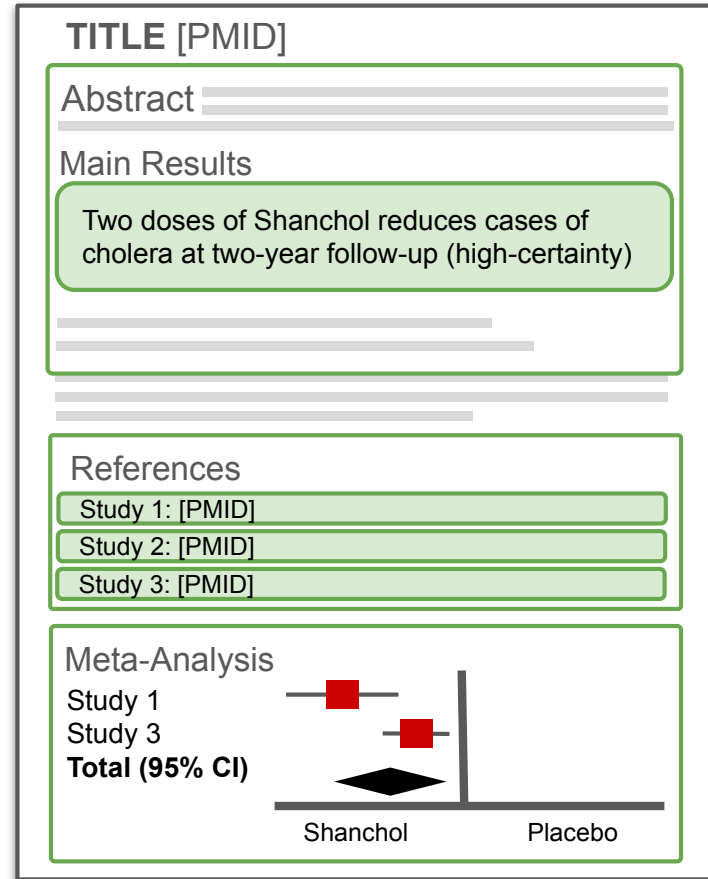
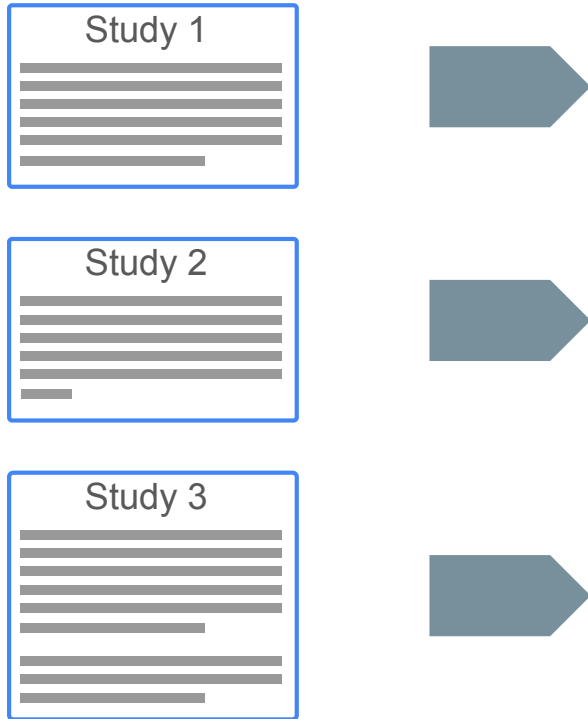




# Can Large Language Models Match the Conclusions of Systematic Reviews?



# What is systematic review?



# Assessing core LLM capabilities

## Domain Expertise

*Oncology & Hematology*

**Q:** Is the overall survival rate higher, lower, or the same when comparing percutaneous ethanol injection (PEI) to percutaneous acetic acid injection (PAI)?

*Psychiatry & Neurology*

**Q:** Is fatigue severity higher, lower, or the same when comparing doxepin to placebo?

*Internal Medicine & Subspecialties*

**Q:** Is patient function measured by HAQ score higher, lower, or the same when comparing biologic monotherapy to placebo?

*Pediatrics & Neonatology*

**Q:** Is the incidence of bronchopulmonary dysplasia (BPD) higher, lower, or the same when comparing prophylactic CPAP to very early CPAP?

## Conflict Resolution

Studies

**Q:** Given these 3 studies, is survival to hospital discharge higher, lower, or the same when comparing untrained bystander CPR with continuous chest compression to untrained bystander CPR with chest compression interrupted with pauses for rescue breathing?

LLM

**Rationale:** survival to hospital discharge is **not significantly different...as shown by similar survival rates in all three studies**  
**Answer:** no difference (x)

**Human Answer:** higher

**Human Rationale:** "There were no significant differences between the two groups in the trials. The **pooled result showed better survival** for the continuous chest compression alone group (RR 1.21, 95% CI 1.01 to 1.46..."

## Scientific Skepticism

Studies

**Q:** Given the 2 studies, Is the risk of cryptococcal IRIS events higher, lower, or the same when comparing early ART initiation to delayed ART initiation?

