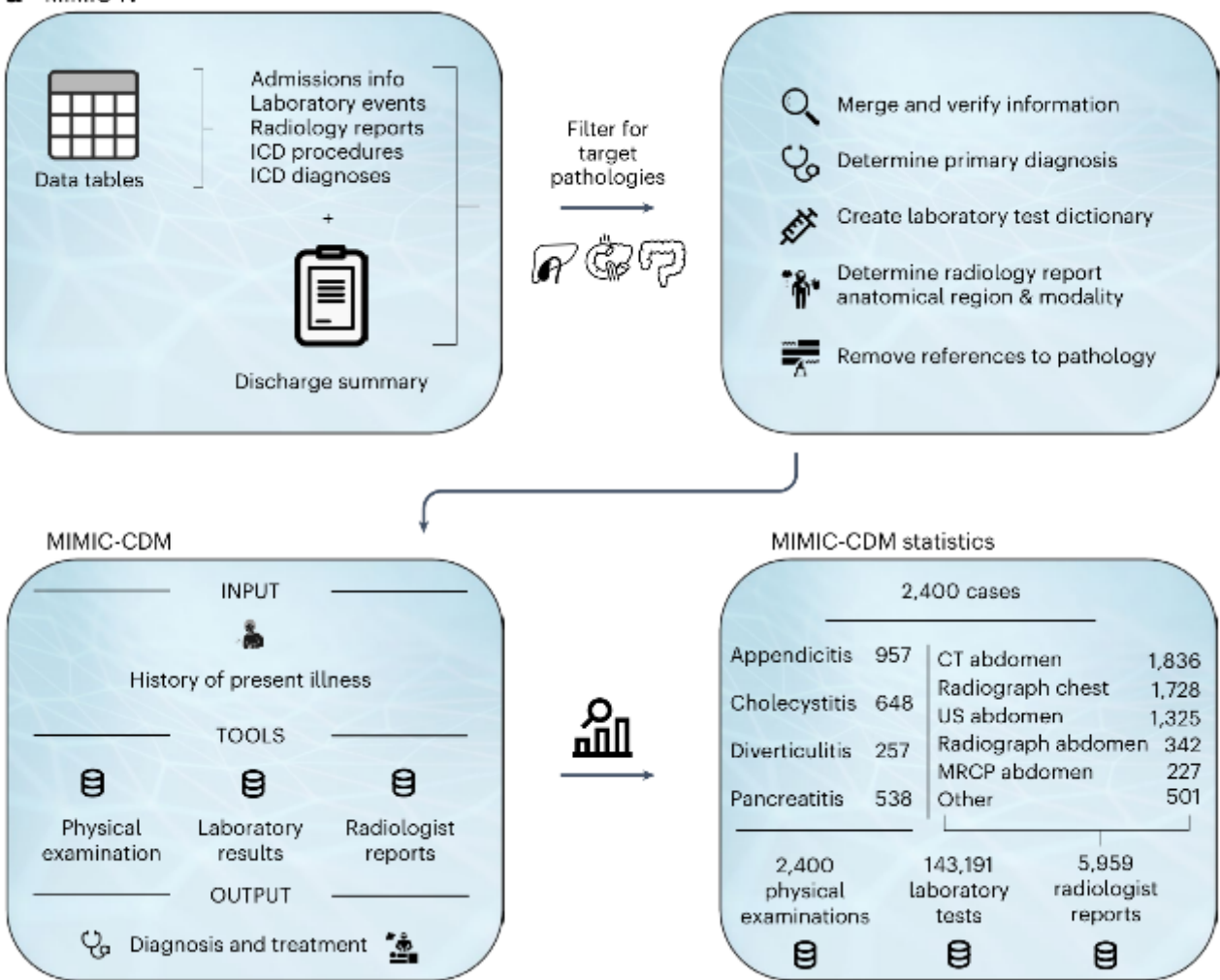
Paul Hager^{1,2,8}, Friederike Jungmann^{1,2,8}, Robbie Holland³, Kunal Bhagat⁴, Inga Hubrecht⁵, Manuel Knauer⁵, Jakob Vielhauer⁶, Marcus Makowski², Rickmer Braren^{2,9}, Georgios Kaissis^{1,2,3,9} & Daniel Rueckert^{1,3,9}

Table 1 | An overview of the considered LLMs and their properties

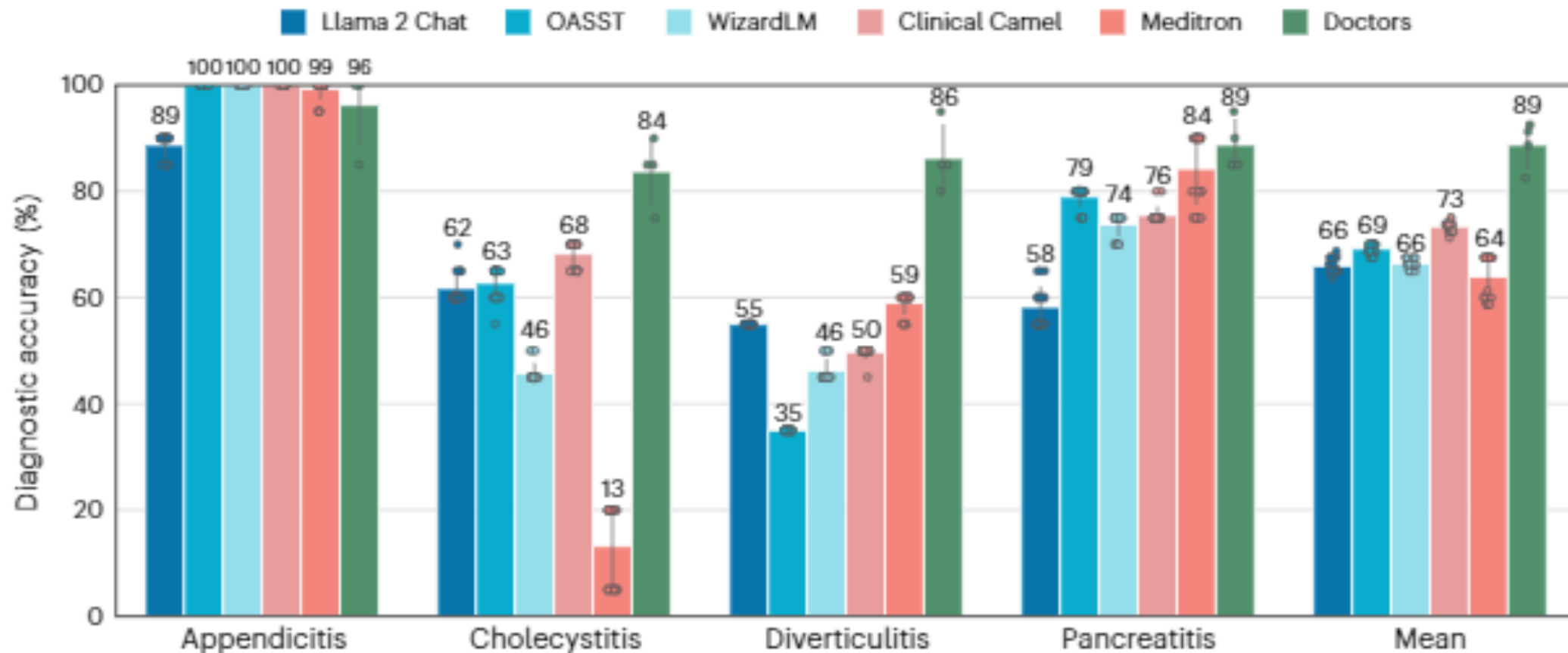
Model	Base	Parameters	Training dataset
Llama 2 Chat ³²	Llama 2 (ref. 32)	70B	Public data ^a
OASST ³³	Llama 2 (ref. 32)	70B	Public data ^a , https://huggingface.co/OASST1 , open-source data
WizardLM ³⁴	Llama 2 (ref. 32)	70B	Public data ^a , Evol-Instruct generated ³⁴
Clinical Camel ¹⁹	Llama 2 (ref. 32)	70B	Public data ^a , https://sharegpt.com/ ; ShareGPT ¹⁹ , MedQA ¹³
Meditron ³⁵	Llama 2 (ref. 32)	70B	Public data ^a , https://huggingface.co/datasets/medtronic , guidelines, public PubMed abstracts ³⁵
Chat-GPT ³⁶	GPT3.5 (ref. 60)	???	User conversations ^b , Common Crawl ⁶¹ , Books2 (ref. 63), Wikipedia
GPT-4 (ref. 64)	???	???	???
Med-PaLM ⁹	Flan-PaLM ⁶²	540B	Webpages ^b , Wikipedia ^b , social media ^b , 473 instruction fine-tuning datasets ⁶⁵ , HealthSearchQA ⁹ , MedicationQA ⁶⁶ , LiveQA ⁶⁷
Med-PaLM 2 (ref. 8)	PaLM 2 (ref. 68)	340B	Web Documents ^b , books ^b , code ^b , mathematics ^b , conversational data ^b , MedQA ¹³ , HealthSearchQA ⁹ , MedicationQA ⁶⁶ , LiveQA ⁶⁷

65K ICU patients → 2400 ICU patients with 4 common Dx

a MIMIC-IV

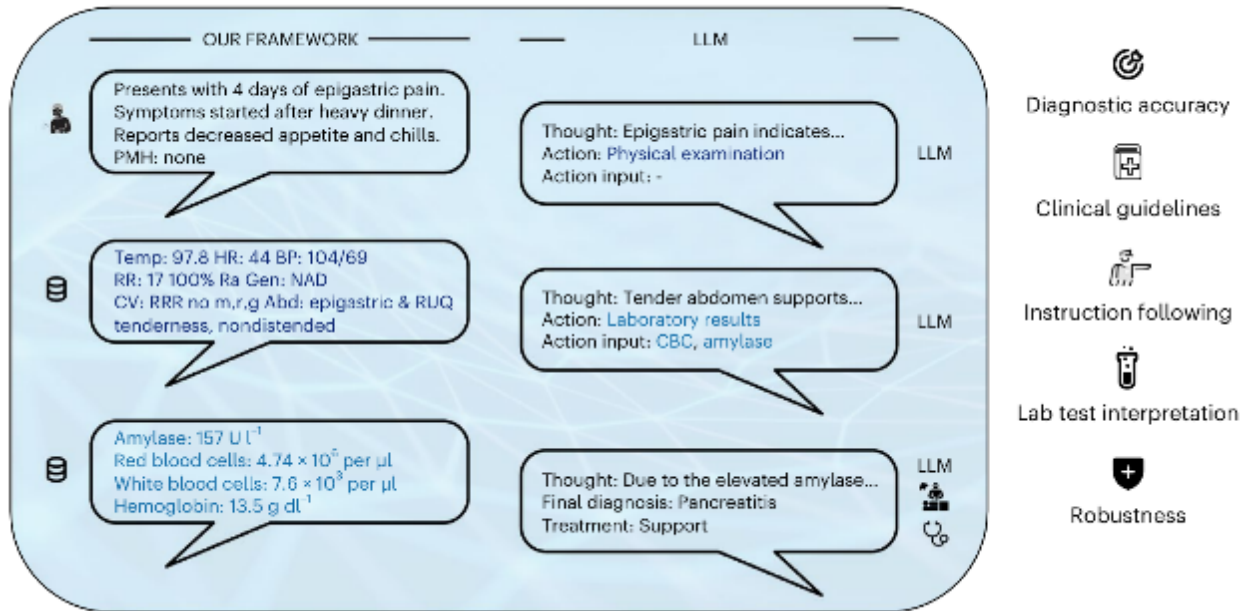


Evaluation of LLMs



When provided with all information on a subset ($n=80$), LLMs diagnose significantly worse than doctors.
Mean diagnostic accuracy of LLMs over multiple seeds ($n=20$) compared to clinicians ($n=4$)

