# AI in Medicine: The Resident Expert That Never Sleeps

Agustina Saenz, MD, MPH

---

## Disclosures

- Senior Clinical Informaticist at Curai Health

# Agenda

📄 Definitions     🔲 Use cases     ⚠️ Challenges

⚖️ Ethics     🔨 Liability     ▭ Model Cards

Mass General Brigham          HARVARD MEDICAL SCHOOL

---

# What is Artificial Intelligence (AI)?

- "Artificial Intelligence refers to the development of computer algorithms that can perform tasks that typically require human intelligence, such as learning, reasoning, perception, and decision making" (ChatGPT)

- AI is based on machine learning algorithms and other computational techniques:

**Probabilistic Reasoning**
Machine Learning
Predictive Modeling
Deep Learning
Decision Trees

**Computational Logic**
Rule-Based Systems
Logic Programming
Heuristic Techniques
Case Based Reasoning

**Optimization Techniques**
Constraint Satisfaction
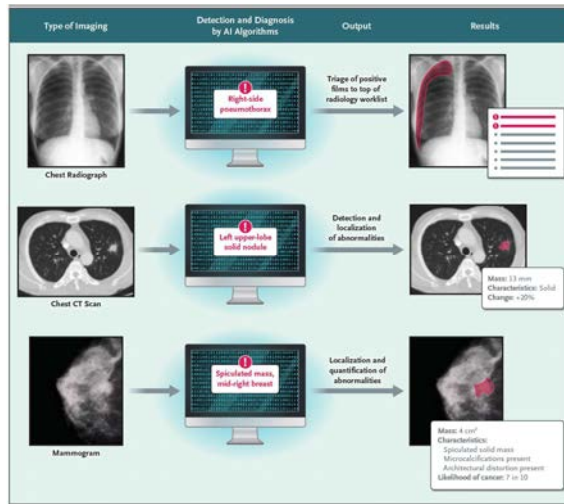Constraint-Based Reasoning
Linear Programming
Genetic Algorithms

**Perception Systems/Ambient Intelligence**
Computer Vision
Olfactive & Haptic
Auditory/Speech

**Knowledge Representation, Learning/Search**
Knowledge Graphs
Semantic Networks

**Natural Language Processing**
Text Analytics
Natural Language Understanding
Chatbots
Natural Language Generation
Dialog Management

Credit: Adam Landman

Mass General Brigham          HARVARD MEDICAL SCHOOL
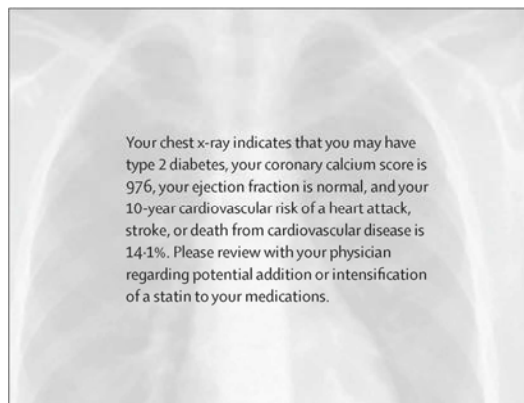
# Medical Image Interpretation



Rajpurkar, Pranav, and Matthew P. Lungren. "The Current and Future State of AI Interpretation of Medical Images." *New England Journal of Medicine* 388.21 (2023): 1981-1990.

---

## Digital medicine
### AI-enabled opportunistic medical scan interpretation

Eric Topol



Your chest x-ray indicates that you may have type 2 diabetes, your coronary calcium score is 976, your ejection fraction is normal, and your 10-year cardiovascular risk of a heart attack, stroke, or death from cardiovascular disease is 14·1%. Please review with your physician regarding potential addition or intensification of a statin to your medications.

# Large Language Models

How frequently do you utilize ChatGPT for work-related activities?*

A)  Never used it
B)  I used it a few times, but it was not helpful
C)  I used it a few times, but I recently started to use it more
D)  I use it almost every day, can't work without it!