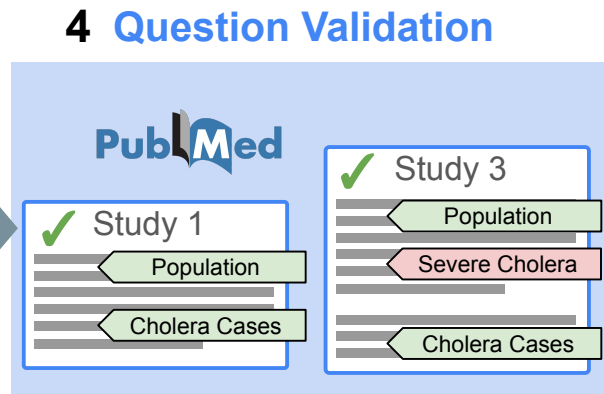
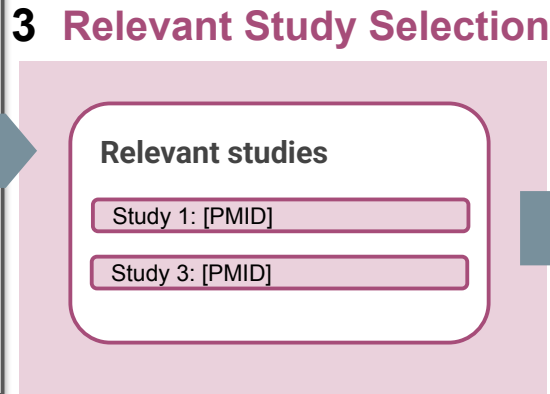
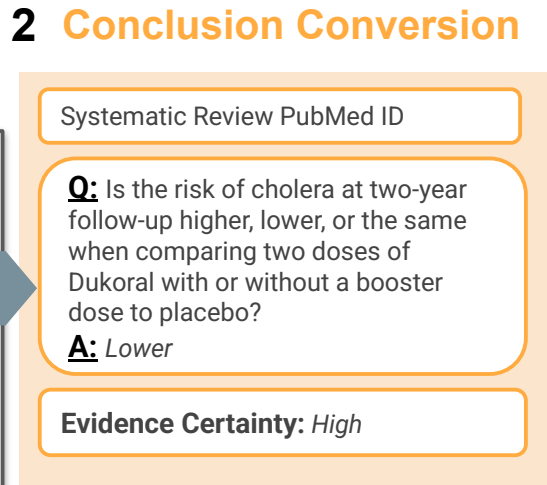
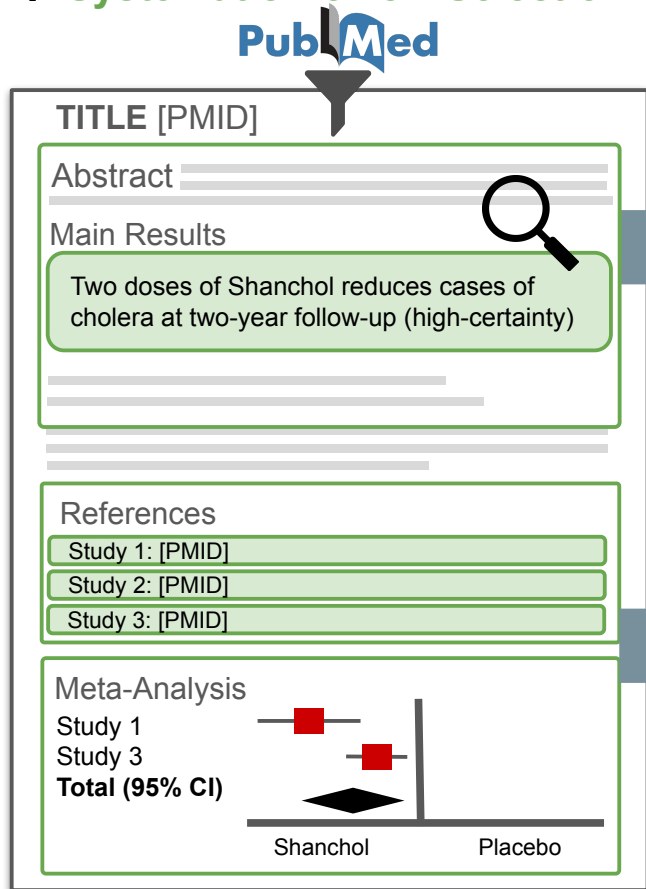
LLM

**Rationale:** ...Article 1 found no statistically significant difference...However, article 2 found a significantly higher risk...the **evidence suggests that the risk is higher**  
**Model Answer:** higher (x)

**Human Answer:** uncertain effect

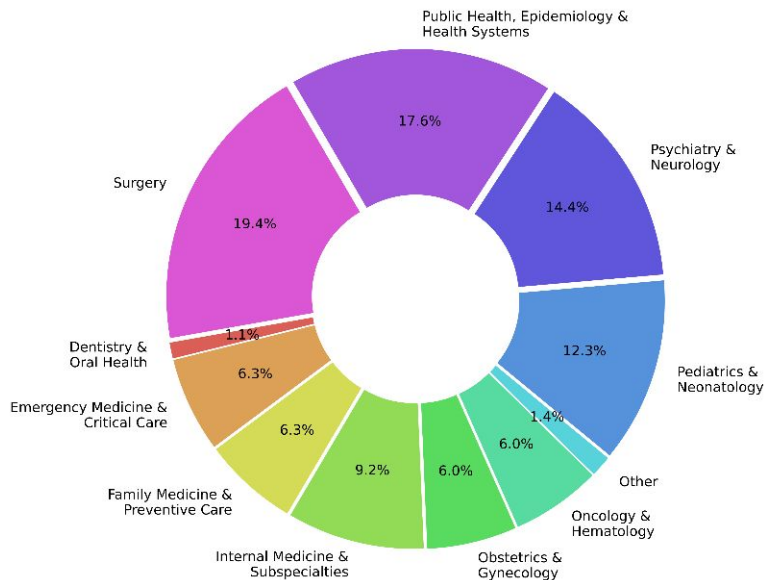
**Human Rationale:** "We are uncertain...[due to] **wide CIs** and **very few clinical events**...IRIS outcome assessors were **unblinded**... **diagnosing** IRIS can be **very subjective**"

