

Data Science and AI for Medicine Training School

TRAINING: LLMs in Medicine

SPEAKER: Sanddhya Jayabalan

GEFÖRDERT VOM



Bundesministerium
für Forschung, Technologie
und Raumfahrt



SACHSEN Diese Maßnahme wird gefördert durch die Bundesregierung aufgrund eines Beschlusses des Deutschen Bundestages. Diese Maßnahme wird mitfinanziert durch Steuermittel auf der Grundlage des von den Abgeordneten des Sächsischen Landtags beschlossenen Haushaltes.

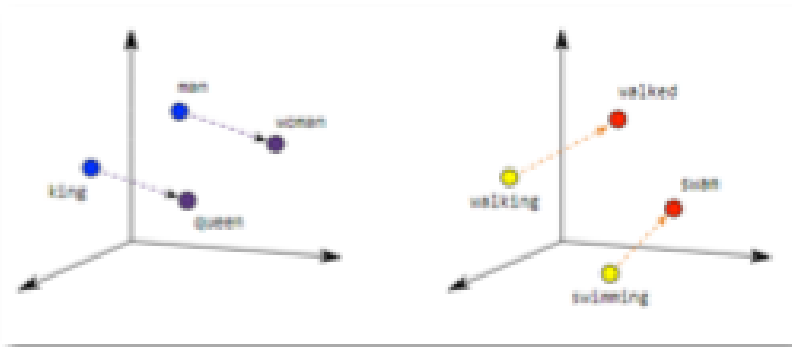


Come2Data
Kompetenzzentrum für
interdisziplinäre Datenwissenschaften

Data Science and AI for Medicine Training School
Training: Python Basics

Slide 1

ScaDS.AI
DRESDEN LEIPZIG



LLMs are the best
model of human
language we have!

Large Language Models.

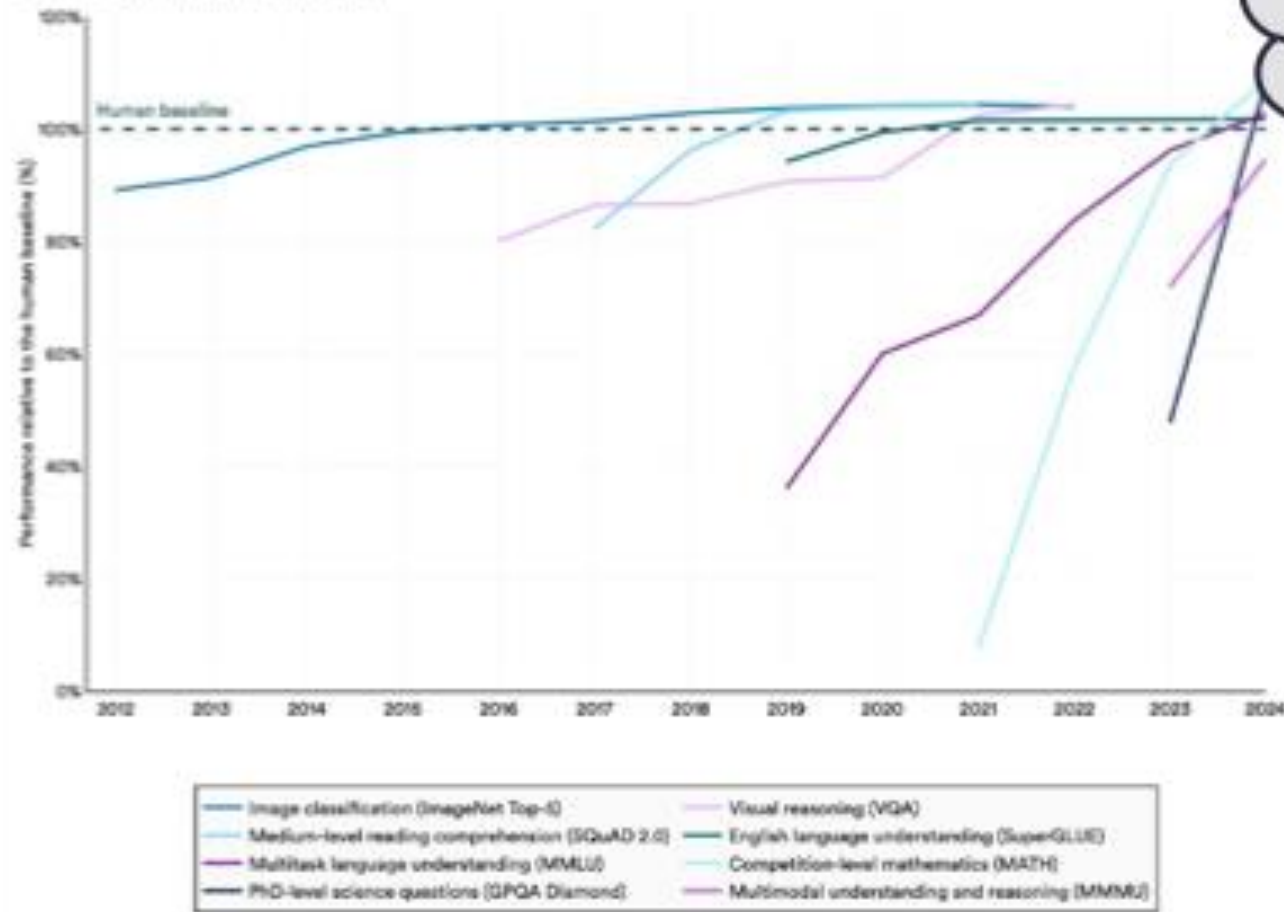


LLMs are highly performant – zero-shot!

LLMs can replace
medical doctors

Select AI Index technical performance benchmarks vs. human performance

Source: AI Index, 2020 | Chart: 2025 AI Index report



<https://hai.stanford.edu/ai-index/2025-ai-index-report>



Source: <https://www.tum.de/en/news-and-events/all-news/press-releases/details/chatgpt-gruender-sam-altman-an-der-tum>

Medicine.