---

# Opportunities for LLMs



Meskó, Bertalan, and Eric J. Topol. "The imperative for regulatory oversight of large language models (or generative AI) in healthcare." *npj Digital Medicine* 6.1 (2023): 120.
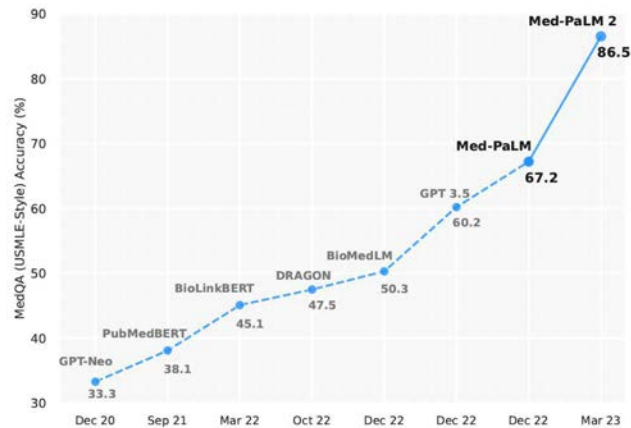
Mass General Brigham

HARVARD
MEDICAL SCHOOL

---

# Medical question answering

| USMLE Sample Exam | GPT-4 (5 shot) | GPT-4 (zero shot) | GPT-3.5 (5 shot) | GPT-3.5 (zero shot) | ChatGPT[†] (zero shot) |
|---|---|---|---|---|---|
| Step 1 | **85.71** | 80.67 | 52.10 | 51.26 | 55.1 |
| Step 2 | **83.33** | 81.67 | 58.33 | 60.83 | 59.1 |
| Step 3 | **90.71** | 89.78 | 64.96 | 58.39 | 60.9 |
| Overall Average* | **86.70** | 84.31 | 58.78 | 56.91 | – |

* Calculated as $\frac{\#correct}{\#questions}$ across all three steps. Each step has slightly different sample size.

Nori, Harsha, et al. "Capabilities of gpt-4 on medical challenge problems." *arXiv preprint arXiv:2303.13375* (2023).

Mass General Brigham

HARVARD
MEDICAL SCHOOL

# Medical question answering



Singhal, Karan, et al. "Towards expert-level medical question answering with large language models." *arXiv preprint arXiv:2305.09617* (2023).

Mass General Brigham · HARVARD MEDICAL SCHOOL

---

# Curbside consult



Lee, Peter, Sebastien Bubeck, and Joseph Petro. "Benefits, limits, and risks of GPT-4 as an AI chatbot for medicine." *New England Journal of Medicine* 388.13 (2023): 1233-1239.

Mass General Brigham · HARVARD MEDICAL SCHOOL

# Documentation

Lee, Peter, Sebastien Bubeck, and Joseph Petro. "Benefits, limits, and risks of GPT-4 as an AI chatbot for medicine." *New England Journal of Medicine* 388.13 (2023): 1233-1239.

Mass General Brigham — HARVARD MEDICAL SCHOOL

---

## Large language models encode clinical knowledge

nature

**a**

| | Evidence of correct comprehension |
|---|---|
| Flan-PaLM | 90.5% |
| Med-PaLM | 97.5% |
| Clinician | 97.8% |

| | Evidence of correct retrieval |
|---|---|
| Flan-PaLM | 76.3% |
| Med-PaLM | 95.4% |
| Clinician | 97.8% |

| | Evidence of correct reasoning |
|---|---|
| Flan-PaLM | 85.7% |
| Med-PaLM | 92.5% |
| Clinician | 97.7% |

**b**

| | Evidence of incorrect comprehension |
|---|---|
| Flan-PaLM | 9.2% |
| Med-PaLM | 5.0% |
| Clinician | 2.2% |

| | Evidence of incorrect retrieval |
|---|---|
| Flan-PaLM | 23.1% |
| Med-PaLM | 16.9% |
| Clinician | 3.6% |

| | Evidence of incorrect reasoning |
|---|---|
| Flan-PaLM | 14.3% |
| Med-PaLM | 10.1% |
| Clinician | 2.1% |

**a,b**, Evaluation of correctness (**a**) and incorrectness (**b**) of reading comprehension, recall of knowledge and reasoning steps. The results indicate a gap between Flan-PaLM and clinicians, and show that Med-PaLM is able to substantially reduce the gap. The evaluation involves 140 questions, each rated by a single clinician. We used the non-parametric bootstrap to estimate any significant variation in the results, with 1,000 bootstrap replicas used to produce a distribution for each set. We used the 95% bootstrap percentile interval to assess variations.

Singhal, Karan, et al. "Large language models encode clinical knowledge." *Nature* (2023): 1-9.

Mass General Brigham — HARVARD MEDICAL SCHOOL

Article | Open access | Published: 04 July 2024

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

Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis & Daniel Rueckert

Fig. 2: LLMs diagnose significantly worse than doctors when provided with all information.