# 1 Systematic Review Selection 2 Conclusion Conversion

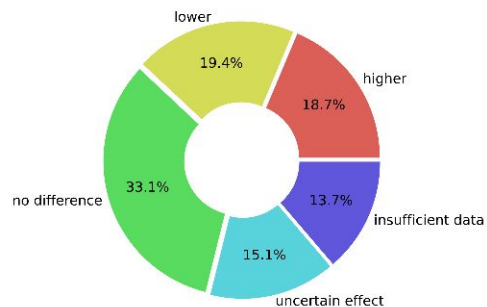


# Dataset statistics

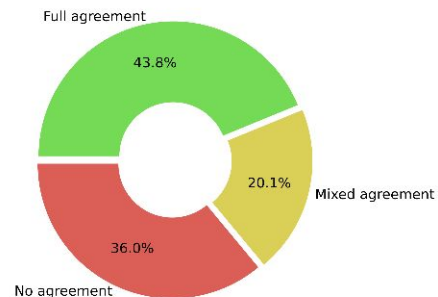
(a) Dataset stratified by medical specialty

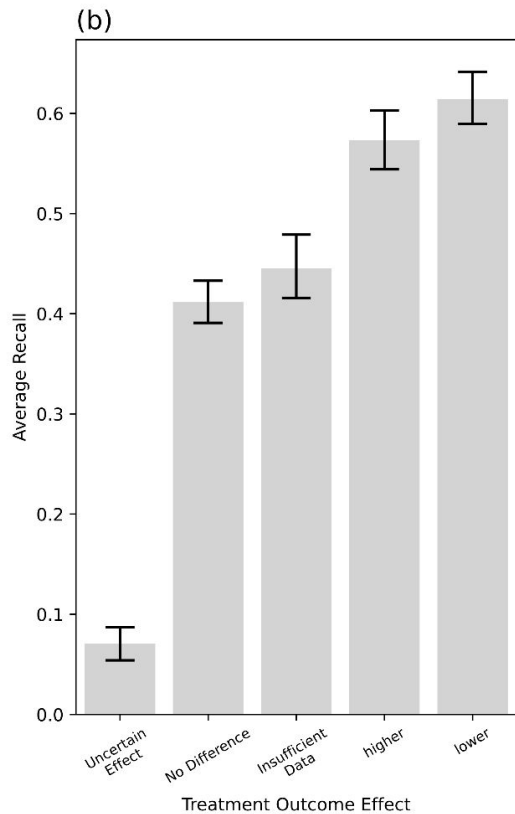
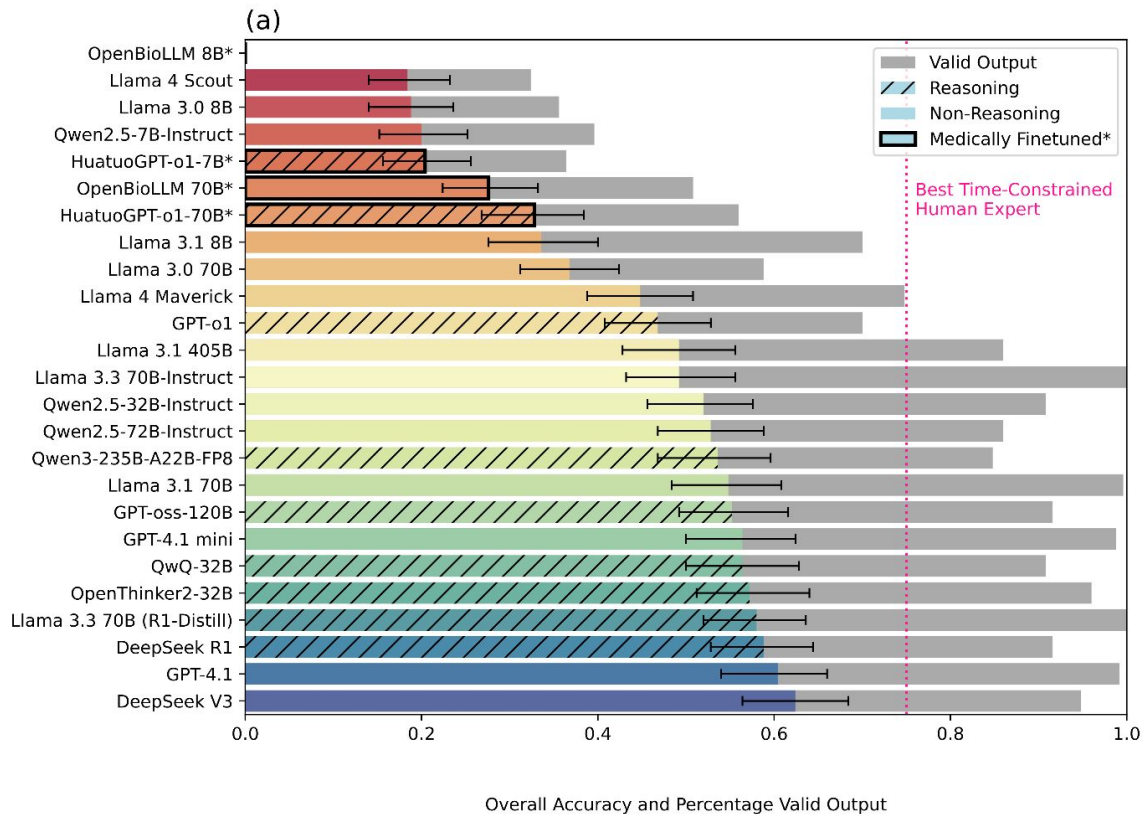


