



# Frontier Large Language Models (LLMs) for the Practicing MFM and Genetics Clinician: A Practical Guide to Model Selection in 2026

Christian Macedonia, M.D.

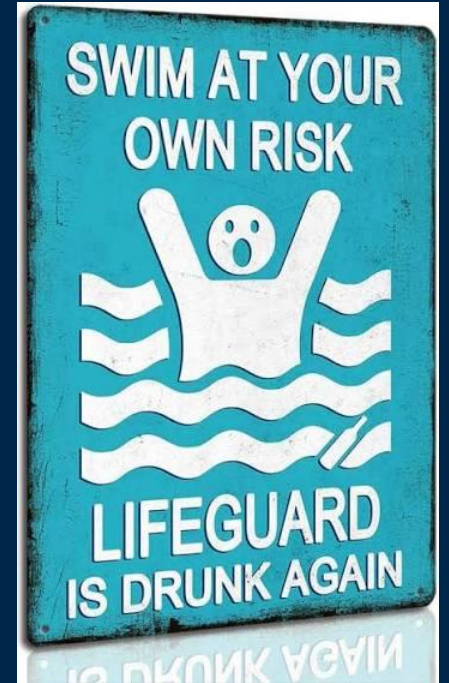
University of Michigan

USA

CEO Lancaster Maternal Fetal Medicine LLC

# Disclosure and Safety Statement

- No financial relationships or conflicts of interest to disclose.
- The following discussion is educational and does not constitute endorsement of any commercial product.
- Audience members remain responsible for compliance with institutional policies, HIPAA/GDPR requirements, and applicable privacy regulations.
- LLM outputs should be independently verified before use in clinical care.
- **Never upload identifiable patient data**
- **Assume every prompt is discoverable**



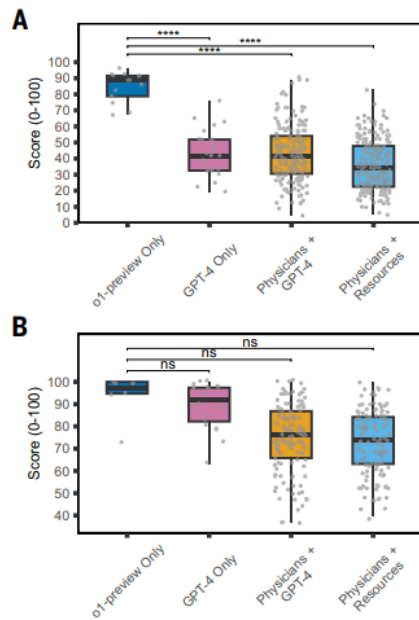
This talk can be found at:

<https://thektproject.org/>

# On many standardized medical reasoning tasks, frontier models now outperform most physicians

**Fig. 4. Comparison of o1-preview, GPT-4, and physicians for management and diagnostic reasoning.**

**(A)** Box plot of normalized management reasoning points by LLMs and physicians on Grey Matters Management Cases. Five cases were included. We generated three o1-preview responses for each case. The prior study collected five GPT-4 responses to each case, 178 completed cases from 46 physicians with access to GPT-4, and 197 completed cases from 46 physicians with access to conventional resources. \* $P \leq 0.05$ , \*\* $P \leq 0.01$ , \*\*\* $P \leq 0.001$ , \*\*\*\* $P \leq 0.0001$ . **(B)** Box plot of normalized diagnostic reasoning points by LLMs and physicians. Six diagnostic challenges were included. We generated one o1-preview response for each case. The prior study collected three GPT-4 responses to all cases, 125 cases completed by 25 physicians with access to GPT-4, and 119 cases completed by 25 physicians with access to conventional resources.