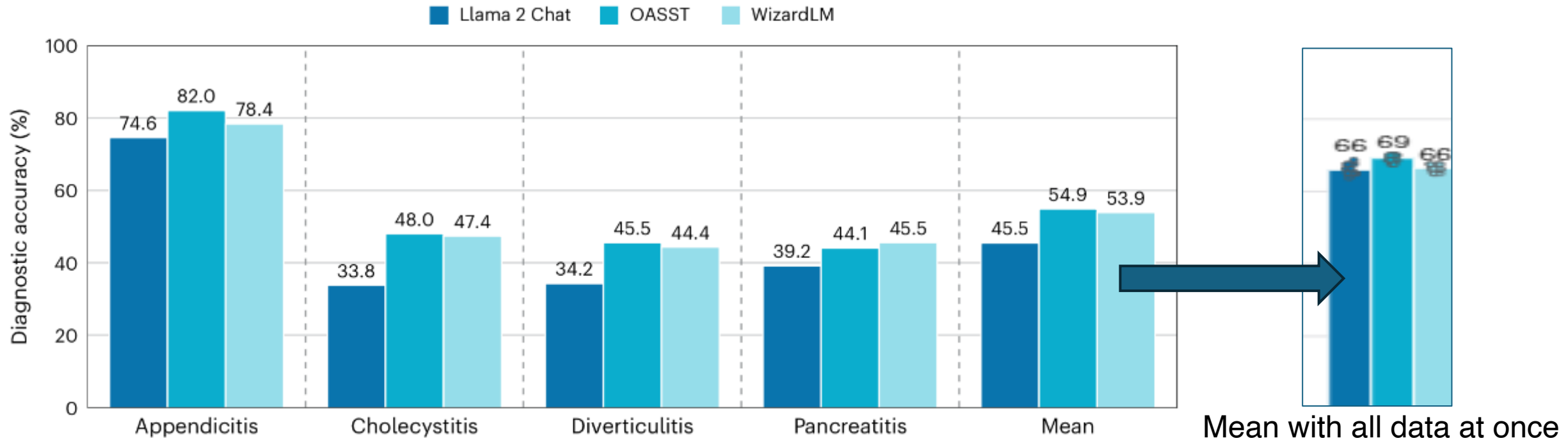
b EVALUATION



Worse performance if autonomous (20% loss)

- Inability to consistently follow clinical guidelines.
- Struggles with interpreting laboratory results.
- Sensitivity to the order and amount of information provided.

minor changes in instructions can greatly change diagnostic accuracy such as asking for the 'main diagnosis' or 'primary diagnosis' instead of 'final diagnosis'



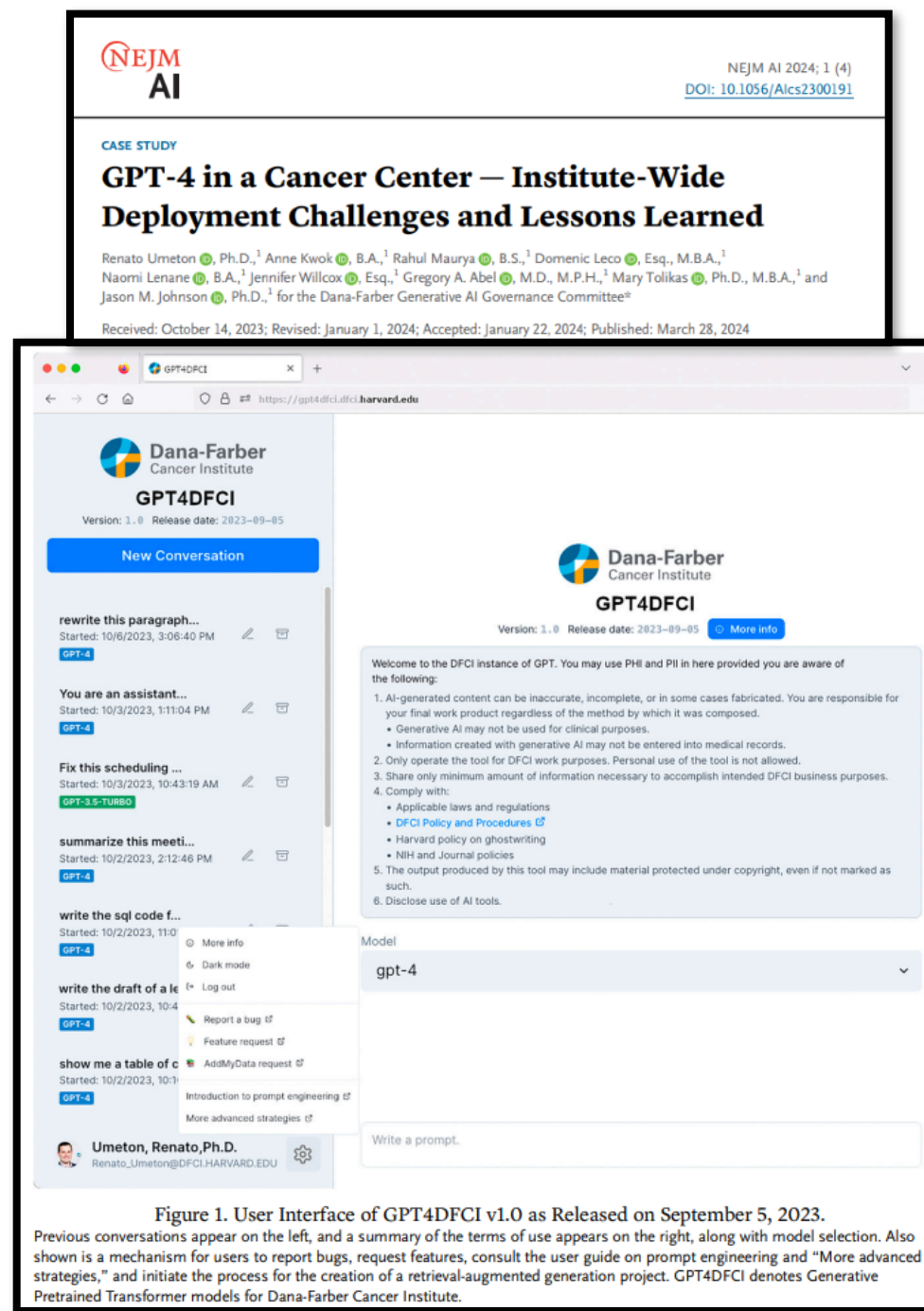
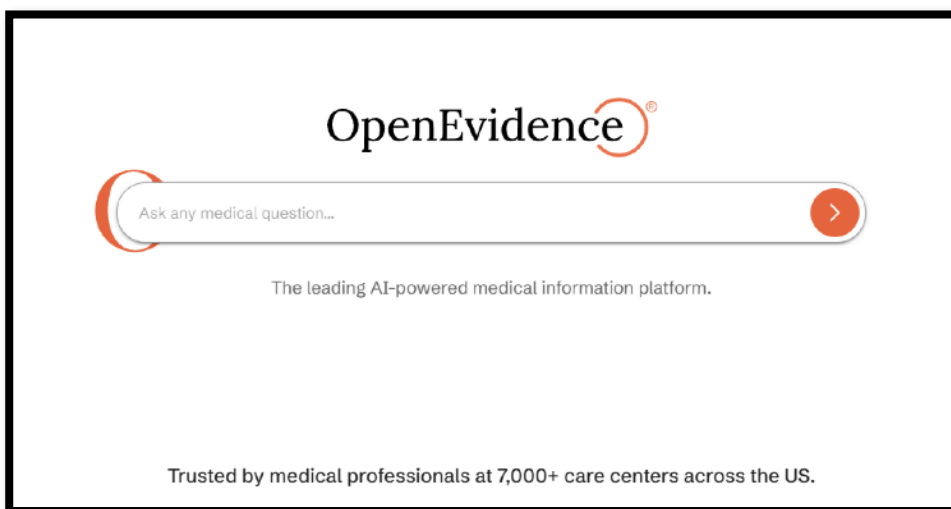
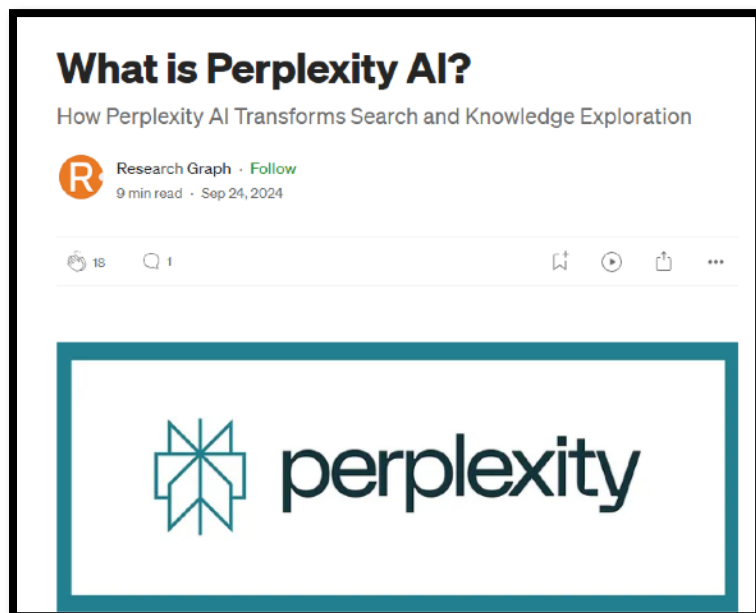
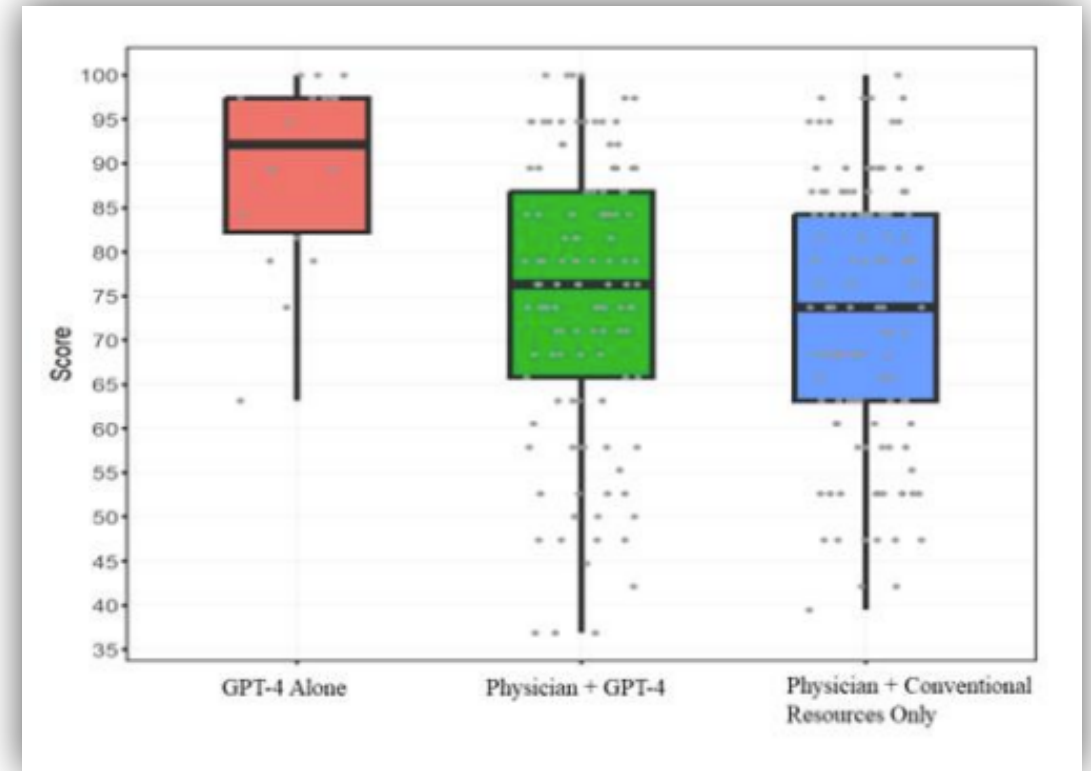
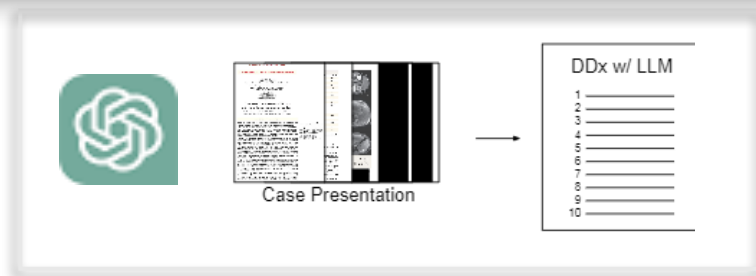
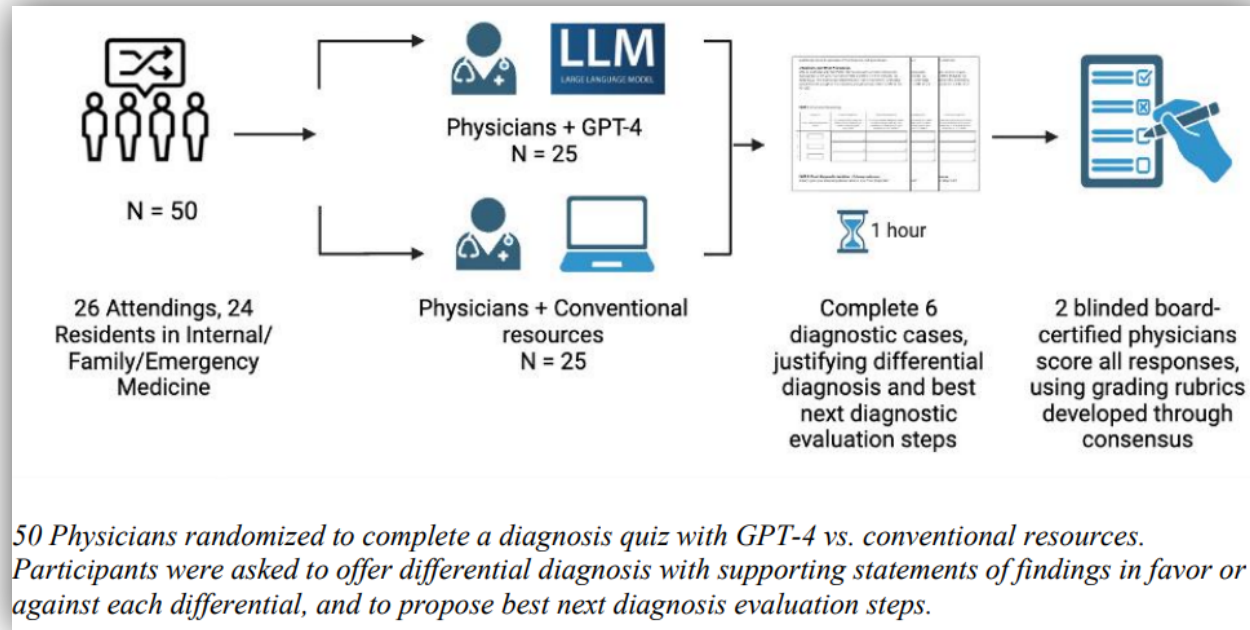