80% of medical data is unstructured

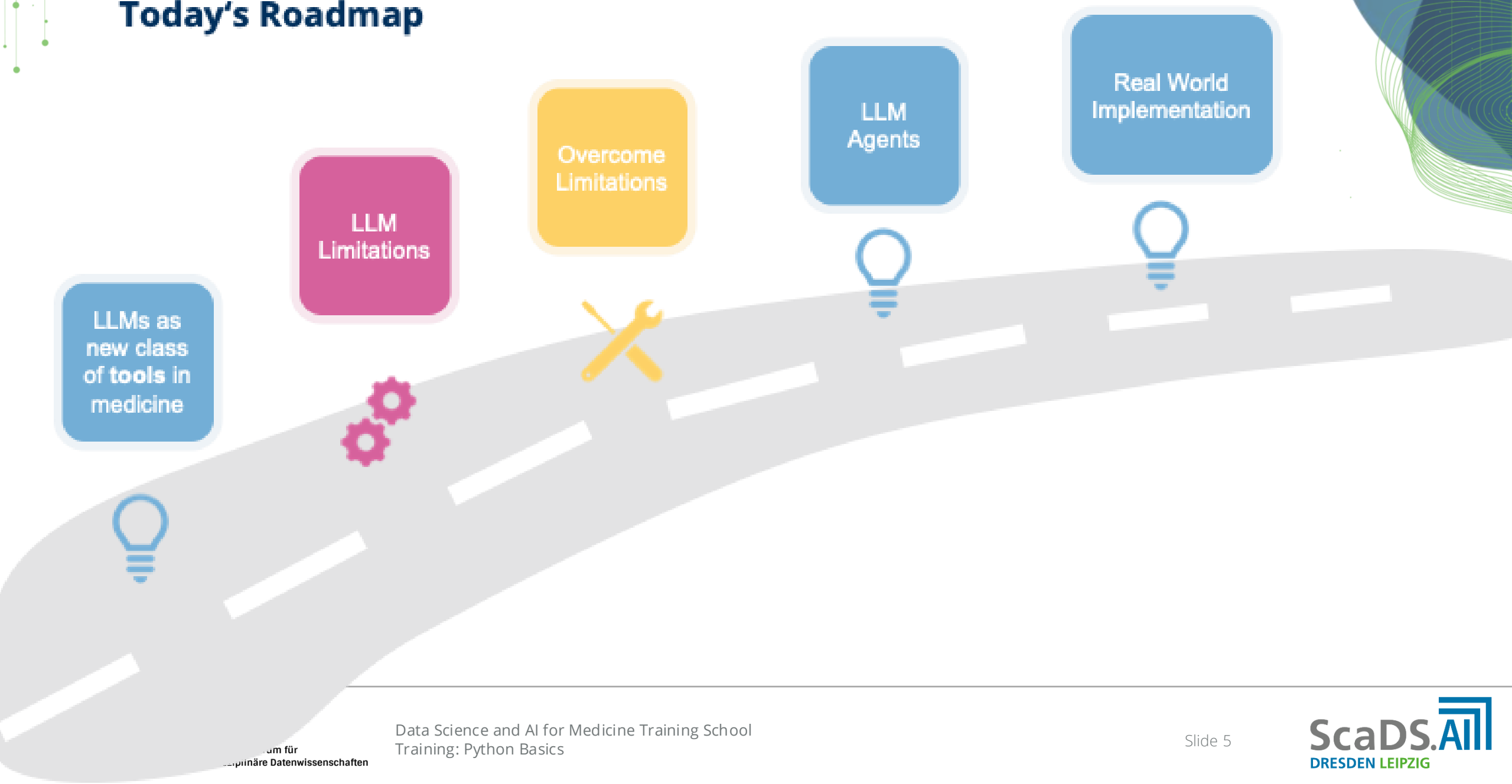
Growing knowledge

Narrative, free text contains valuable insights

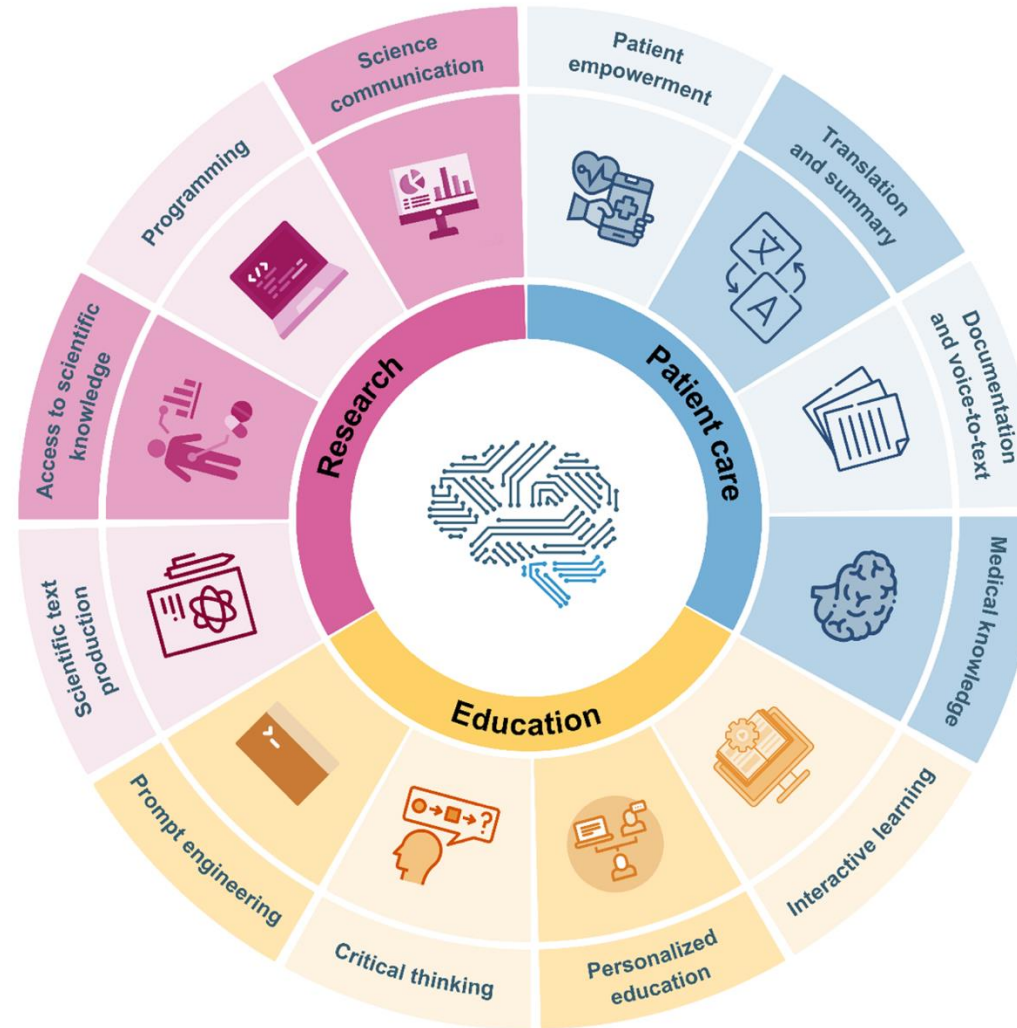
Need for personalized medicine



Today's Roadmap



Large Language Models: A new class of tools in medicine



Large Language Models: A new class of tools in medicine

npj | digital medicine

LLMs structure unstructured data



- Cancer registries
- Electronic Patient Records
- Medical Coding
- ...

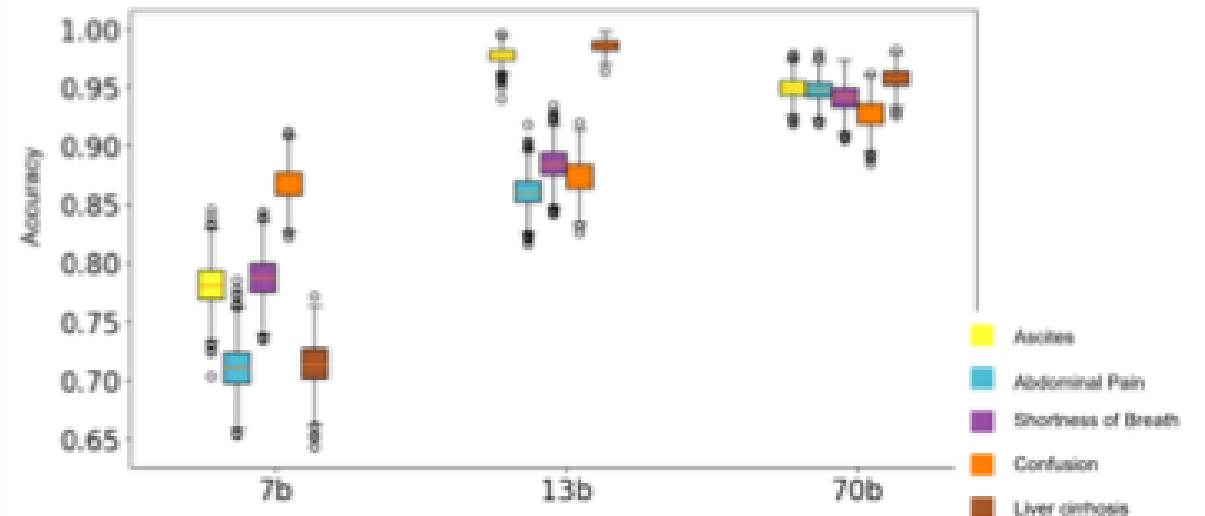
LLM-AIx



Privacy-preserving large language models for structured medical information retrieval