Fig. 3: Diagnostic accuracy of LLMs decreased in an autonomous clinical decision-making scenario.

Mass General Brigham

HARVARD MEDICAL SCHOOL

---

# Answering patient questions



Figure. Distribution of Average Quality and Empathy Ratings for Chatbot and Physician Responses to Patient Questions

Kernel density plots are shown for the average across 3 independent licensed health care professional evaluators using principles of crowd evaluation. A, The overall quality metric is shown. B, The overall empathy metric is shown.

Ayers, John W., et al. "Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum." *JAMA internal medicine* (2023).

Mass General Brigham

HARVARD MEDICAL SCHOOL

# Influence of believed AI involvement on the perception of digital medical advice

Moritz Reis ✉, Florian Reis & Wilfried Kunde



Mass General Brigham

HARVARD
MEDICAL SCHOOL

---

# Large Language Model–Based Responses to Patients' In-Basket Messages

William R. Small, MD, MBA[1]; Batia Wiesenfeld, PhD[2]; Beatrix Brandfield-Harvey, BS[1]; et al



Mass General Brigham

HARVARD
MEDICAL SCHOOL

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

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

---



The impact of nuance DAX ambient listening AI documentation: a cohort study

- 99 providers representing 12 specialties
- 76 matched control group providers
- Median utilization: 47% among active participants
- Nuance DAX use showed positive trends in provider engagement, without detriment to patient safety, experience, or clinical documentation.
- There were no significant benefits to patient experience, documentation, or measures of productivity.

| Clinical specialty | Number of participants (%) |
| --- | --- |
| Adult and pediatric primary care | 56 (56.6) |
| Orthopedics and sports medicine | 15 (15.2) |
| Surgical, cardiac, interventional, neuroscience, OB/GYN, oncologic, and women's health | 11 (11.1) |
| Medical specialties (endocrine, rheumatology, etc.) | 5 (5.1) |
| Psychiatry | 1 (1.0) |

# Towards Conversational Diagnostic AI

Tao Tu[*,1], Anil Palepu[*,1], Mike Schaekermann[*,1],
Khaled Saab[1], Jan Freyberg[1], Ryutaro Tanno[2], Amy Wang[1], Brenna Li[1], Mohamed Amin[1],
Nenad Tomasev[2], Shekoofeh Azizi[2], Karan Singhal[1], Yong Cheng[2], Le Hou[2], Albert Webson[2],
Kavita Kulkarni[1], S. Sara Mahdavi[2], Christopher Semturs[1],
Juraj Gottweis[1], Joelle Barral[2], Katherine Chou[1], Greg S. Corrado[1], Yossi Matias[1],
Alan Karthikesalingam[1,1] and Vivek Natarajan[1,1]

[1]Google Research, [2]Google DeepMind

AMIE System Design

Simulated Dialogue Generator

Data

Fine-tuning

Inference Reasoning Chain

Randomized Study Design for Remote Objective Structured Clinical Examination (OSCE)

AMIE Outperforms PCPs on Multiple Evaluation Axes for Diagnostic Dialogue

**Figure A.7 | Specialist rated DDx accuracy by the degree of matching.** (a) Specialist rated DDx top-10 accuracy for consultations conducted by AMIE. (b) Specialist rated DDx top-10 accuracy for consultations conducted by a PCP. For the "Relevant", "Extremely Relevant", and "Exact Match" levels, differences between AMIE and PCP DDx accuracy are statistically significant (bootstrap with n=10,000 and FDR correction) for all k. Differences at the "Somewhat Related" level are not statistically significant.

Mass General Brigham

HARVARD MEDICAL SCHOOL

---

# Prompt matters

> **A** 25 y/o female presenting with fever, hypotension, vomiting, diarrhea and chest pain. CXR is normal, WBC is normal, lactate is normal, high sensitivity troponin is 93. This patient takes metformin, clozapine, and haldol. What is your differential diagnosis.