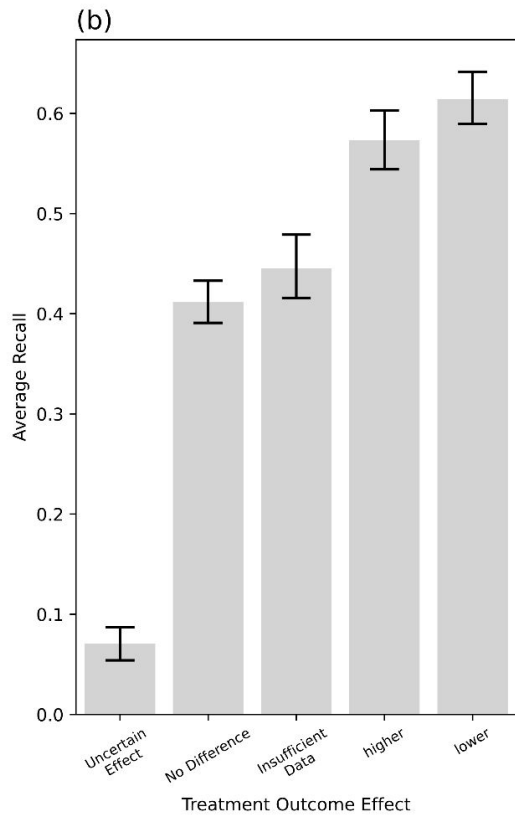
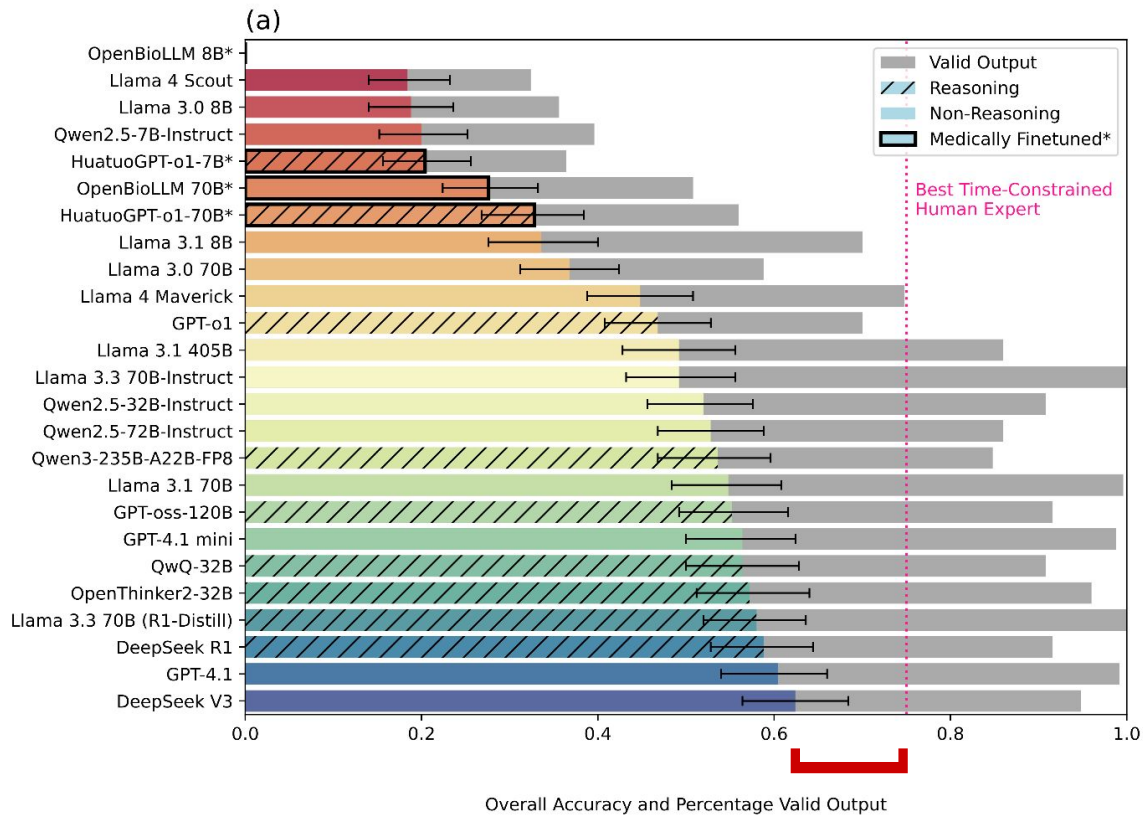
(b) Dataset stratified by treatment outcome effect