[Isabella Catharina Wiest](#), [Dyke Ferber](#), [Jiefu Zhu](#), [Marko van Treeck](#), [Sonja K. Meyer](#), [Radhika Juglan](#),

[Zunaimo I. Carrero](#), [Daniel Paech](#), [Jens Kleesiek](#), [Matthias P. Ebert](#), [Daniel Truhn](#) & [Jakob Nikolas Kather](#) 



Large Language Models: A new class of tools in medicine



LLMs summarize medical text



- Chart review
- Actionable guideline summaries
- ...

nature medicine

Article

<https://doi.org/10.1038/s41591-024-02855-0>

Adapted large language models can outperform medical experts in clinical text summarization

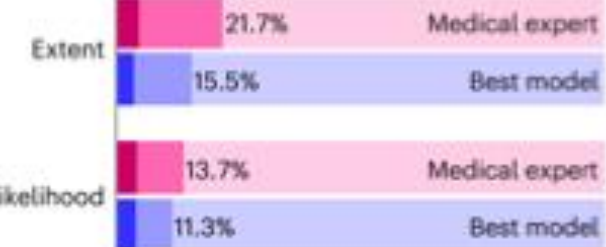
Received: 23 October 2023

Accepted: 2 February 2024

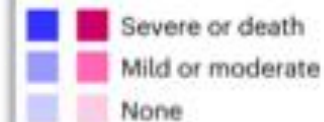
Published online: 21 February 2024

[Check for updates](#)

Daniël Van Veen^{1,2,3,4}, Caro Van Uden^{1,2}, Louis Blomkerker^{1,2}, Jean-Benoît Dufresne⁵, Asad Razi⁶, Christian Blumhagen^{1,2,3,4}, Anuj Pareek^{1,2,3,4}, Malgorzata Polacin⁷, Eduardo Portes Reis^{1,2}, Anna Seidenharn^{1,2,3,4}, Mihai Rohatgi^{1,2,3,4}, Poonam Hoxamand⁸, William Collins^{1,2}, Neema Khajep⁹, Curtis P. Langholz^{1,2,3,4}, Jason Hwang^{1,2}, Sergio Carido^{1,2}, John Pebody^{1,2} & Ashley S. Chaudhuri^{1,2,3,4}



Extent of harm



Likelihood of harm



Doctor LLM is ready to see you now.



<https://greater.com/healthungproblems/>

NEJM
AI

NEJM AI 2025;2(4)
DOI: 10.1056/AIa.2400832

ORIGINAL ARTICLE

Randomized Trial of a Generative AI Chatbot for Mental Health Treatment

Michael V. Heinz M.D.,^{1,2} Daniel M. Madsen Ph.D.,^{1,2} Brianna M. Trudeau B.A.,¹ Sukanya Bhattacharya B.A.,² Yufei Wang M.S.,² Haley A. Santa Abi O. Jemini B.A.,² Abigail J. Saltsman B.A.,² Tess Z. Griffin Ph.D.,² and Nicholas C. Jacobson Ph.D.^{1,2,3,4}

Received: August 11, 2024; Revised: November 18, 2024; Accepted: February 2, 2025; Published March 27, 2025