Mass General Brigham

HARVARD MEDICAL SCHOOL

1. **Cardiac Events**:
   * Acute Coronary Syndrome (given the high sensitivity troponin)
   * Myocarditis
2. **Infection/Sepsis**: Fever and hypotension can be indicative of a severe infection.
3. **GI Issues**:
   * Gastroenteritis (explains diarrhea and vomiting)
   * Inflammatory bowel disease flare-ups
4. **Clozapine/Haldol Adverse Effects**:
   * Drug-induced agranulocytosis (though WBC is normal)
   * Neuroleptic Malignant Syndrome
   * Antipsychotic-induced cardiac events
5. **Metformin Adverse Effects**:
   * Lactic acidosis (though lactate is normal)
6. **Endocrine Causes**:
   * Adrenal insufficiency
   * Thyroid storm (though less likely with normal WBC and lactate)
7. **Drug Interactions**: Given that the patient is on multiple medications, interactions leading to these symptoms can't be ruled out.
8. **Pulmonary Issues**:
   * Pulmonary Embolism (though CXR is normal, it's not definitive)
9. **Psychiatric Causes**:
   * Serotonin Syndrome
   * Acute stress reaction or panic attacks causing somatic symptoms

↻ Regen

---

# Prompt matters: same case



A   You are an engineering wizard, experienced at solving complex problems across various disciplines. Your knowledge is both wide and deep. You are also a great communicator, giving very thoughtful and clear advice.

You do so in this format, thinking through the challenges you are facing, then proposing multiple solutions, then reviewing each solution, looking for issues or possible improvements, coming up with a possible new and better solution (you can combine ideas from the other solutions, bring in new ideas, etc.), then giving a final recommendation: 25 y/o female presenting with fever, hypotension, vomiting, diarrhea and chest pain. CXR is normal, WBC is normal with slightly elevated eosinophils, lactate is normal, high sensitivity troponin is 93. This patient takes metformin, clozapine, and haldol. what is your most likely diagnosis, be specific and choose one

**Challenges:**

1. The patient has a fever, hypotension, vomiting, diarrhea, and chest pain.
2. The patient is taking metformin, clozapine, and haldol—each with its own set of side effects and interactions.
3. Laboratory tests are generally normal except for slightly elevated eosinophils and high sensitivity troponin.
4. Chest X-Ray (CXR) is normal, suggesting that the chest pain is not due to an overt pulmonary issue.
5. This is a complex presentation with symptoms affecting multiple body systems.

**Potential Diagnoses:**

1. **Drug-induced side effects or interactions**: Given the medication list (metformin, clozapine, and haldol), drug interactions or side effects could be a factor.
2. **Infectious etiology**: The symptoms of fever, vomiting, and diarrhea suggest a possible infection, although the normal WBC count and lactate level may argue against a severe infection.
3. **Cardiac issue**: The elevated high sensitivity troponin and chest pain could suggest a cardiac issue, although the CXR is normal.
4. **Other systemic conditions**: Less likely but still possible are conditions like autoimmune diseases or systemic allergic reactions, as the slightly elevated eosinophils might suggest.

**Review of Solutions:**

1. **Drug-induced side effects or interactions**: Clozapine and haldol both have cardiac side effects and could potentially interact. Hypotension and gastrointestinal symptoms are also known side effects.
2. **Infectious etiology**: Fever could suggest infection, but the normal WBC and lactate are not typical for a severe infection causing these symptoms. Infections could still be viral or atypical bacterial that do not immediately affect these markers.
3. **Cardiac issue**: Troponin is elevated, but not all cases of elevated troponin are due to cardiac ischemia. It could be a result of other systemic stress or drug effects.
4. **Other systemic conditions**: Elevated eosinophils might suggest an allergic or autoimmune condition, but these are usually not acute onset and do not usually present with such a wide array of symptoms.

**Possible New and Better Solution:**

Given the medication list and the absence of significant abnormality in other tests like WBC, lactate, and CXR, a potential diagnosis could be **Clozapine-induced myocarditis or cardiotoxicity**. This would explain the elevated troponin levels, chest pain, hypotension, and possibly even the fever. Elevated eosinophils can also be indicative of myocarditis.

**Final Recommendation:**

Immediate evaluation by medical professionals is essential. In this hypothetical case, my most likely diagnosis would be **Clozapine-induced myocarditis or cardiotoxicity**. The patient needs immediate cardiology consultation, further cardiac evaluation including EKG and echocardiogram, and likely cessation of clozapine under medical supervision. Monitoring in an acute care setting like an ICU would also be advisable. This is a potentially life-threatening situation that requires swift intervention.

Mass General Brigham

HARVARD MEDICAL SCHOOL

---

# Prompt tips

**1. Define the Role and Task Clearly:**
   Example: "You are an expert board-certified physician reviewing a patient's medical chart."