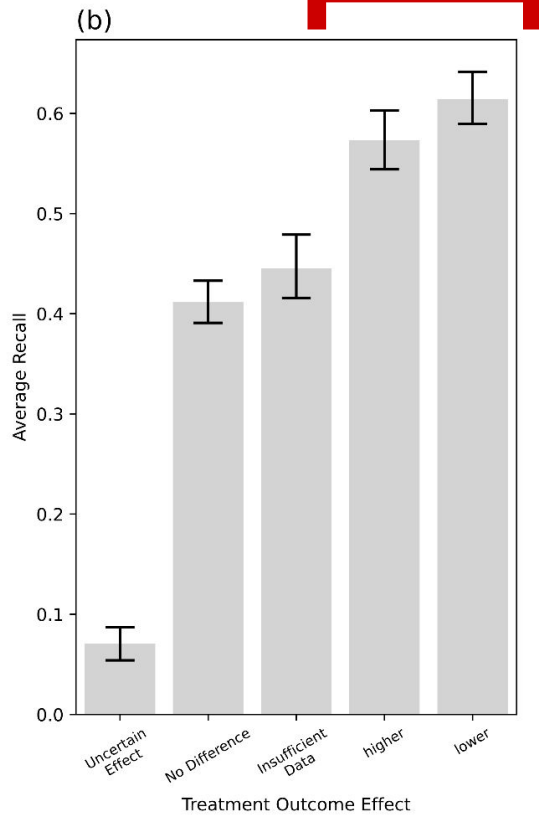
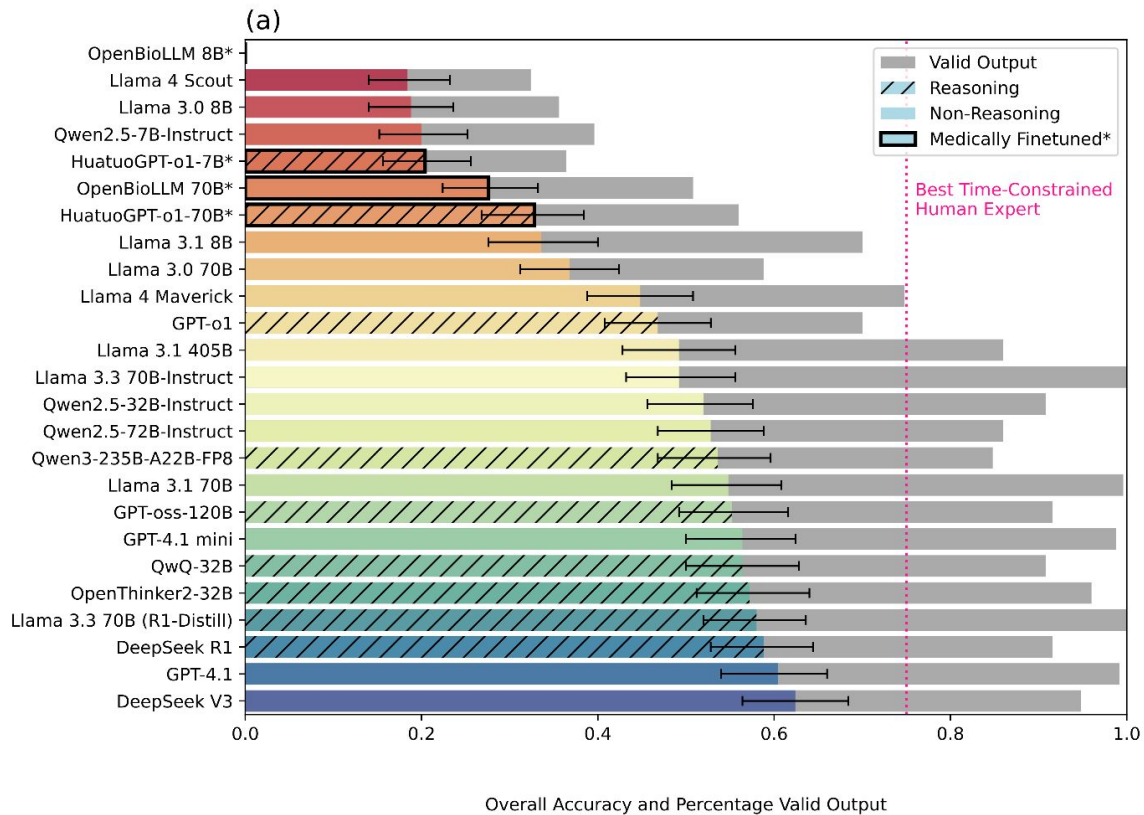