Science 30 APRIL 2026

Peter G. Brodeur et al., *Performance of a large language model on the reasoning tasks of a physician*. **Science** 392,524-527 (2026). DOI:10.1126/science.adz4433



DOI:10.1056/Alcs2501343

## CASE STUDY

### LLM-Assisted Reanalysis of Unsolved Rare Disease Genomes Increases Diagnostic Yield

Aaron Jaech<sup>1</sup>, Ph.D.,<sup>1</sup> Morgan Cheatham<sup>2</sup>, M.D.,<sup>3</sup> Suyash S. Shringarpure<sup>4</sup>, Ph.D.,<sup>1</sup> Casie A. Genetti<sup>5</sup>, M.S.,<sup>4</sup> Pratiksha Pradhan<sup>6</sup>, M.S.,<sup>4</sup> Aarti Bagul<sup>7</sup>, M.S.,<sup>1</sup> Trishan Panch<sup>8</sup>, M.D.,<sup>5</sup> Benjamin Rader<sup>9</sup>, M.D., M.P.H.,<sup>6,7</sup> Kara Sewalk<sup>10</sup>, M.P.H.,<sup>6</sup> Rebecca Distler<sup>11</sup>, M.P.H.,<sup>1</sup> Katherine N. Anderson<sup>12</sup>, M.S.,<sup>4</sup> Melissa Stewart<sup>13</sup>, M.P.H.,<sup>6</sup> Sarah Leifer<sup>14</sup>, M.P.H.,<sup>6</sup> David C. Glahn<sup>15</sup>, Ph.D.,<sup>16,17</sup> Richard D. Goldstein<sup>18</sup>, M.D.,<sup>19</sup> Monica Wojcik<sup>20</sup>, M.D., M.P.H.,<sup>3,4</sup> Alan H. Beggs<sup>21</sup>, Ph.D.,<sup>3,4</sup> John S. Brownstein<sup>22</sup>, Ph.D.,<sup>3,4</sup> and Catherine A. Brownstein<sup>23</sup>, Ph.D., M.P.H.<sup>3,4,8,11</sup>

Received: December 4, 2025; Revised: February 19, 2026; Accepted: March 25, 2026; Published: June 18, 2026

#### Abstract

**BACKGROUND** Rare and undiagnosed genetic disorders affect millions of patients globally, and many patients endure years of inconclusive testing. Conventional genomic interpretation can be insufficiently sensitive and costly and is rarely repeated as knowledge evolves.

**METHODS** We conducted a retrospective multicohort reanalysis using a large language model (LLM)-assisted workflow that ingests clinician notes, Human Phenotype Ontology (HPO) terms, and a filtered variant table to propose explanation-rich candidate hypotheses for expert adjudication under American College of Medical Genetics and Genomics and Association for Molecular Pathology criteria. A diagnosis was defined a priori as a variant classified as pathogenic or likely pathogenic, confirmed in a Clinical Laboratory Improvement Amendments–certified laboratory, and returned to families. Secondary outputs included “rediscoveries” of externally established diagnoses not yet available locally and hypothesis generation signals.