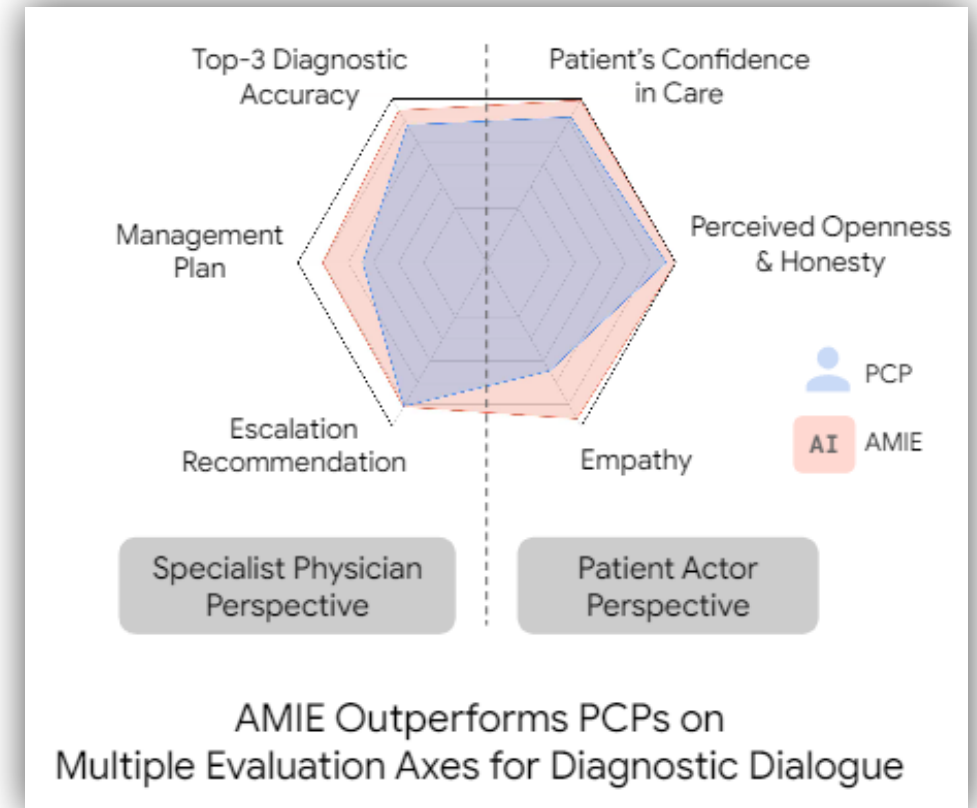
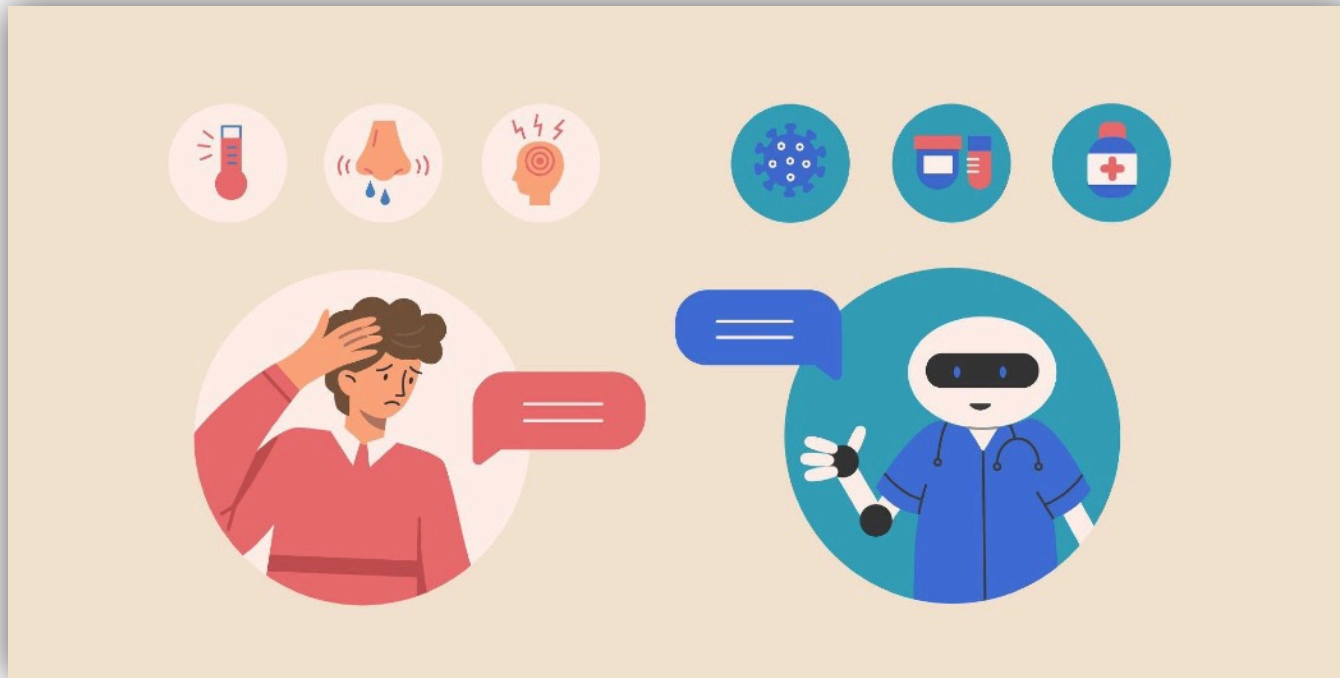
Figure 1. User Interface of GPT4DFCI v1.0 as Released on September 5, 2023. Previous conversations appear on the left, and a summary of the terms of use appears on the right, along with model selection. Also shown is a mechanism for users to report bugs, request features, consult the user guide on prompt engineering and “More advanced strategies,” and initiate the process for the creation of a retrieval-augmented generation project. GPT4DFCI denotes Generative Pretrained Transformer models for Dana-Farber Cancer Institute.

LLM Influence on Diagnostic Reasoning



Distribution of Diagnostic Performance Scores

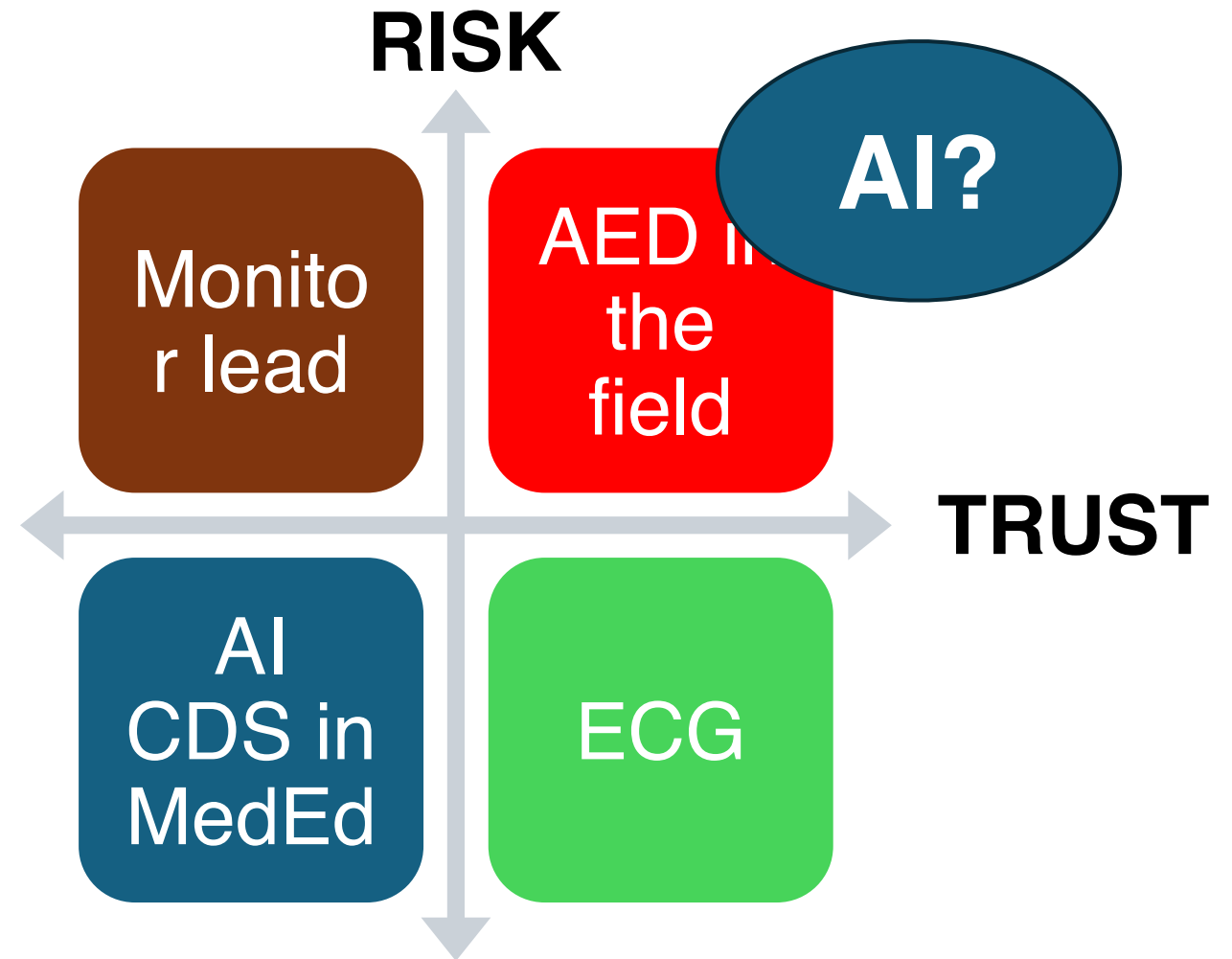
Conversational Diagnostic AI



Tu T, Palepu A, Schaekermann M, Saab K, Freyberg J, Tanno R, *et al.* Towards Conversational Diagnostic AI. ArXiv 2024.