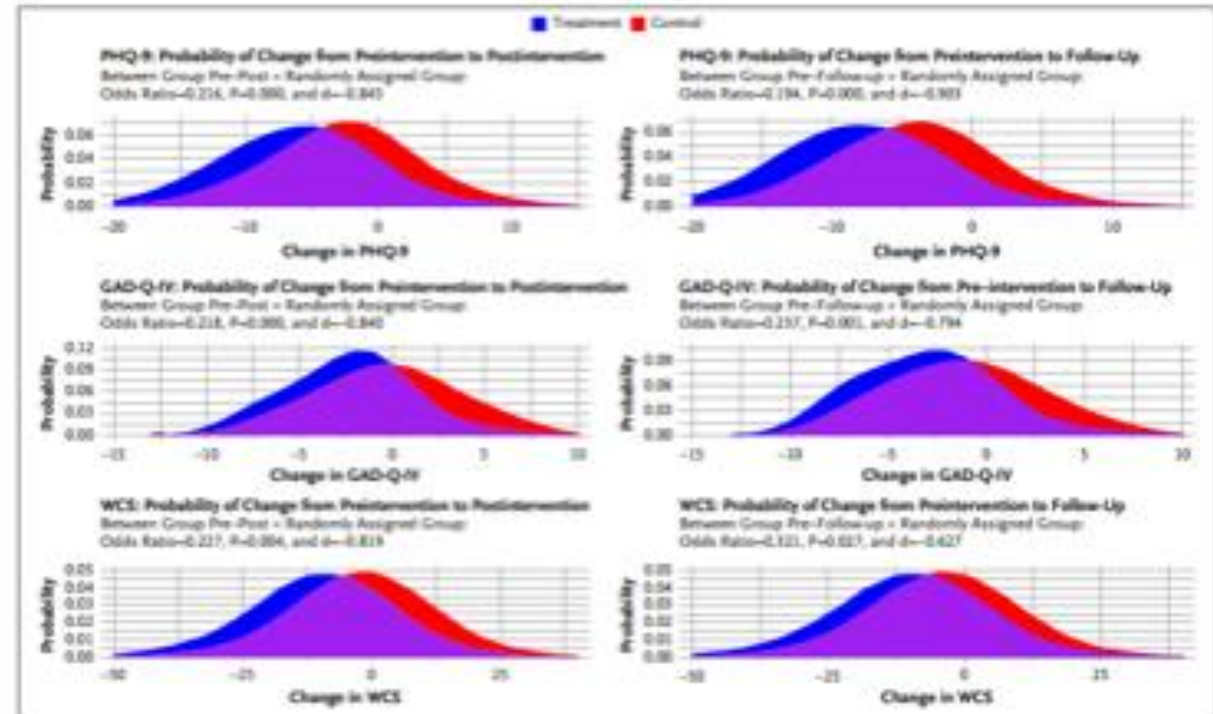


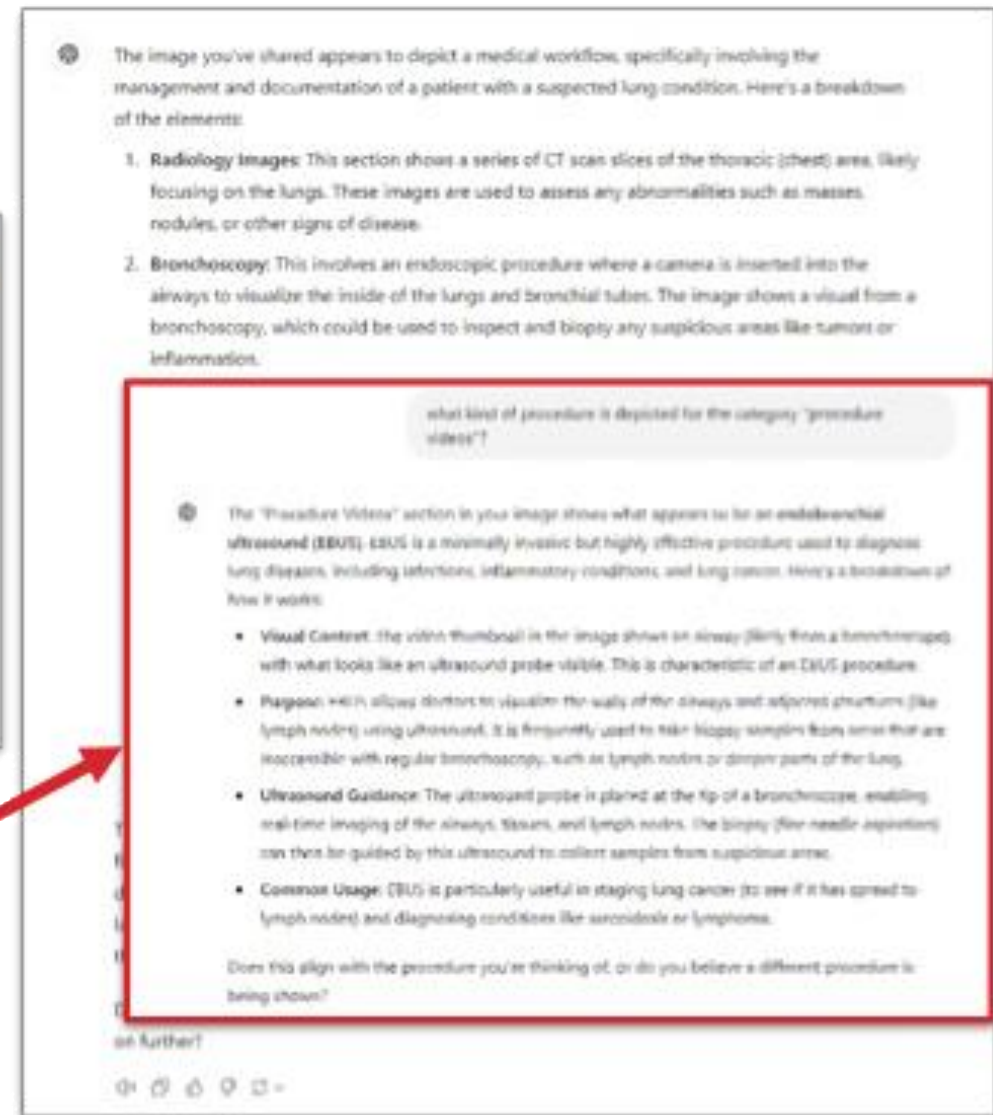
Figure 3. Distributions Representing Smoothed Probability of Changes in Clinical Outcomes (Depression, Anxiety, Weight Concerns, Row-Wise) Postintervention (4 Weeks, Left Column) and at Follow-Up (8 Weeks, Right Column).

The probabilities shown in these plots are derived directly from the CLMMs through predicted probabilities for each possible change score under treatment and control conditions. The treatment group is visualized in blue, and the control group is visualized in red. CLMM denotes cumulative-link mixed model; GAD-Q-IV, Generalized Anxiety Disorder Questionnaire for the Diagnostic and Statistical Manual of Mental Disorders, Fourth Edition; PHQ-9, Patient Health Questionnaire 9; and WCS, Weight Concerns Scale.

LLMs have impressive zero-shot abilities



Zero-shot!





LLMs come with limitations



LLMs hallucinate or present factual inaccuracies



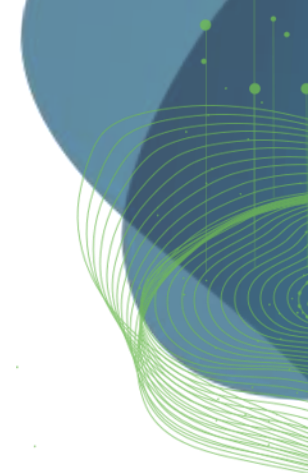
LLMs are better reasoning engines, do not use them as knowledge bases