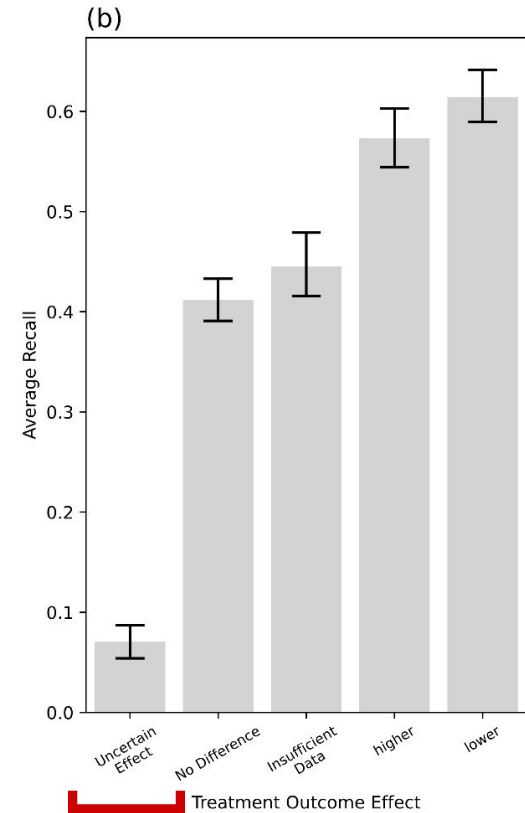
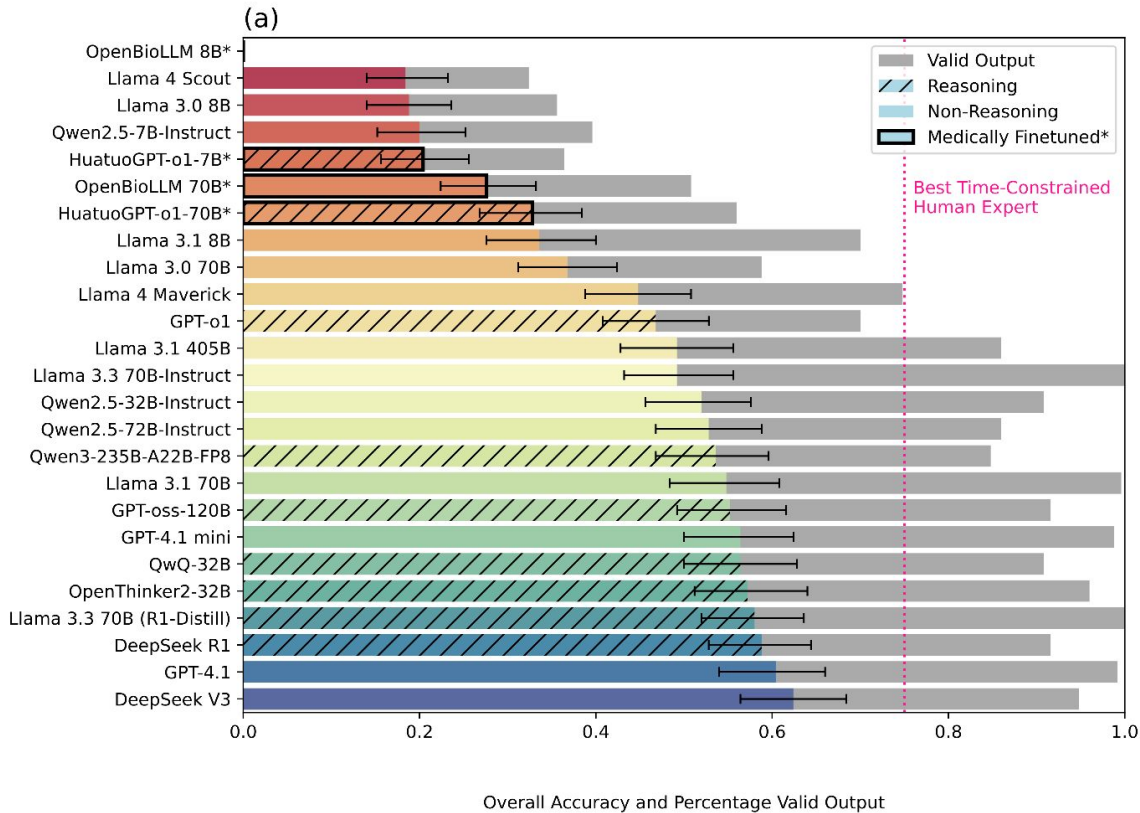
(c) Dataset stratified by source concordance with correct answer

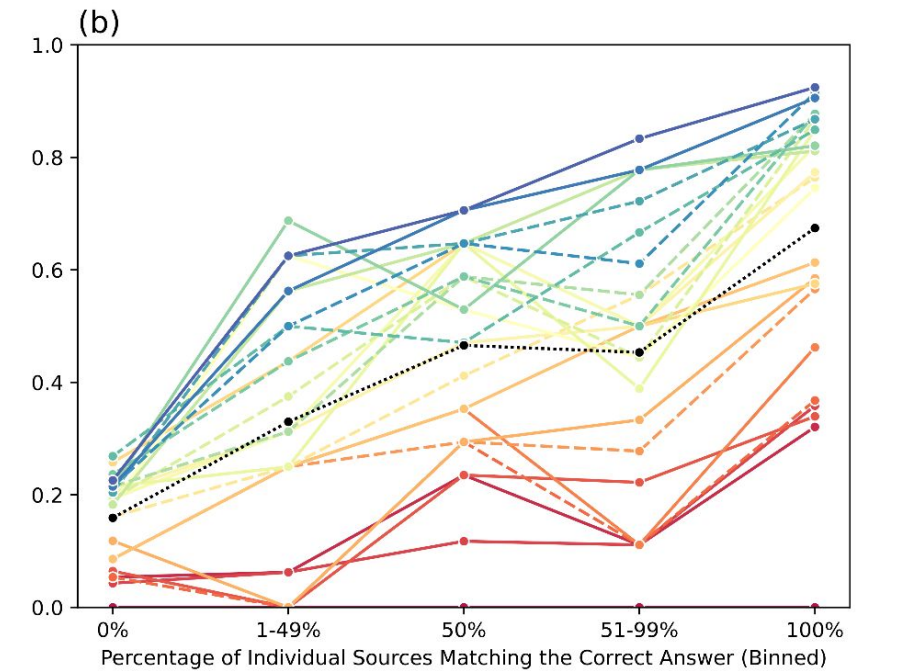
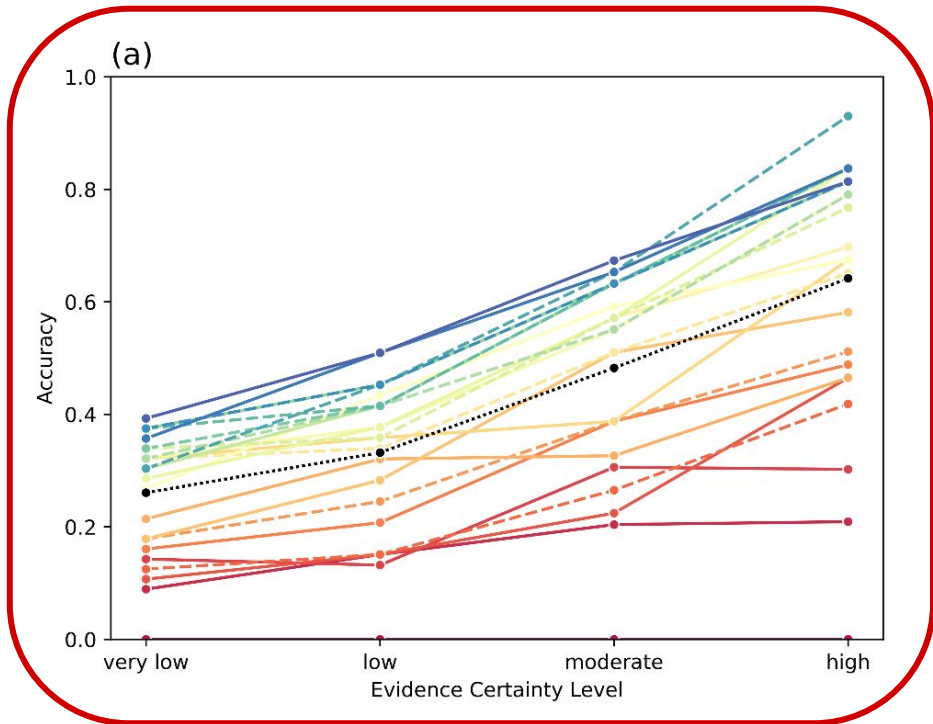