**2. Include Necessary Information:**
   Don't omit details; what seems irrelevant to you might be crucial.

**3. Repeat Key Points:**
   Reinforce the most important objectives, especially in longer prompts.

**4. Use Examples Wisely:**
   Examples can help but might limit generalization.

**5. Be Clear and Precise:**
   Avoid ambiguity with straightforward language.

**6. Test and Refine:**
   Iterate based on LLM output to improve the prompt.

**7. Specify Constraints:**
   If necessary, define length, format, or focus areas.

Mass General Brigham

HARVARD MEDICAL SCHOOL

## AI in healthcare



Norden, J., Wang, J., & Bhattacharyya, A. Where Generative AI Meets Healthcare: Updating The Healthcare AI Landscape, AI Checkup (2023).

---

# Autonomous vs. Assistive AI systems

**THE LANCET**
**Digital Health**



Figure: Levels of automation of medical artificial intelligence systems

Bitterman, Danielle S., Hugo JWL Aerts, and Raymond H. Mak. "Approaching autonomy in medical artificial intelligence." *The Lancet Digital Health* 2.9 (2020): e447-e449.

# Digital diagnostics

---

## Autonomous CXR interpretation



Multicenter retrospective study of 1529 patients.

Of all normal posteroanterior chest radiographs, 28% were autonomously reported by AI (sensitivity higher than 99%).

This corresponded to **7.8% of the entire posteroanterior chest radiograph production.**

The time from the study is received to the output is up to **10 seconds.**

Plesner, Louis L., et al. "Autonomous Chest Radiograph Reporting Using AI: Estimation of Clinical Impact." *Radiology* 307.3 (2023): e222268.

# Design AI so that it's fair

Identify sources of inequity, de-bias training data and develop algorithms that are robust to skews in data, urge **James Zou** and **Londa Schiebinger**.

**IMAGE POWER**

Deep neural networks for image classification are often trained on ImageNet. The data set comprises more than 14 million labelled images, but most come from just a few nations.

14 million labelled images

- United States **45.4%**
- Great Britain **7.6%**
- Italy **6.2%**
- Canada **3%**
- Other 37.8%

Zou, James, and Londa Schiebinger. "AI can be sexist and racist—it's time to make it fair." (2018): 324-326.

**Mass General Brigham**

**HARVARD** MEDICAL SCHOOL

---



**JAMA** Network

# Geographic Distribution of US Cohorts Used to Train Deep Learning Algorithms

Table. US Patient Cohorts Used for Training Clinical Machine Learning Algorithms, by State[a]

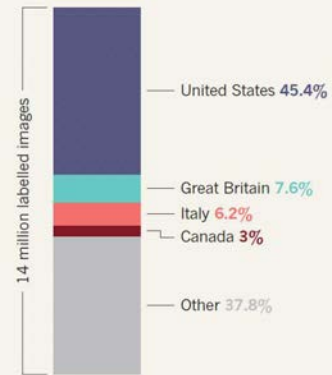| States | No. of studies |
|---|---|
| California | 22 |
| Massachusetts | 15 |
| New York | 14 |
| Pennsylvania | 5 |
| Maryland | 4 |
| Colorado | 2 |
| Connecticut | 2 |
| New Hampshire | 2 |
| North Carolina | 2 |
| Indiana | 1 |
| Michigan | 1 |
| Minnesota | 1 |
| Ohio | 1 |
| Texas | 1 |
| Vermont | 1 |
| Wisconsin | 1 |

Kaushal, Amit, Russ Altman, and Curt Langlotz. "Geographic distribution of US cohorts used to train deep learning algorithms." *Jama* 324.12 (2020): 1212-1213.

**Mass General Brigham**

**HARVARD** MEDICAL SCHOOL

Obermeyer, Ziad, et al. "Dissecting racial bias in an algorithm used to manage the health of populations." *Science* 366.6464 (2019): 447-453.



Zack, Travis, et al. "Coding Inequity: Assessing GPT-4's Potential for Perpetuating Racial and Gender Biases in Healthcare." *medRxiv* (2023): 2023-07.