If your model is not hosted locally, you can never be sure about data privacy

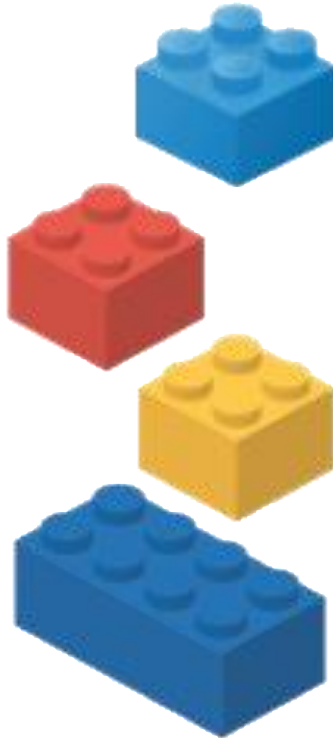


Always be sensitive about bias and fairness issues





How to build a great prompt



Use examples

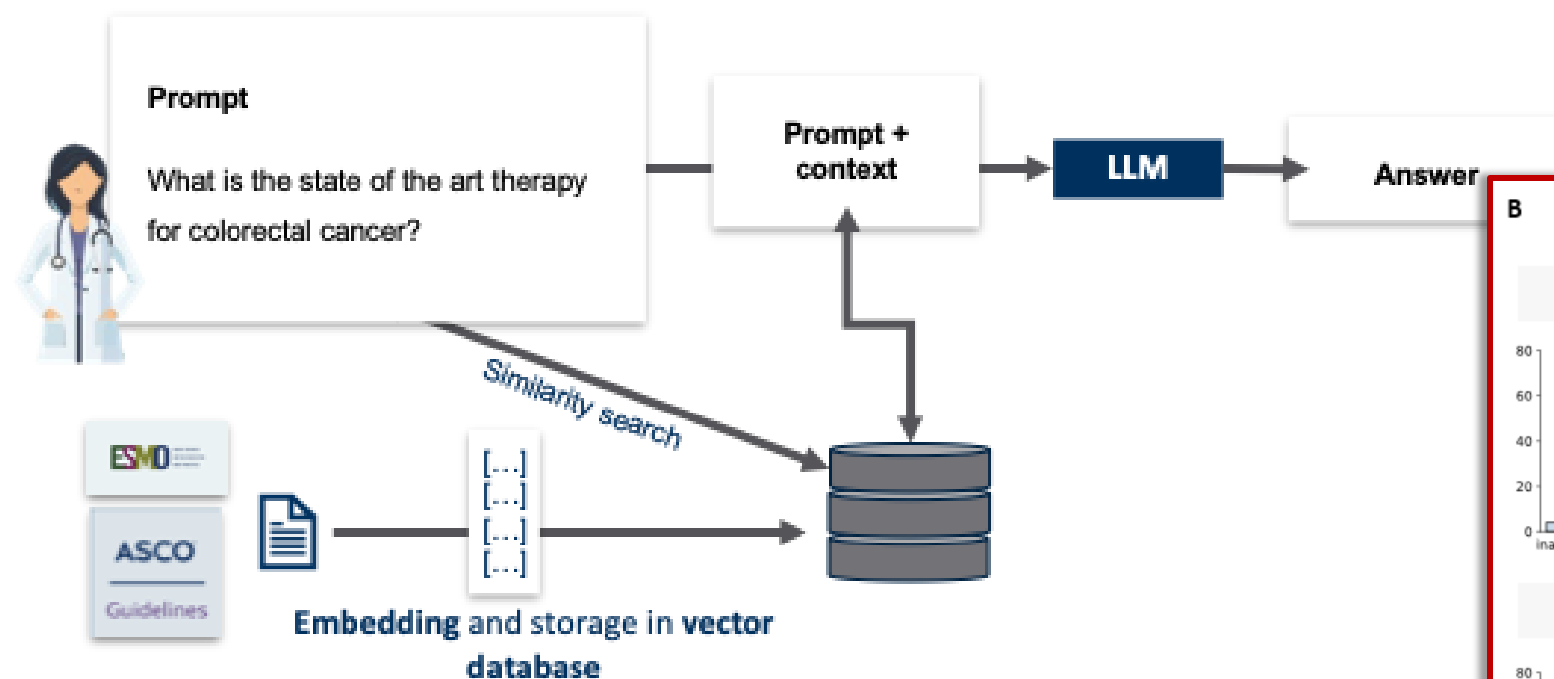
Define style, audience, role

Give structure

Be clear and precise



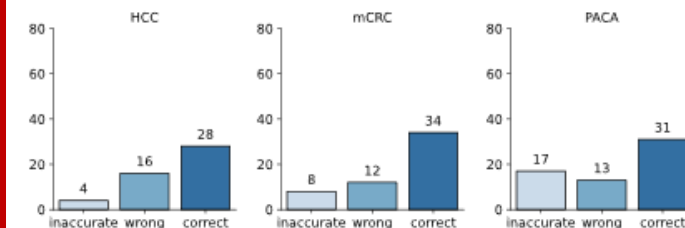
LLMs can be augmented with information from solid knowledge databases



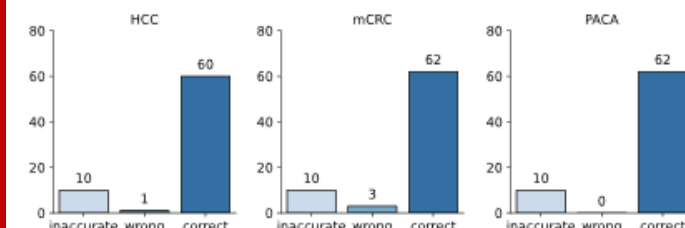
Ferber, D., Wiest, I. C., Wollstein, G., Ebert, M. P., Beutel, G., Eckardt, J. N., ... & Kothner, J. N. (2024). GPT-4 for Information Retrieval and Comparison of Medical Oncology Guidelines. *NPJ AI*, 8(200235).