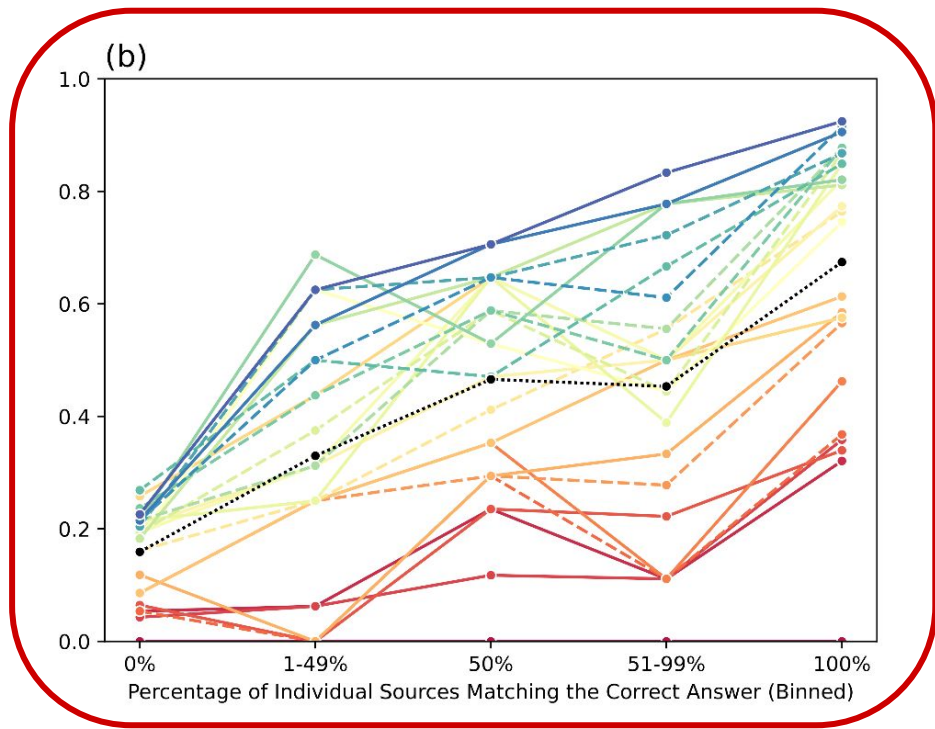
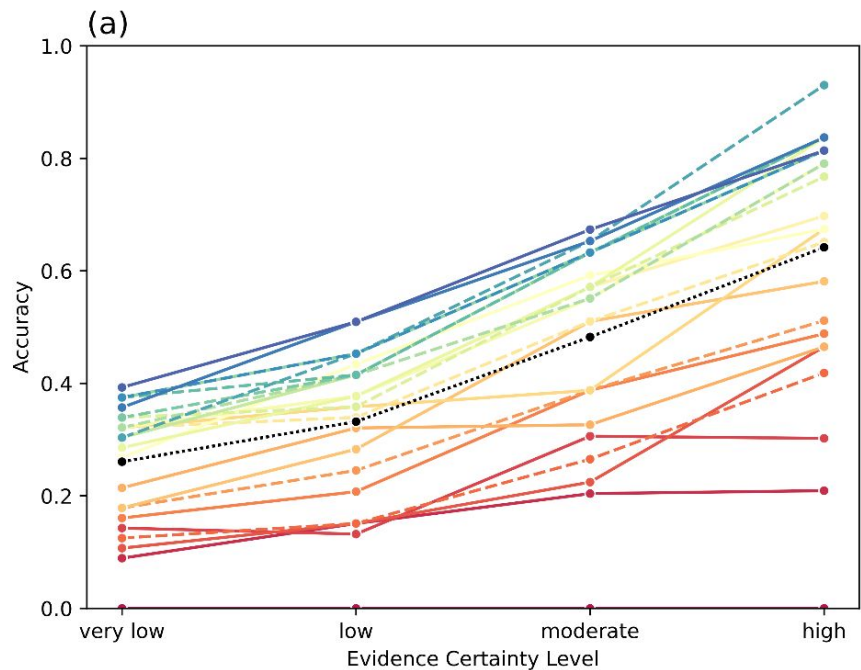