Entrustment In Patient Care

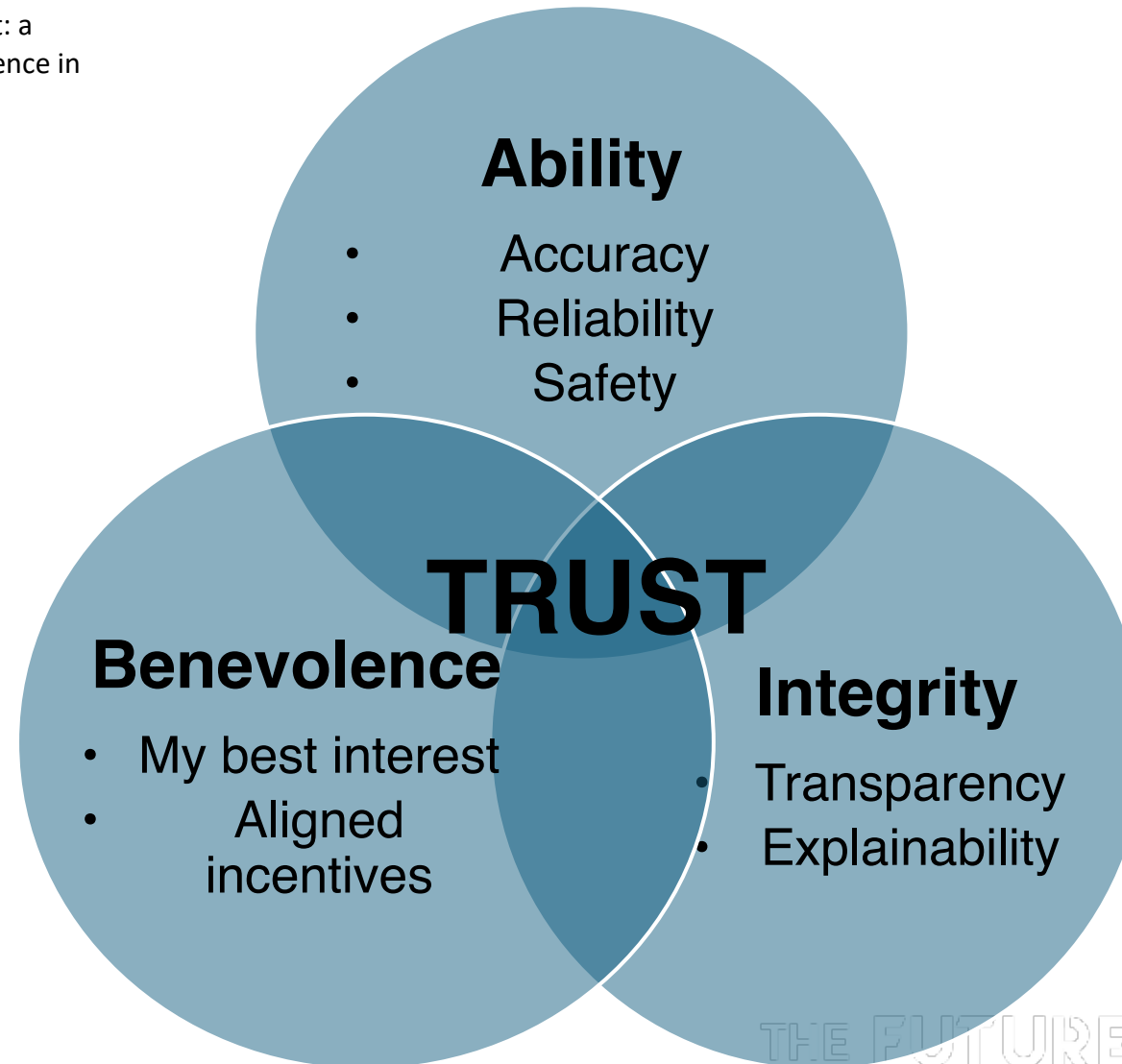
Is this patient having ventricular tachycardia?



Entrustment In Patient Care



Brian C. Gin, M., PhD, et al. (In press). "Entrustment: a framework to safeguard the use of artificial intelligence in health professions education." [Acad Med.](#)



Is AI safe?

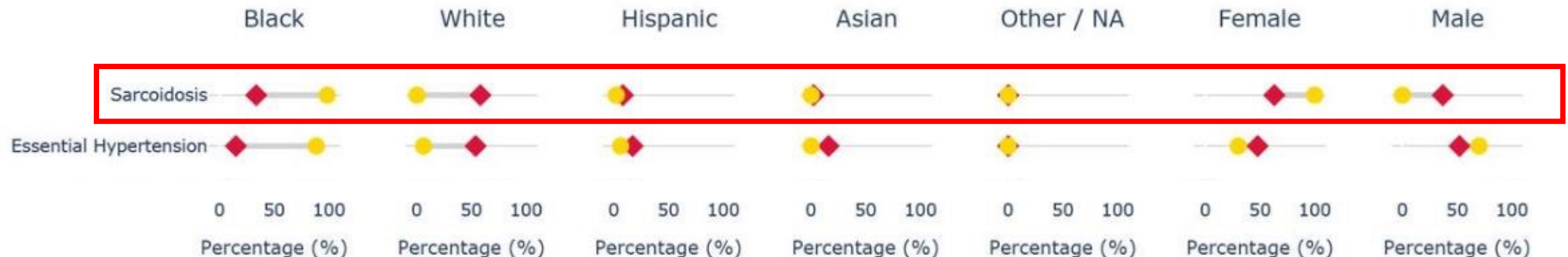


LLMs contain the bias of their training



- Asked GPT-4 to create clinical vignettes
 - Over-represented demographic stereotypes of diseases
- Asked GPT-4 to give management plans for cases while substituting gender and race/ethnicity
 - Less likely to recommend advanced imaging for Black patients compared to White patients

GPT-4-Estimated and True Patient Demographic Distribution of Patients with Each Condition



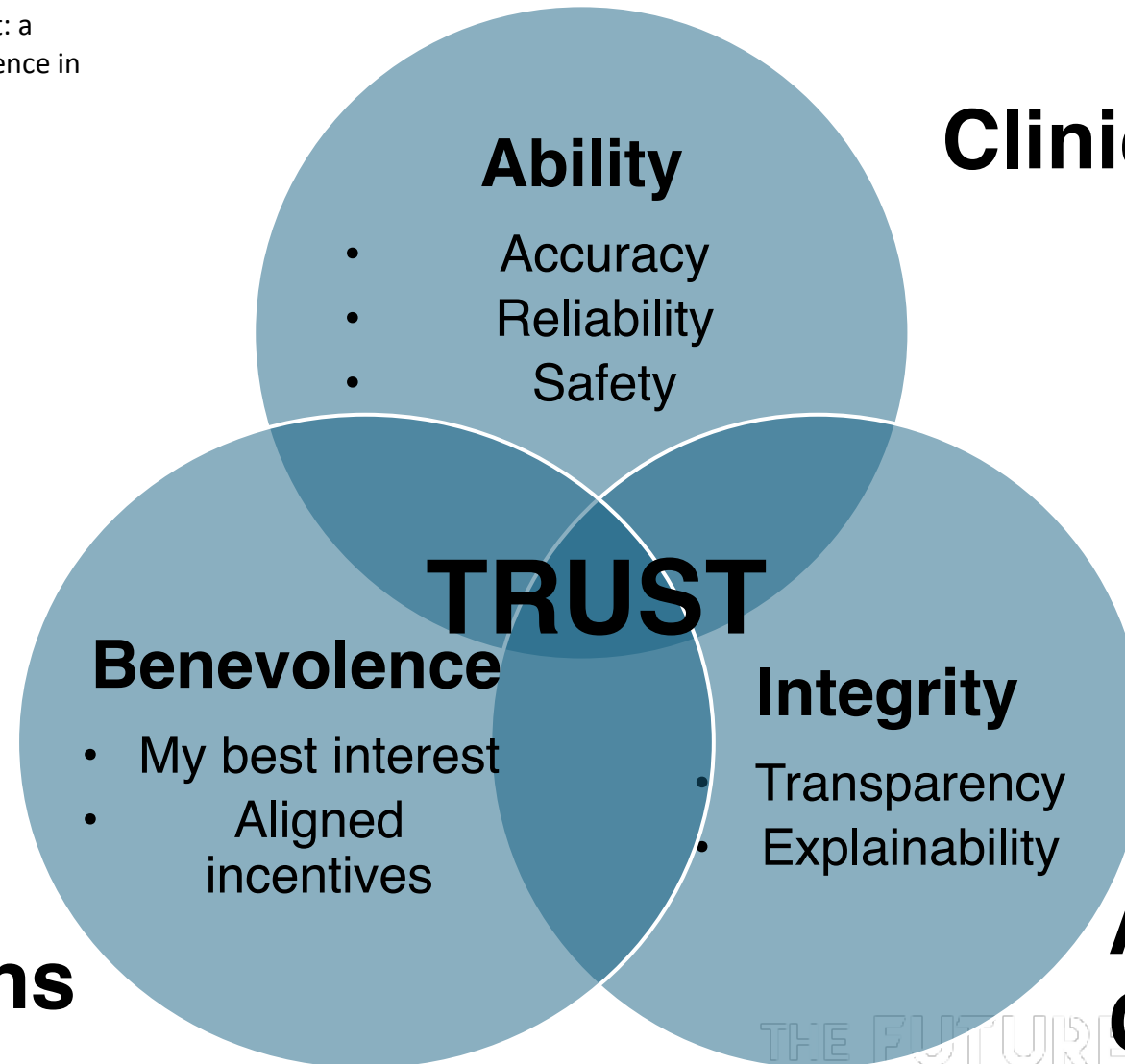
Legend: ◆ True ● GPT-4 Estimated

Entrustment In Patient Care



Brian C. Gin, M., PhD, et al. (In press). "Entrustment: a framework to safeguard the use of artificial intelligence in health professions education." [Acad Med.](#)

**Regulation by
third-party (e.g.,
FDA) and
patients/clinicians**



Clinical trials

**AI report cards
Chain-of-Thought**

Do you want more?



- **Option1:** Thank you!!!
- **Option2:** Next slide.

Clinical Trials

- Summarize existing evaluations of LLMs in health care
- A systematic search of PubMed and Web of Science was performed for studies published between January 1, 2022, and February 19, 2024
- 519 studies reviewed, published between January 1, 2022, and February 19, 2024

Health care tasks							
Enhancing medical knowledge	222	91	44	33	16	10	3
Making diagnoses	100	38	11	11	14	4	0
Educating patients	88	68	32	22	18	3	2
Making treatment recommendations	47	22	9	8	3	1	0
Communicating with patients	35	29	8	15	22	1	0
Care coordination and planning	36	24	4	5	7	1	0
Triaging patients	24	7	5	2	8	8	0
Carrying out a literature review	18	7	3	2	2	2	0
Synthesizing data for research	16	7	2	3	2	2	0
Generating medical reports	8	8	2	0	3	0	0
Conducting medical research	8	7	3	3	3	0	0
Providing asynchronous care	8	5	3	3	1	1	0
Managing clinical knowledge	5	5	1	1	0	0	0
Clinical note-taking	6	2	1	1	0	0	1
Generating clinical referrals	3	0	0	0	0	0	0
Enhancing surgical operations	3	3	1	1	0	0	0
Biomedical data mining	2	0	0	0	0	0	0
Generating billing codes	1	0	0	0	0	0	0
Writing prescriptions	1	0	0	0	0	0	0
NLP and NLU tasks							
Question answering	398	194	71	61	54	14	5
Text classification	29	10	6	5	10	2	0
Information extraction	29	12	8	5	4	6	0
Summarization	29	21	7	3	8	0	1
Conversational dialogue	6	6	1	1	5	1	0
Translation	5	1	2	2	1	2	0
Accuracy				Dimension of evaluation			
Comprehensiveness				Fairness, bias, and toxicity evaluation			
Factuality				Deployment metrics			
Robustness				Calibration and uncertainty			

Bedi, S., et al. (2024). "Testing and Evaluation of Health Care Applications of Large Language Models." JAMA.