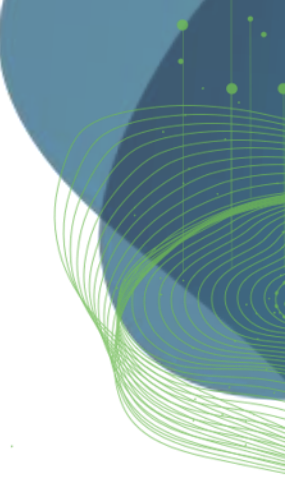
B Human Response Evaluation

Without RAG



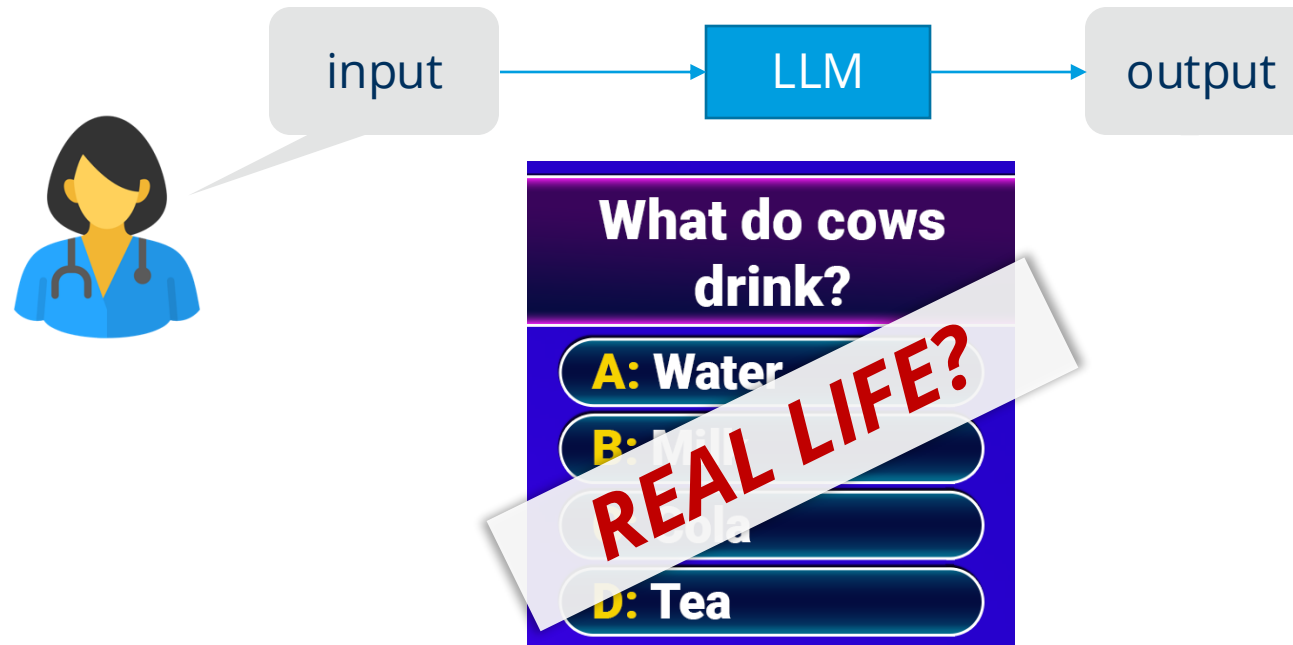
With RAG



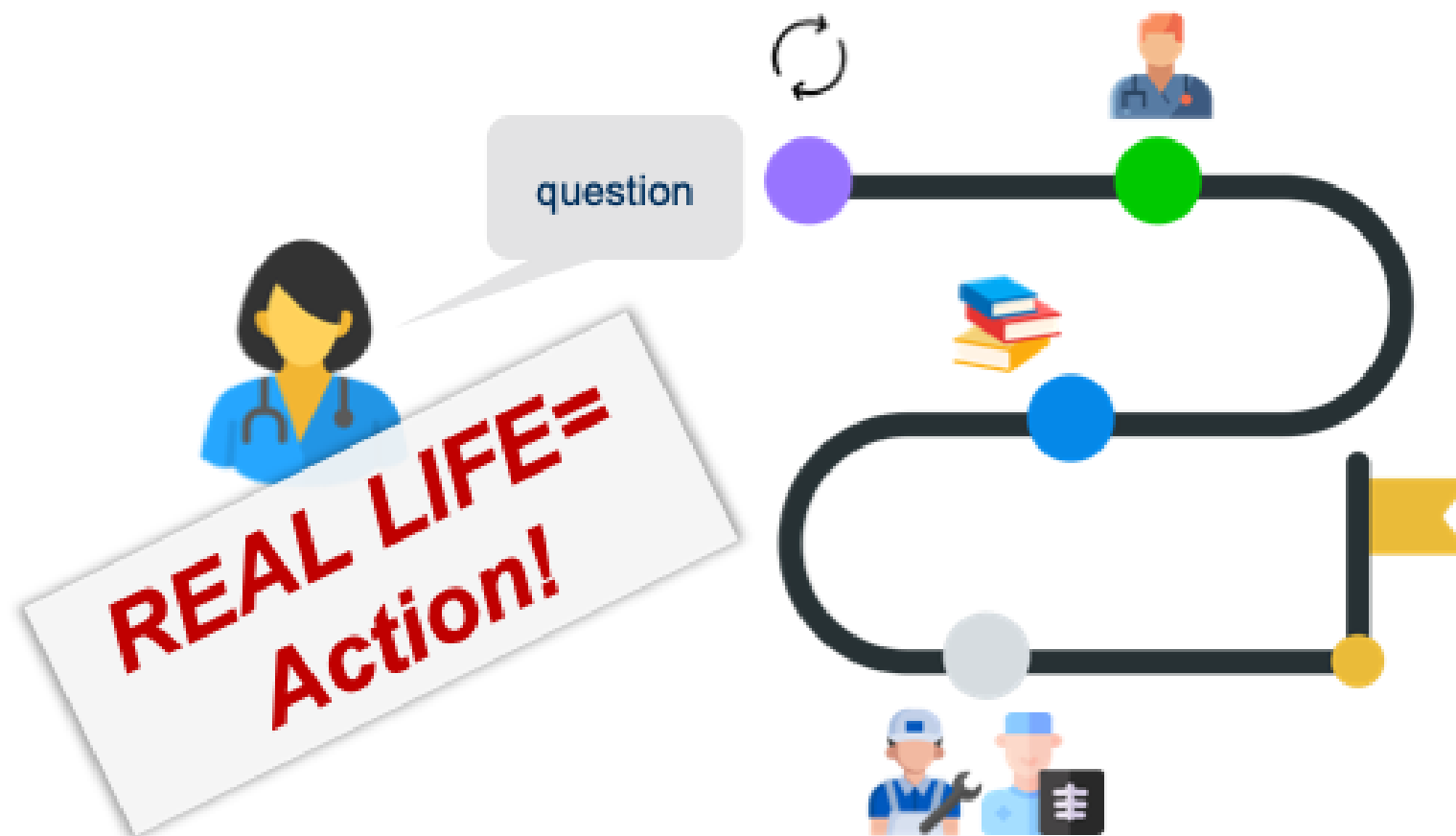


- **Prompt engineering:** clear instructions, structure, examples
- **Use as reasoning engines,** not static knowledge bases
- **Reliable knowledge base:** retrieval, curated corpora, RAG
- **Stepwise reasoning:** “think step by step,” chain-of-thought
- **Tool use & multimodality:** calculators, EHRs, APIs, websearch
- **Human-in-the-loop:** expert validation & oversight
- **Domain adaptation:** fine-tuning or instruction tuning
- **Cross-checking:** ensembles, multiple prompts, consistency
- **Safety guardrails:** rules, filters, clinical guidelines
- **Awareness & training:** educate users about risks/limits