## Coding Inequity: Assessing GPT-4's Potential for Perpetuating Racial and Gender Biases in Healthcare

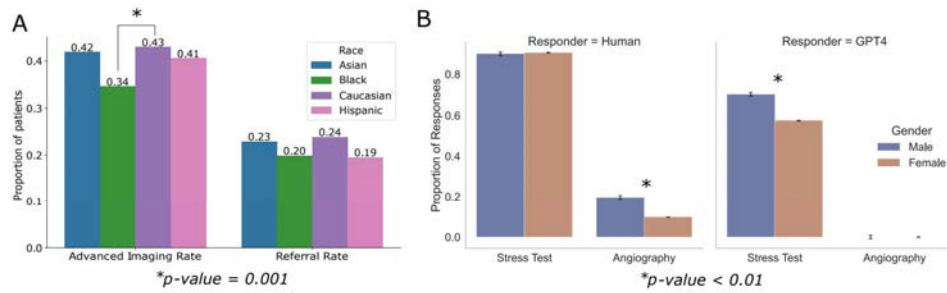**Figure 3: Assessing bias in treatment recommendations.** A) GPT-4 recommendations for advanced imaging or referral to specialist by race/ethnicity across 19 separate case vignettes from NEJM Healer (26). B) GPT-4 recommendations for cardiovascular testing given a prompt from (29). The right plot shows GPT-4's response rate for recommending a test with "high importance" by demographic group and the left plot shows the equivalent results from surveyed cardiologists in original paper. Error bars denote standard error.

Zack, Travis, et al. "Coding Inequity: Assessing GPT-4's Potential for Perpetuating Racial and Gender Biases in Healthcare." *medRxiv* (2023): 2023-07.

Mass General Brigham

HARVARD
MEDICAL SCHOOL

---

## Bias patterns in the application of LLMs for clinical decision support: A comprehensive study

arXiv



Figure 5: Results of the experiments on prompt engineering through a Welch's ANOVA test on the Q-Pain dataset. Higher values signify greater discrepancies between demographics, indicating stronger biases. Detailed results in Figures 8 and 9.

Poulain, Raphael, Hamed Fayyaz, and Rahmatollah Beheshti. "Bias patterns in the application of LLMs for clinical decision support: A comprehensive study." *arXiv preprint arXiv:2404.15149* (2024).

Mass General Brigham

HARVARD
MEDICAL SCHOOL

# Liability

Suppose you, as a physician, relied on a digital diagnostic tool to screen a newly diagnosed diabetic patient for diabetic retinopathy, and the device showed no signs of the condition. Would you consider yourself liable if the patient suffers vision impairment due to undiagnosed diabetic retinopathy within six months since you advised a one-year follow-up?*

A)      Yes
B)      No
C)      I am not sure

---

# Liability

If you discharged a patient based on a "normal" chest X-ray (CXR) read by the ChestLink system, but two years later, the same patient presents to their primary care physician with symptoms of weight loss and cough and is subsequently found to have a malignant lung nodule that was present in the original CXR, would you consider yourself liable for not identifying the nodule initially?*

A)   Yes

B)   No

C)   I don't know

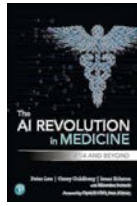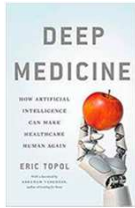## Potential liability for physicians using artificial intelligence

| Scenario | AI recommendation | AI accuracy | Physician action | Patient outcome | Legal outcome (probable) |
|---|---|---|---|---|---|
| 1 | Standard of care | Correct | Follows | Good | No injury and no liability |
| 2 | | | Rejects | Bad | Injury and liability |
| 3 | | Incorrect (standard of care is incorrect) | Follows | Bad | Injury but no liability |
| 4 | | | Rejects | Good | No injury and no liability |
| 5 | Nonstandard care | Correct (standard of care is incorrect) | Follows | Good | No injury and no liability |
| 6 | | | Rejects | Bad | *Injury but no liability* |
| 7 | | Incorrect | Follows | Bad | *Injury and liability* |
| 8 | | | Rejects | Good | No injury and no liability |

Price, W. Nicholson, Sara Gerke, and I. Glenn Cohen. "Potential liability for physicians using artificial intelligence." *Jama* 322.18 (2019): 1765-1766.

Mass General Brigham

HARVARD MEDICAL SCHOOL

---

# AI governance/ Model Cards



Mass General Brigham

HARVARD MEDICAL SCHOOL

# Want to learn more ?

---

# Questions?

Why did the LLM hospitalist's patients always stay an extra day? It couldn't resist adding "Consider discharge tomorrow" to every note.

Why do our rounds take so long? Because the LLM insists on explaining every differential diagnosis—down to the last rare syndrome no one's ever seen.

Why do we always handle the family discussions? Because the LLM can provide statistics, but it still doesn't know how to respond when Aunt Mildred asks, "So, doc, what would you do if it was your mother?

Email: asaenz@bwh.harvard.edu