Chain-of-Thought prompting



What is the cause of this patient's respiratory failure?

This patient has ARDS



Explain your reasoning

Because their symptoms started three months ago and they like Candy



Chain-of-Thought prompting



What is the cause of this patient's respiratory failure?
Explain your reasoning.

This patient may have ILD because of the chronicity of their symptoms and smoking Hx. ARDS would occur more acutely





How many rs are in "strawberry?"



Which came first: the chicken or the egg?



Is a hot dog a sandwich?



Solve an advanced math problem



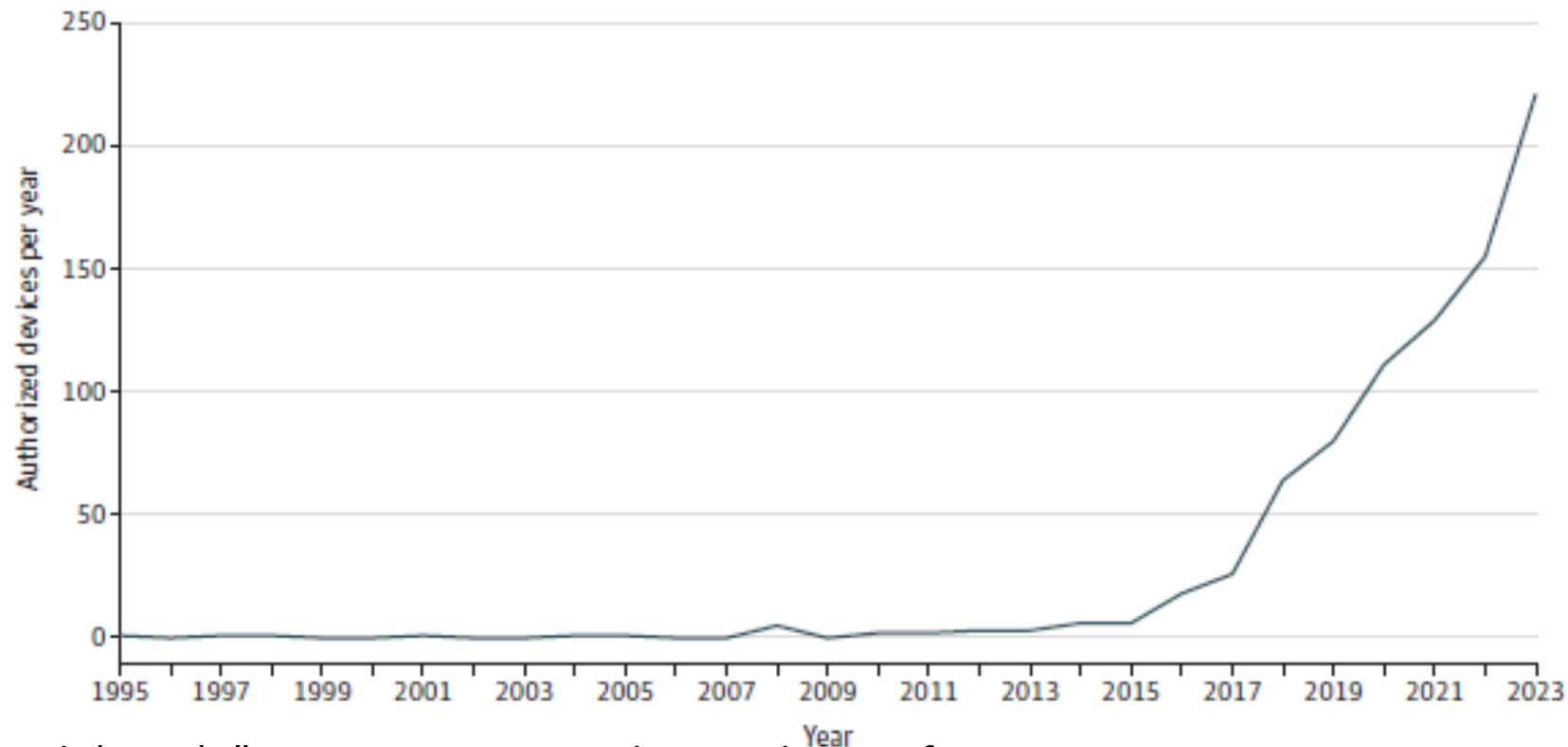
Do you want more?



- **Option1:** Thank you!!!
- **Option2:** Next slide.

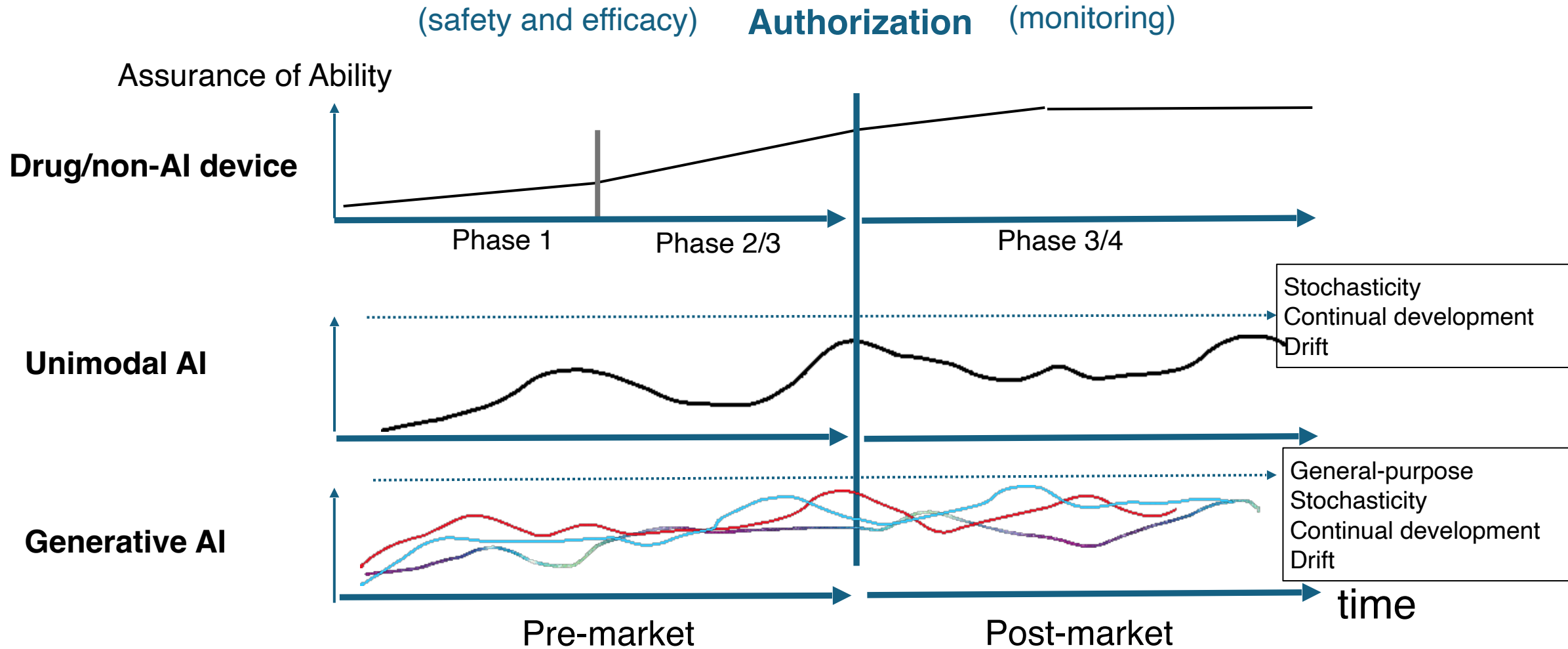
Artificial Intelligence–Enabled Medical Devices

Figure 1. Artificial Intelligence–Enabled Medical Devices Authorized for Marketing by the US Food and Drug Administration, by Year

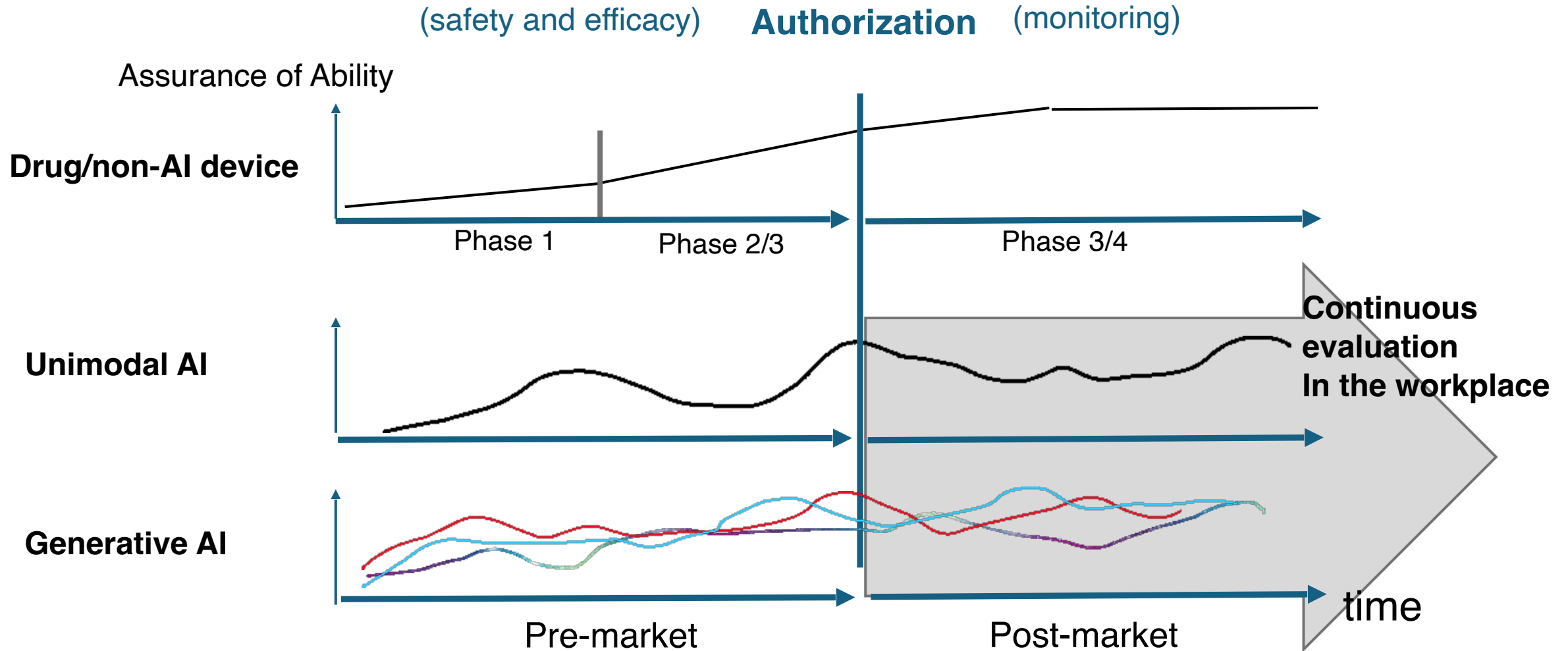


Warraich, H. J., et al. (2024). "FDA Perspective on the Regulation of Artificial Intelligence in Health Care and Biomedicine." JAMA

Life Cycle Regulation



Life Cycle Regulation



Conclusions



“Continuous complementary efforts to better understand how AI performs in the settings in which it is deployed. This will entail a comprehensive approach reaching far beyond the FDA, spanning the consumer and health care ecosystems to keep pace with accelerating technical progress.”

“Strong oversight by the FDA and other agencies aims to protect the long-term success of regulated products by maintaining a high grade of public trust in the regulated space.”

“Regulated industries, academia, and the FDA will need to develop and optimize the tools needed to assess the ongoing safety and effectiveness of AI in health care and biomedicine. The FDA will continue to play a central role with a focus on health outcomes, but all involved sectors will need to attend to AI with the care and rigor this potentially transformative technology merits.”