The era of agentic AI - from models to AI Agents

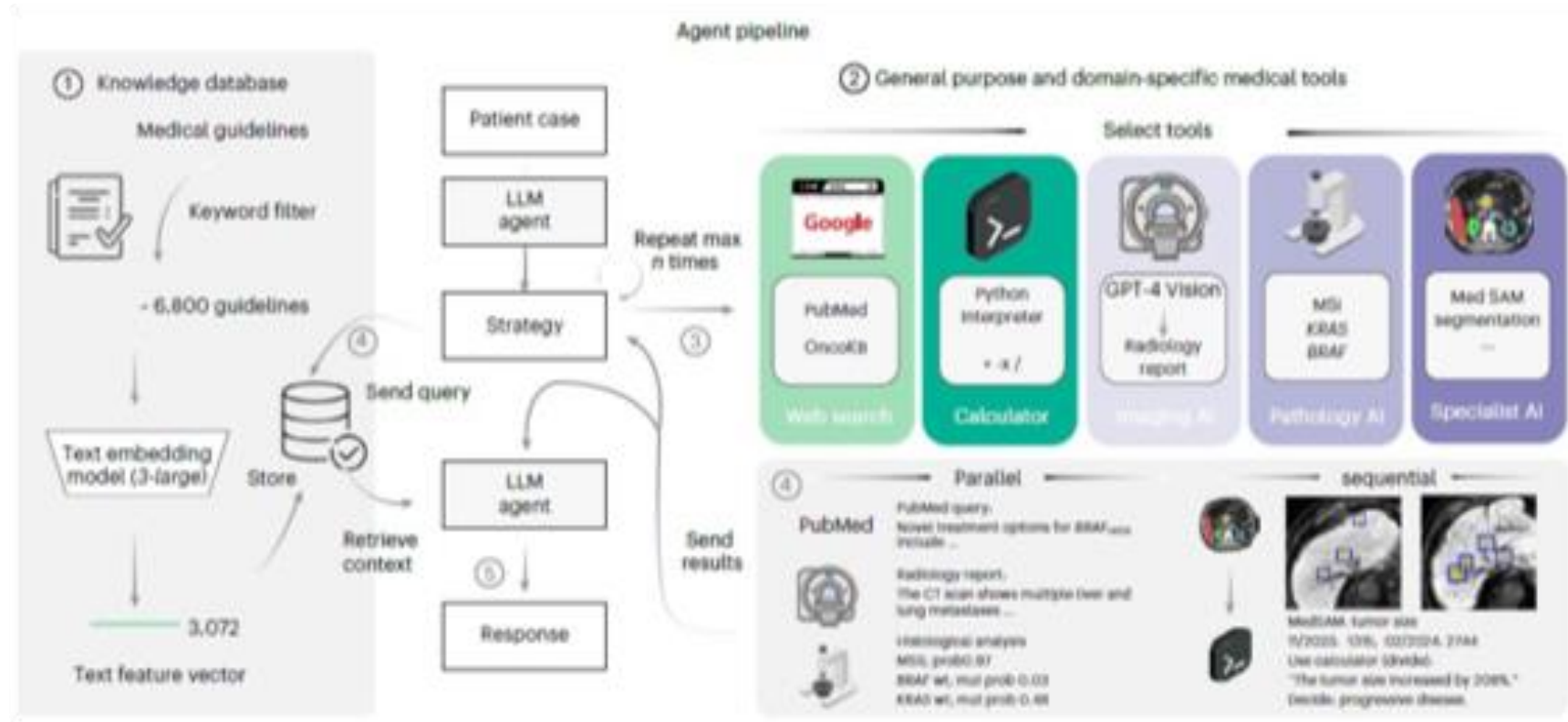


The era of agentic AI - from models to AI Agents



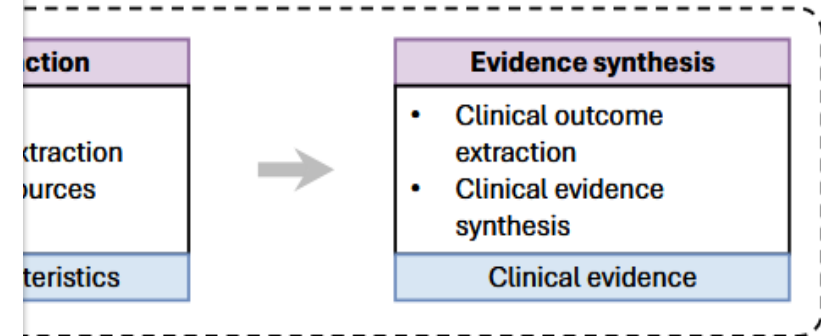
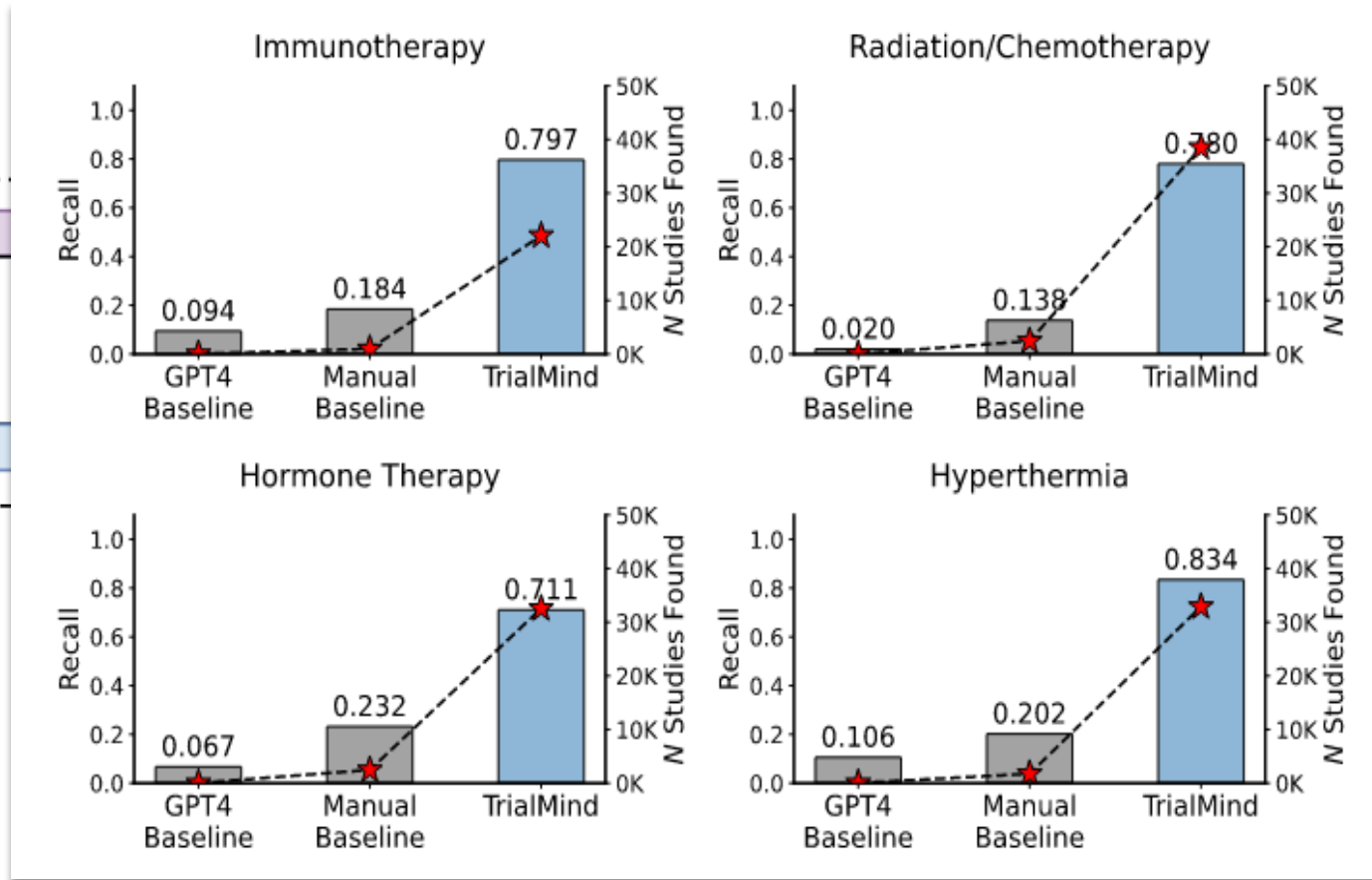
- ✓ Reason
- ✓ Consult experts
- ✓ Knowledge
- ✓ Tools

LLM agents for Clinical Decision Support



Ferber, D., ElNahhas, O.S.M., Wölflert, G. et al. Development and validation of an autonomous artificial intelligence agent for clinical decision-making in oncology. *Nat Cancer* 6, 1337–1349 (2025). <https://doi.org/10.1038/s43018-025-00991-6>

The era of agentic AI - from models to AI Agents



System-designed LLMs outperform humans and simple LLMs

Wang, Zifeng, et al. "Accelerating clinical evidence synthesis with large language models." *npj Digital Medicine* 8.1 (2025): 509.

The era of agentic AI – questions for implementation

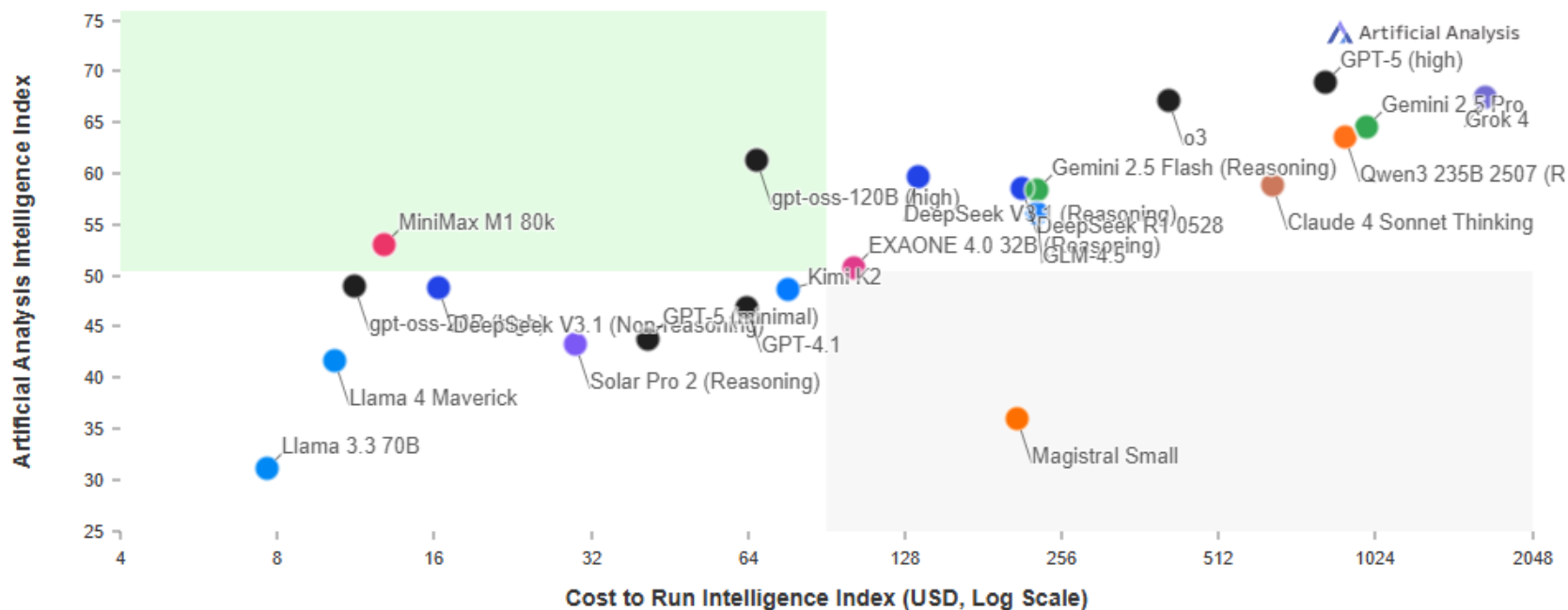


Intelligence vs. Cost to Run Artificial Intelligence Index

Artificial Analysis Intelligence Index; Cost to Run Intelligence Index

Most attractive quadrant

Alibaba Anthropic DeepSeek Google LG AI Research Meta MiniMax Mistral Moonshot AI OpenAI
Upstage xAI Z AI