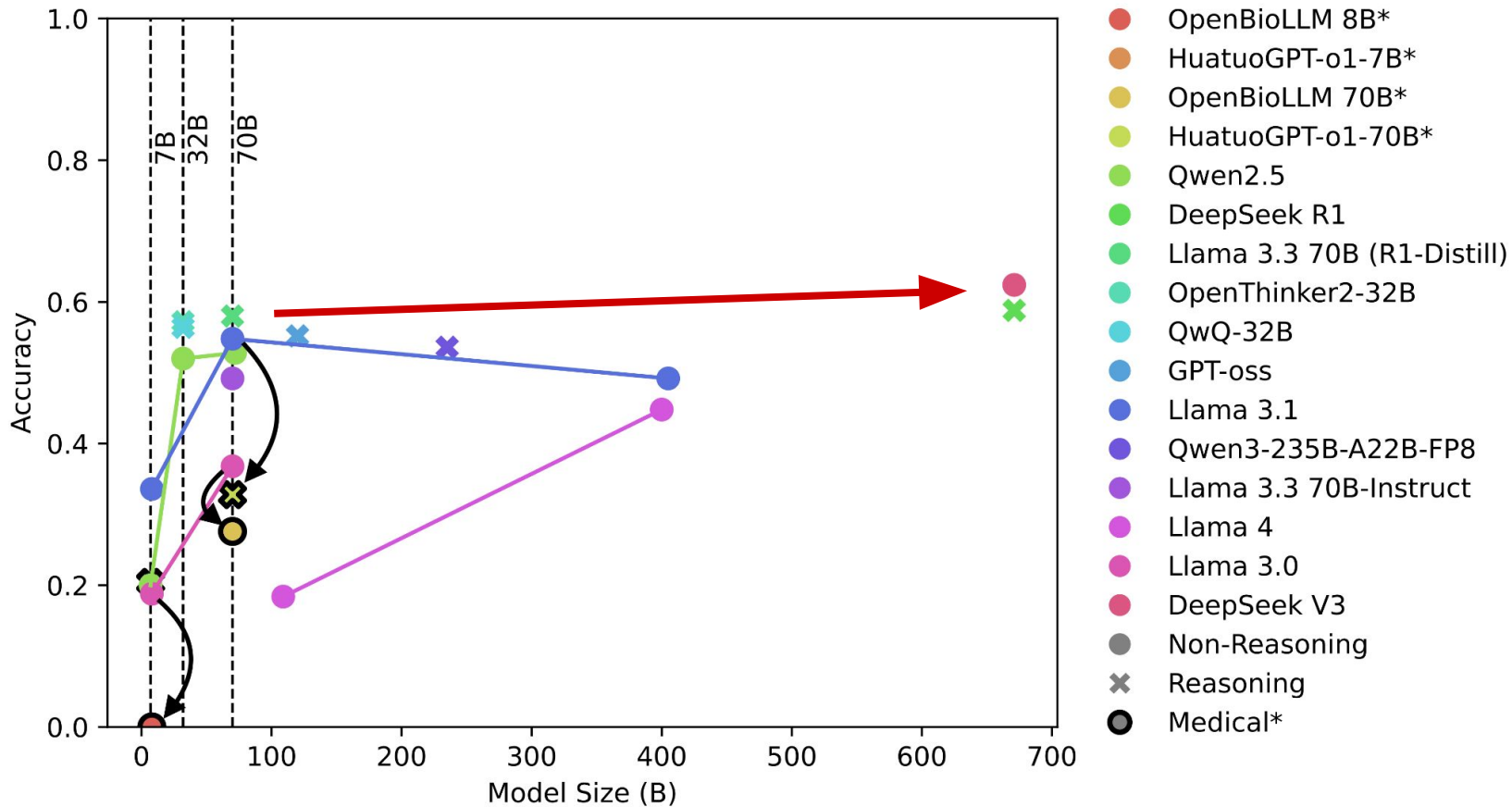






- |                     |                   |                        |                     |                            |               |
|---------------------|-------------------|------------------------|---------------------|----------------------------|---------------|
| OpenBioLLM 8B*      | OpenBioLLM 70B*   | GPT-o1                 | Qwen3-235B-A22B-FP8 | OpenThinker2-32B           | DeepSeek V3   |
| Llama 4 Scout       | HuatuogPT-o1-70B* | Llama 3.1 405B         | Llama 3.1 70B       | Llama 3.3 70B (R1-Distill) | Non-Reasoning |
| Llama 3.0 8B        | Llama 3.1 8B      | Llama 3.3 70B-Instruct | GPT-oss-120B        | DeepSeek R1                | Reasoning     |
| Qwen2.5-7B-Instruct | Llama 3.0 70B     | Qwen2.5-32B-Instruct   | GPT-4.1 mini        | GPT-4.1                    | Average       |
| HuatuogPT-o1-7B*    | Llama 4 Maverick  | Qwen2.5-72B-Instruct   | QwQ-32B             |                            |               |





# CAN LARGE LANGUAGE MODELS MATCH THE CONCLUSIONS OF SYSTEMATIC REVIEWS?

**Christopher Polzak\***<sup>1</sup>

**Alejandro Lozano\***<sup>1</sup>

**Min Woo Sun\***<sup>1</sup>

**James Burgess**<sup>1</sup>

**Yuhui Zhang**<sup>1</sup>

**Kevin Wu**<sup>1</sup>

**Chia-Chun Chiang**<sup>2</sup>

**Jeffrey J. Nirschl**<sup>1</sup>

**Serena Yeung-Levy**<sup>1</sup>

<sup>1</sup>Stanford University

<sup>2</sup>Mayo Clinic

<https://zy-f.github.io/website-med-evidence/>