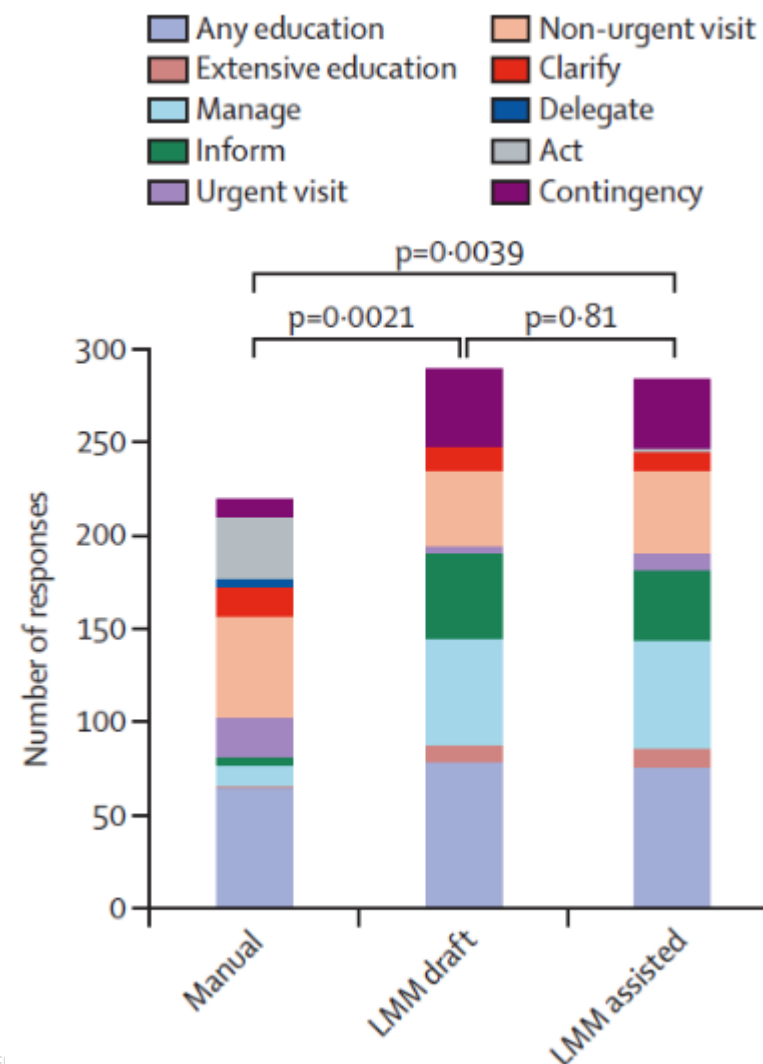
Do you want more?



- **Option1:** Thank you!!!
- **Option2:** Next slide.

Impact of using LLM to reply to patient messages

- The mean manual response (34 words) was **shorter** than the LLM draft (169 words) and LLM-assisted responses (160 words; $p < 0.0001$).
- The assessing physicians felt that the **LLM drafts posed a risk of severe harm in 7.1%** of survey responses and **death in one (0.6%)** survey response.
- Most harmful responses were due to *incorrectly determining or conveying the acuity of the scenario and recommended action*.
- The assessing physicians reported that the **LLM draft improved subjective efficiency in 76.9% of cases**.



COMMENTARY

Ambient Artificial Intelligence Scribes to Alleviate the Burden of Clinical Documentation

Aaron A. Tierney, PhD, Gregg Gayre, MD, Brian Hoberman, MD, MBA, Britt Mattern, MBA, Manuel Ballesca, MD, Patricia Kipnis, PhD, Vincent Liu, MD, MS, Kristine Lee, MD

Vol. 5 No. 3 | March 2024

DOI: 10.1056/CAT.23.0404

Pilot in 10000 physicians

Oct 2023 – Dec 2023

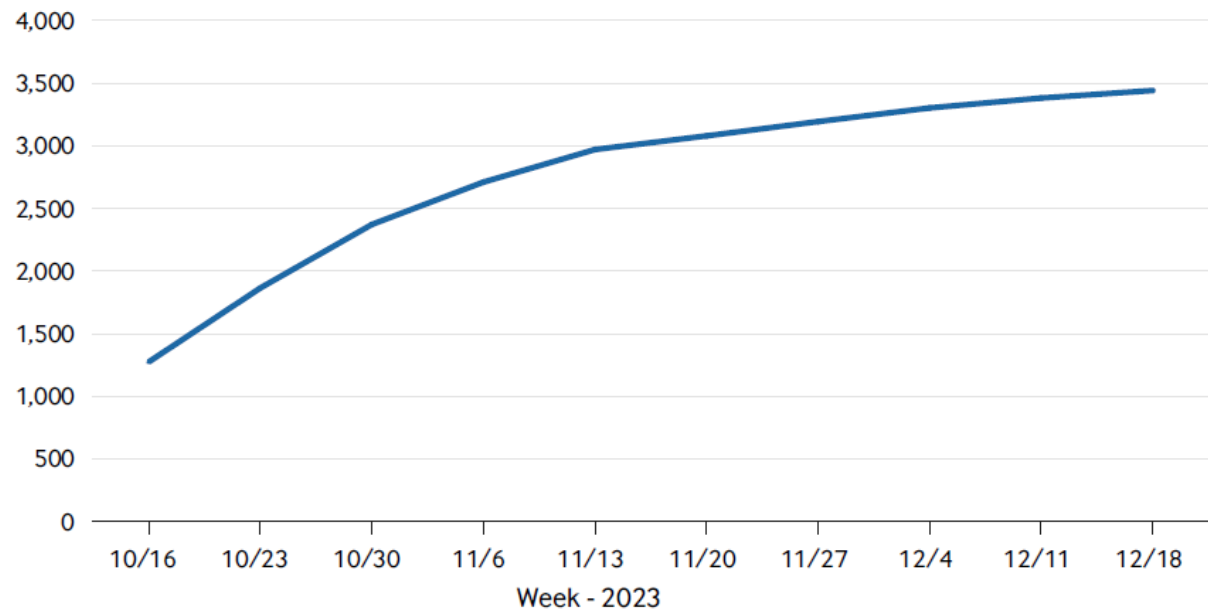
Modified PDQI-9, patient surveys, testimonials

Included assessment of confabulations and bias

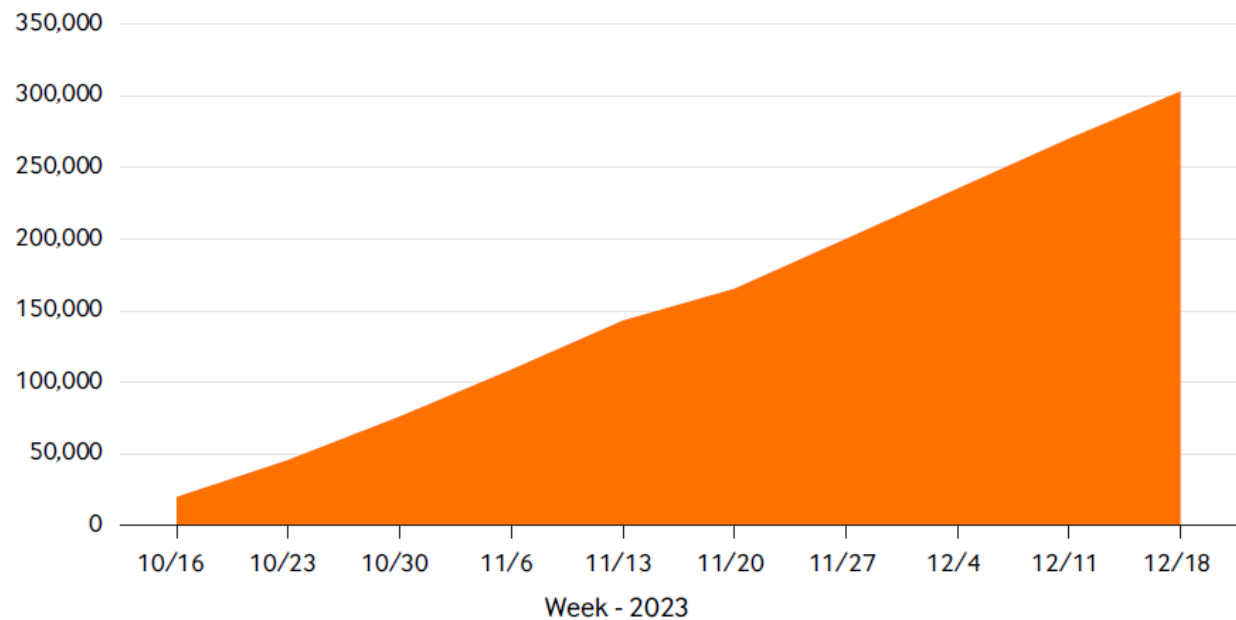
The study had four primary goals:

1. To assess uptake and engagement by both patients and clinicians.
2. To evaluate the effectiveness of the AI scribe in real clinical settings.
3. To determine if the AI scribe enhances the physician-patient relationship.
4. To verify that documentation quality was maintained.

Panel A. Unique Physicians Ever Using AI Scribe

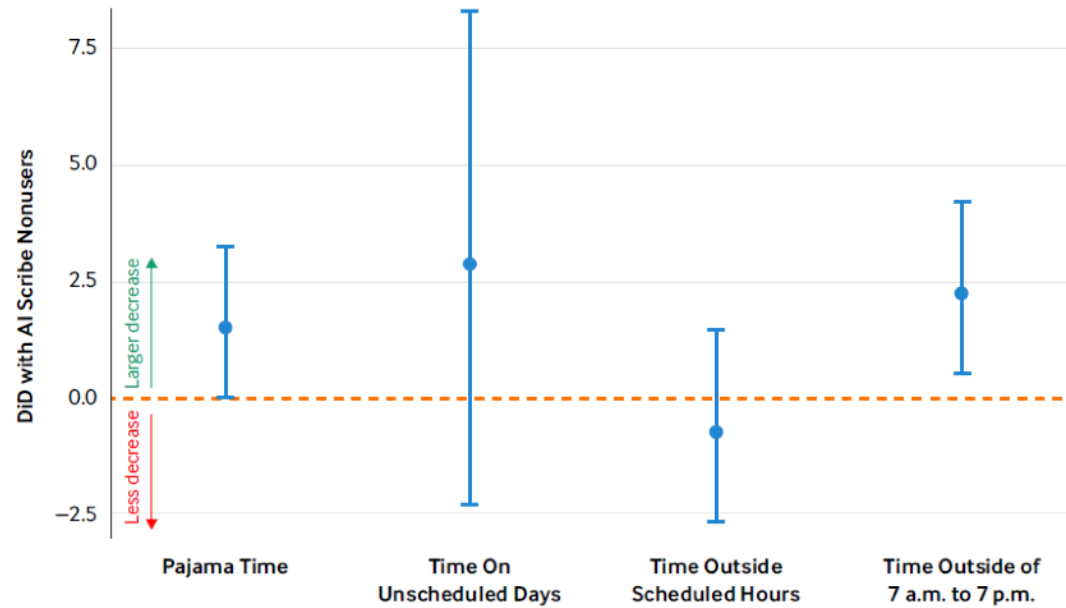


Panel B. Cumulative AI Scribe Visits

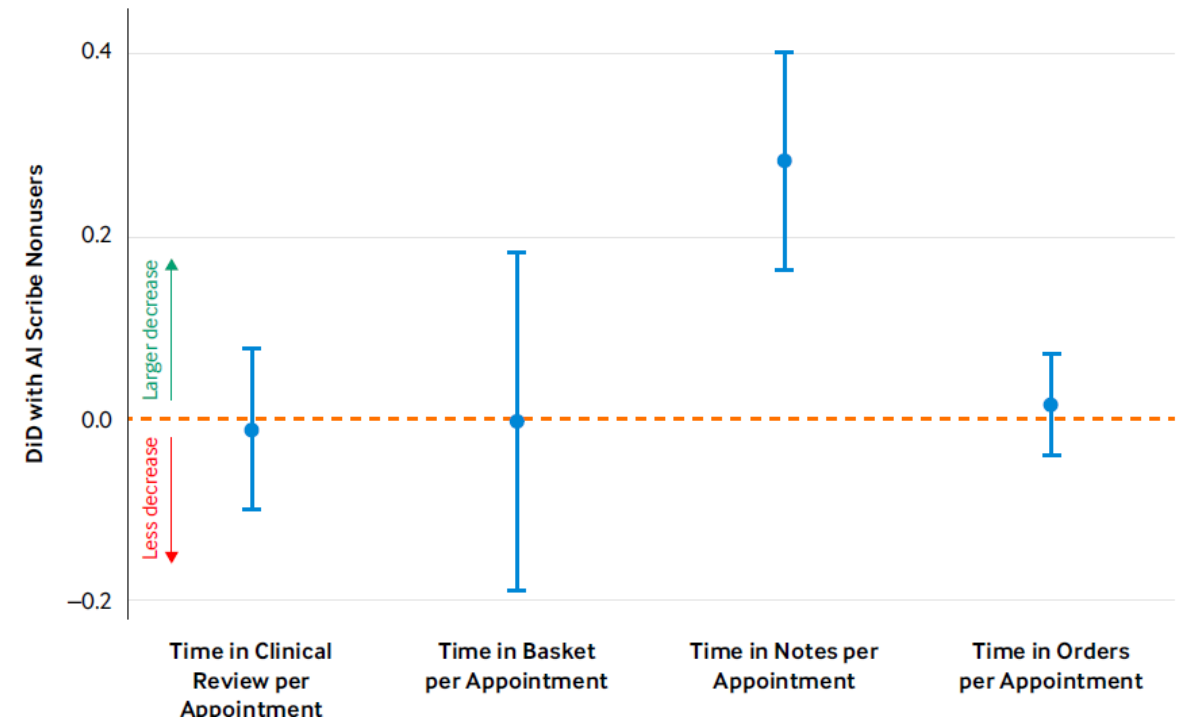




Panel A. Primary Care Physician Time Spent in the EHR-Related Activities



Panel B. Primary Care Physician Time Spent in Appointment-Related Activities



dose-response effect, with higher usage associated with more significant time reductions.