<https://artificialanalysis.ai/>

Model Performance and real world effectiveness might differ

JAMA Network Open

Original Investigation | Health Informatics

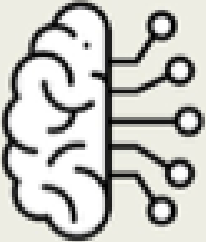
Large Language Model Influence on Diagnostic Reasoning

A Randomized Clinical Trial

Ethan Goh, MBBS, MS^{1,2}, Robert Gallo, MD³, Jason Hom, MD⁴, et al

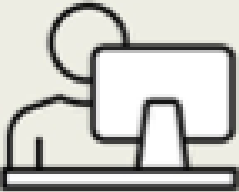
INTERVENTION

50 Participants randomized



25 Generative artificial intelligence (AI) chatbot

Participants with access to AI chatbot were allocated 60 min to review up to 6 clinical vignettes

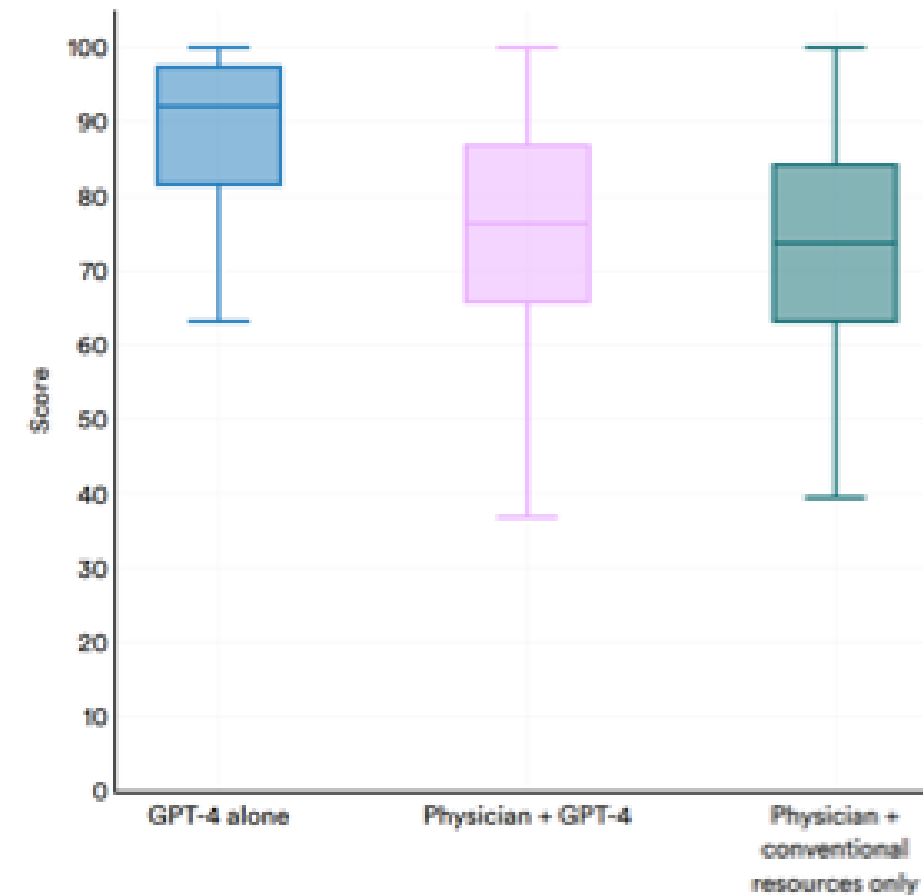


25 Conventional resources

Participants with access to conventional online resources such were allocated 60 min to review up to 6 clinical vignettes

LLM performance in clinical diagnosis

Source: Goh et al., 2024 | Chart: 2025 AI Index report



Find Use Cases that Matter

Impact of AI Scribe on physician EHR usage

Source: Ma et al., 2024 | Chart: 2025 AI Index report

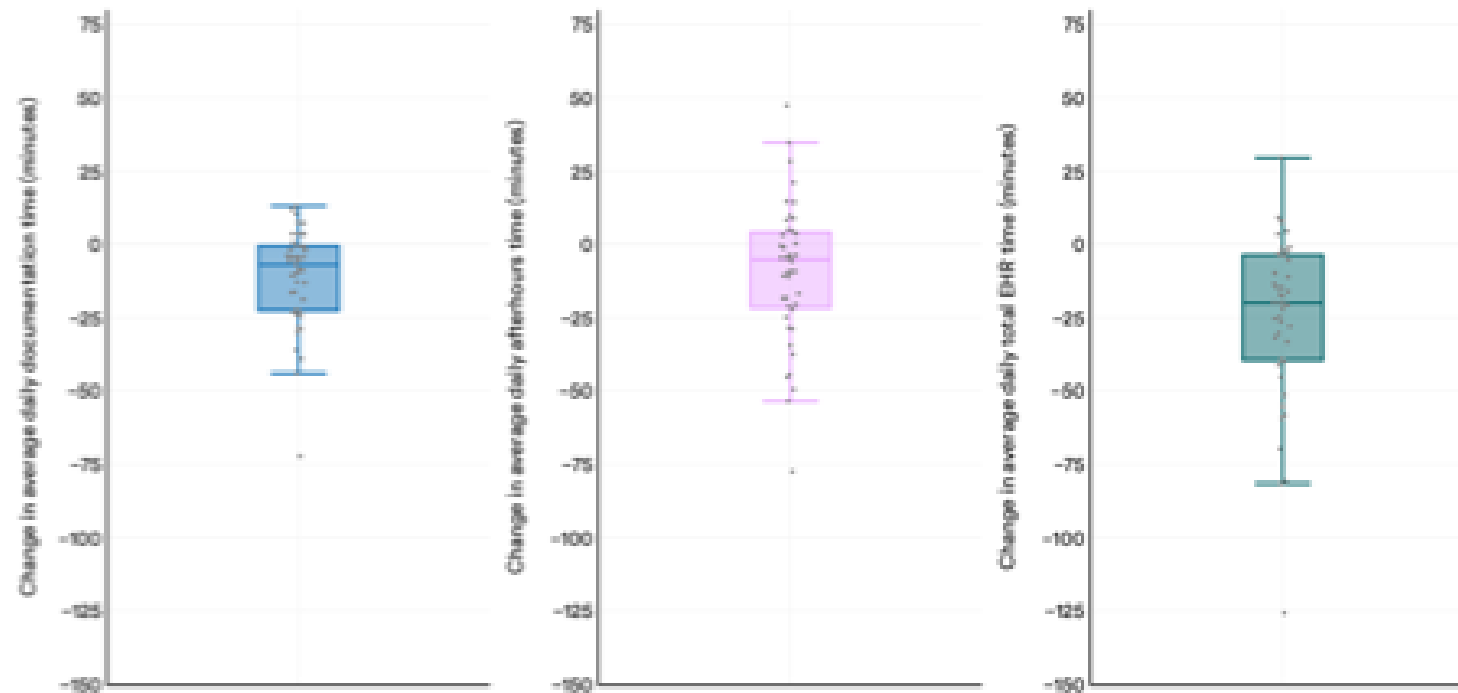


Figure S4.9

AI scribe led to

↓ 30 seconds per note

↓ 20 minutes EHR time per day

↓ ~30% less burden and burnout



Any questions or remarks?

Let's practice – Python Basics in Jupyter Lab