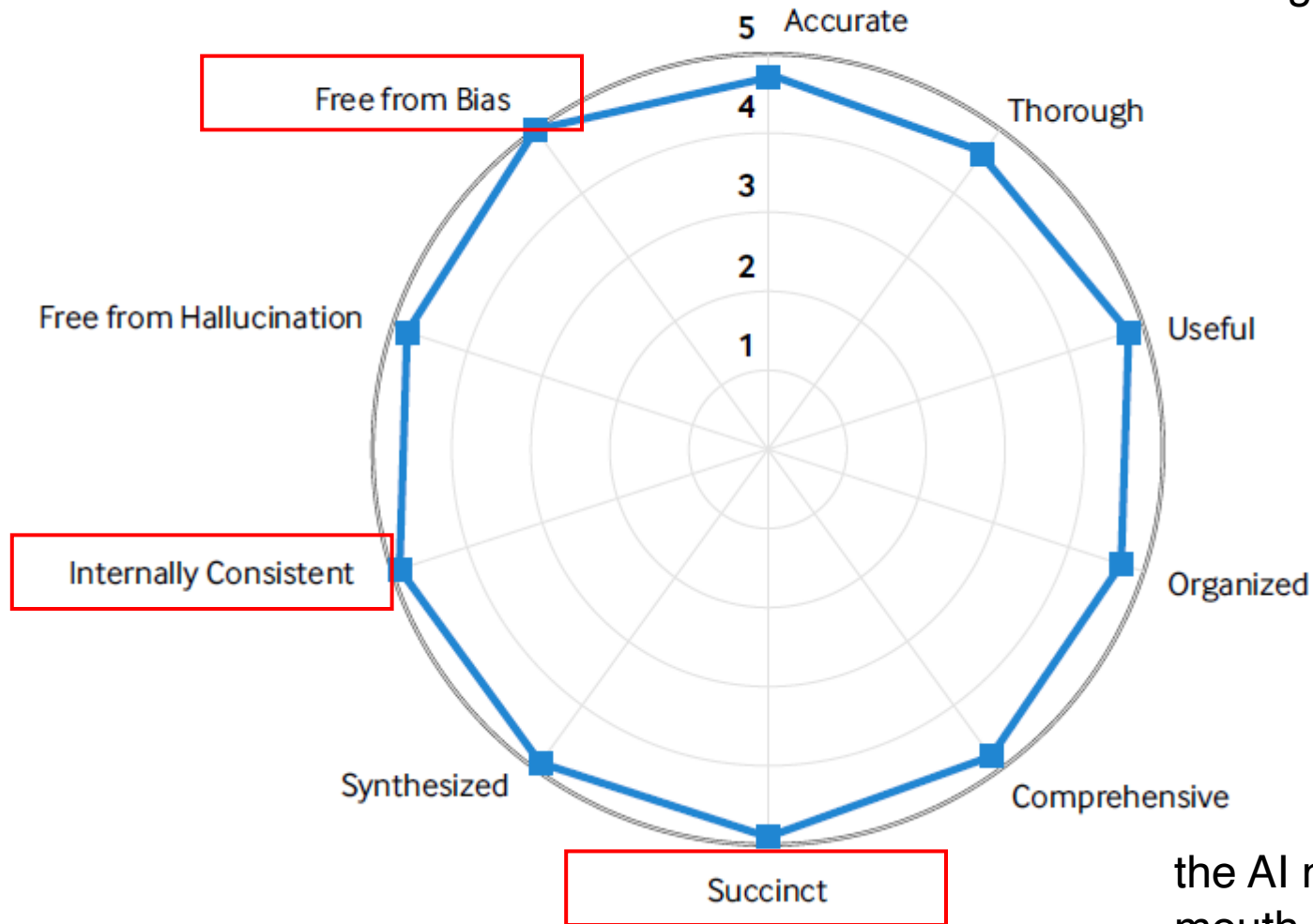
“more engaging and focused conversations with patients, enhancing the visit experience”

“that the scribe was especially helpful during lengthy appointments”

“game changer, it made notes more concise and improved the quality of visits”.

Feedback Category	Percentage
Patients who reported spending more time conversing with their physician	71%
Patients who reported spending less time conversing with their physician	1%
Patients who observed that their physician spent less time looking at the computer	81%
Patients who indicated that the AI scribe had no effect or improved the visit experience (negative effect)	100% (0%)
Patients who felt neutral to very comfortable with AI use in their care (felt uncomfortable)	100% (0%)

AI Summary Quality Metrics



the physician mentioned the need to schedule a prostate exam, yet the AI summarized this as the exam already having been performed

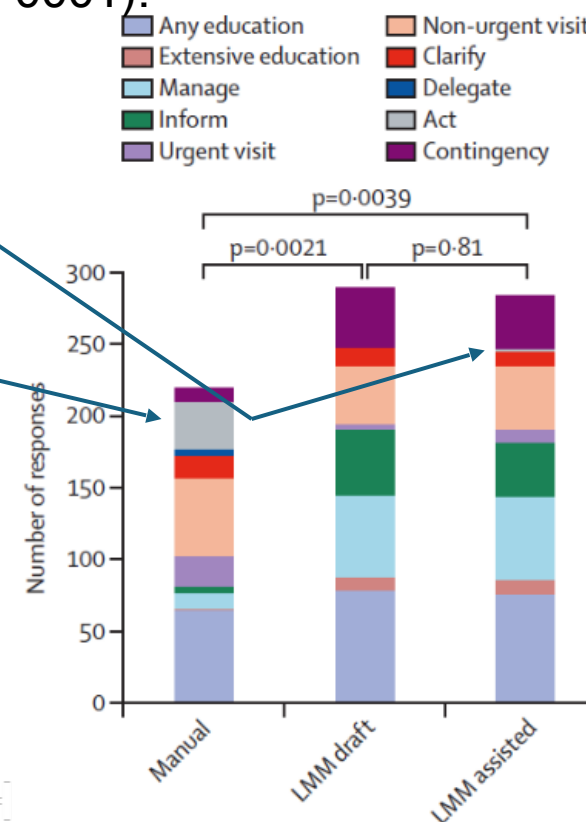
some summaries missed details, such as assessments for chest pain or anxiety.

the AI mistakenly inferred a diagnosis of hand, foot, and mouth disease when the physician had merely listed separate symptoms affecting the hand, feet, and mouth



Risk	Description	Example
Increased chart clutter	LLM note generation adds text volume, leading to need for summarization and more chart bloat.	Multiple team members create lengthy LLM notes. Covering physician requests LLM summary instead of reading all entries.
Decreased information density	LLMs generate verbose outputs that dilute essential clinical information.	A lengthy LLM-generated cardiology note lacks the focused insights of a concise staff cardiologist's note.
Persuasion and automation bias	LLMs may appear authoritative, causing clinicians to over-rely on their recommendations.	Primary team implements a tentative treatment plan directly from the LLM's confident tone without consulting with the original team.
Increased time to verify	Verifying and editing LLM-generated text adds to clinicians' workload.	Aware of confabulation risk, Physician spends extra time verifying a list of past medications generated by LLM to avoid redundant prescriptions.
Model collapse	LLMs trained on LLM-generated data risk "model collapse," losing insight and diversity in outputs.	An LLM trained on repetitive treatment data struggles to handle complex or rare cases due to limited exposure to varied clinical scenarios.

The mean manual response (34 words) was **shorter** than the LLM draft (169 words) and LLM-assisted responses (160 words; $p < 0.0001$).

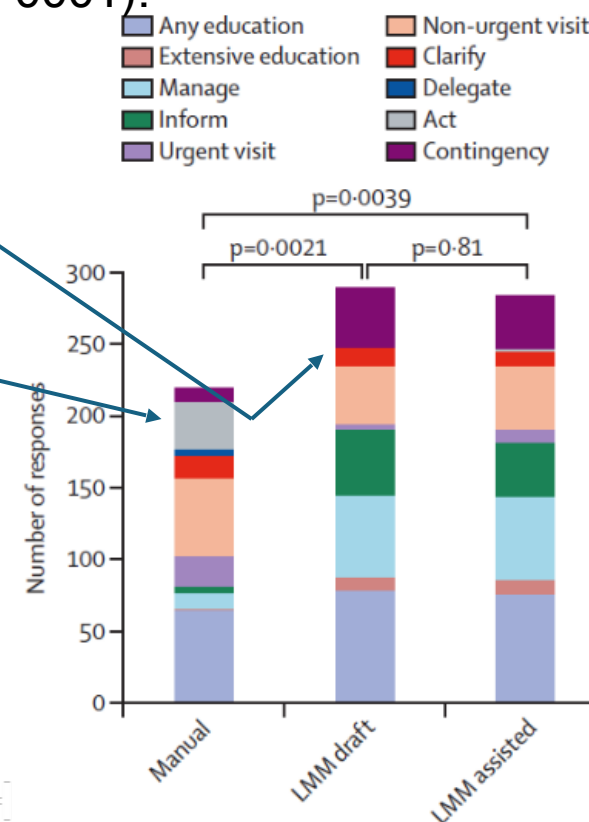


McCoy, L.G., A.K. Manrai, and A. Rodman, *Large Language Models and the Degradation of the Medical Record*. *New England Journal of Medicine*, 2024. **391(17): p. 1561-1564**.



Risk	Description	Example
Increased chart clutter	LLM note generation adds text volume, leading to need for summarization and more chart bloat.	Multiple team members create lengthy LLM notes. Covering physician requests LLM summary instead of reading all entries.
Decreased information density	LLMs generate verbose outputs that dilute essential clinical information.	A lengthy LLM-generated cardiology note lacks the focused insights of a concise staff cardiologist's note.
Persuasion and automation bias	LLMs may appear authoritative, causing clinicians to over-rely on their recommendations.	Primary team implements a tentative treatment plan directly from the LLM's confident tone without consulting with the original team.
Increased time to verify	Verifying and editing LLM-generated text adds to clinicians' workload.	Aware of confabulation risk, Physician spends extra time verifying a list of past medications generated by LLM to avoid redundant prescriptions.
Model collapse	LLMs trained on LLM-generated data risk "model collapse," losing insight and diversity in outputs.	An LLM trained on repetitive treatment data struggles to handle complex or rare cases due to limited exposure to varied clinical scenarios.

The mean manual response (34 words) was **shorter** than the LLM draft (169 words) and LLM-assisted responses (160 words; $p < 0.0001$).



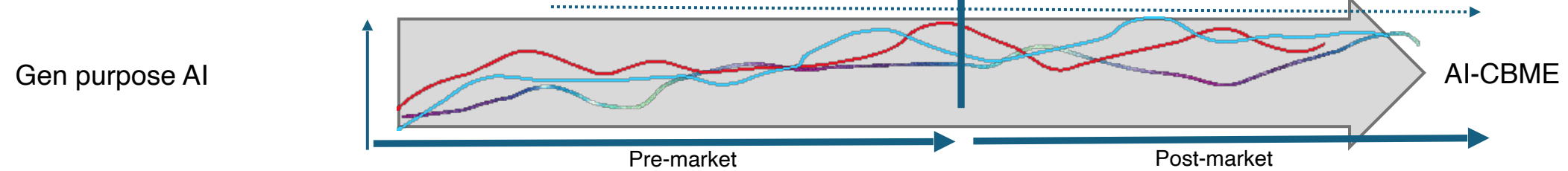
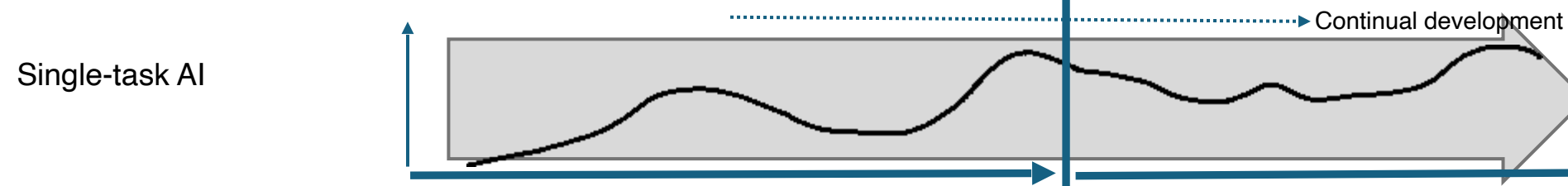
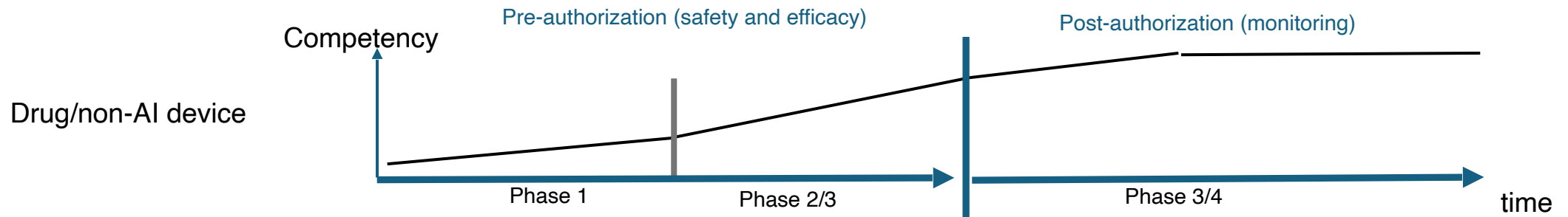
McCoy, L.G., A.K. Manrai, and A. Rodman, *Large Language Models and the Degradation of the Medical Record*. *New England Journal of Medicine*, 2024. **391(17): p. 1561-1564**.

Benevolence: The Need For Regulation

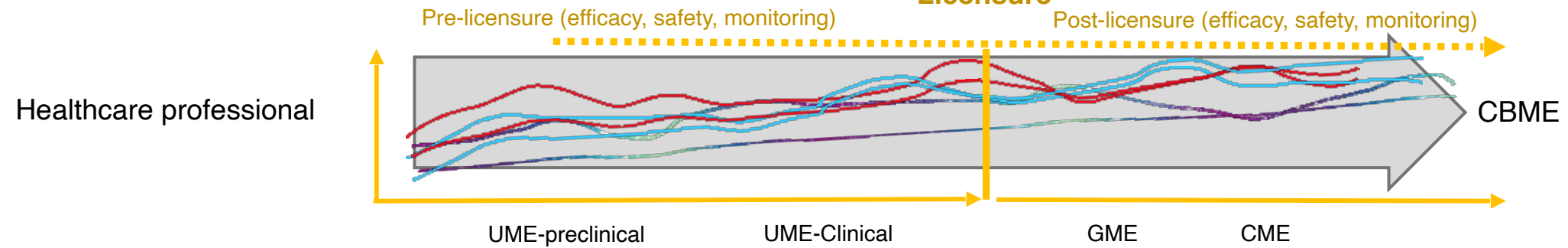


- Societal, non-commercial entities must be involved in the development life cycle and regulation
- FDA
- Non-profit coalitions (e.g., CHAI)
 - FDA stepped away
- The case for AI-CBME: Continuous assessment of multimodal AI by stakeholder “educators”
 - GenAIs are not deterministic but stochastic
 - Real-world GenAIs can drift
 - GenAIs are general purpose and can be used for a wide variety of often unanticipated tasks
 - Large-scale GenAIs, such as large language models (LLMs), are typically kept opaque
 - GenAIs are not fully describable – “dark complexity” make them unpredictable

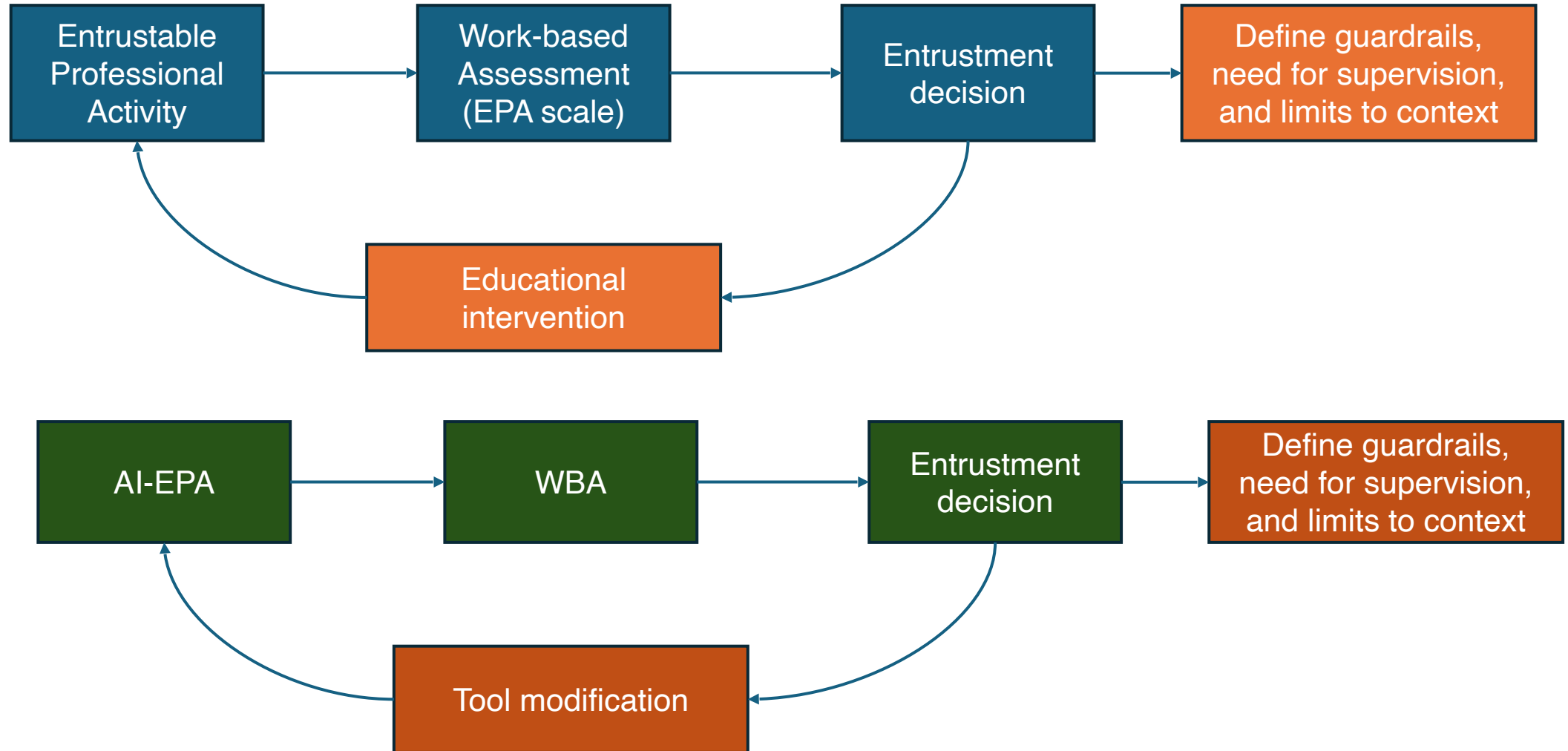
Authorization



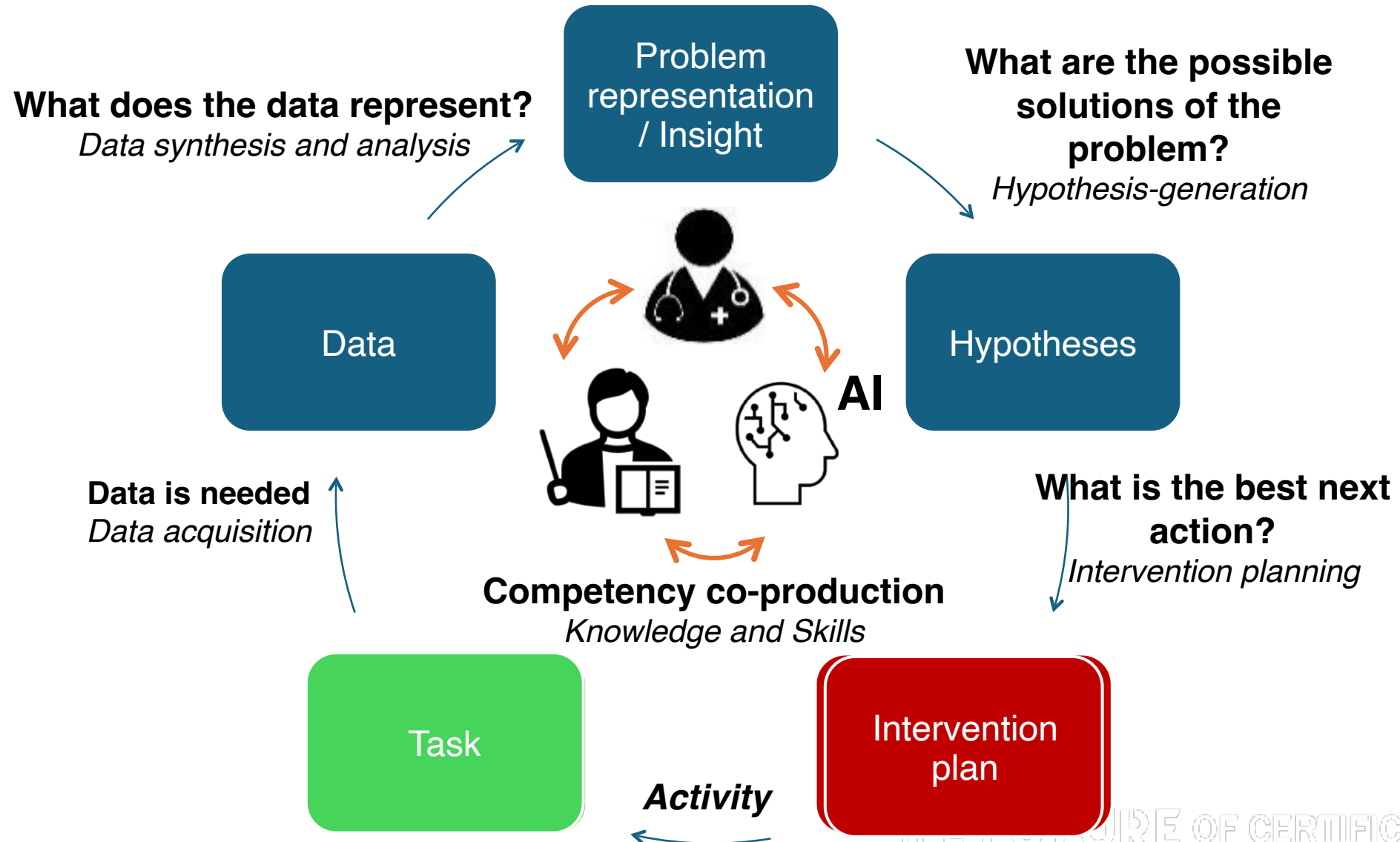
Licensure



Entrustment Framework To Safeguard Use Of AI In Health Professions Education



Conclusion: Co-Production





Thank you!