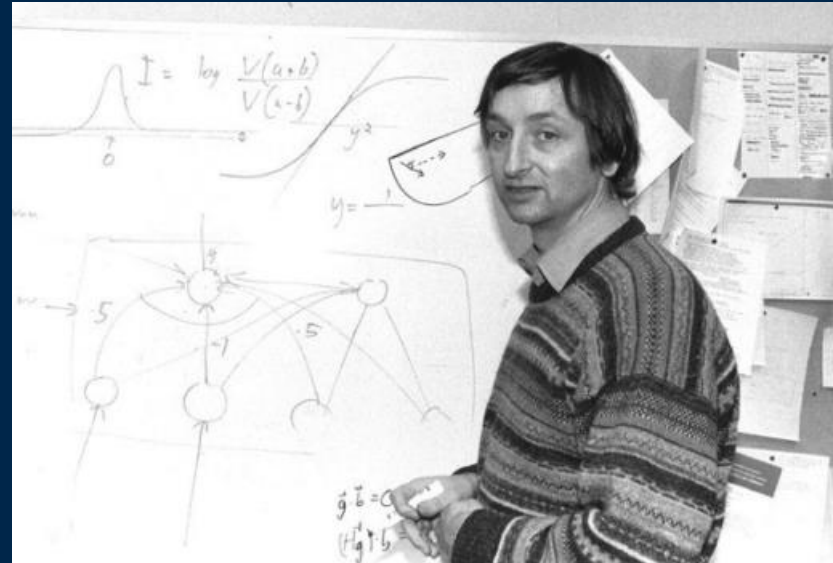
#### Methods

Briefly, annotated variant call files (VCFs) from existing short- and long-read exome and genome datasets were uploaded for patients who had undergone sequencing to diagnose suspected genetic conditions.<sup>11</sup> These were paired with metadata files containing deidentified basic demographic and clinical data, including detailed phenotypic information encoded as HPO codes. The data were then analyzed using OpenAI’s o3-deep-research model via the OpenAI application programming interface (API) or the ChatGPT (Chat Generative Pretrained Transformer) (web/app) user interface (selecting Research/Deep Research mode). Analyses used only standard platform tools (sandboxed Python/Code Interpreter and web search enabled) and were limited to user-provided files and publicly available web resources. LLM prompt development was iterative, eventually settling on two versions used to generate the data, discussed below. A sample output is provided in

Figure 1.

# Which of the Models to Choose From?

- 2022: One major *frontier* model ... GPT3
- 2026:
  - GPT
  - Claude
  - Gemini
  - Grok
  - Le Chat
  - DeepSeek
  - OpenEvidence
  - Replit
  - OpenClaw / local models (LLAMA/Meta -- Gemma/Google --DeepSeek R1)

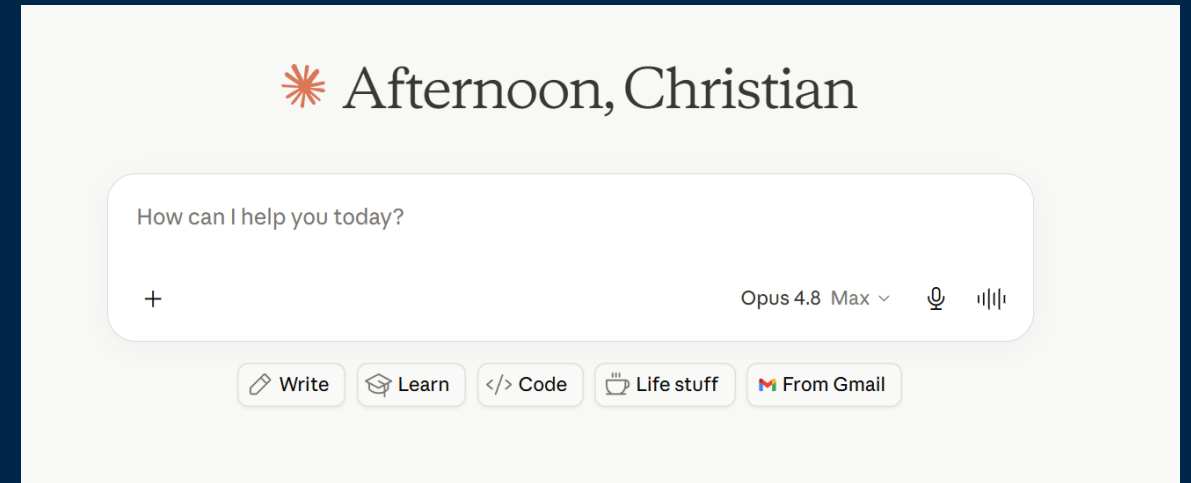


Key message: *"The best model depends on the job."*











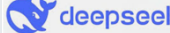





Def. Frontier Model: A big AI model that is general purpose and sits at the frontier of capability

# Common Medical Uses of AI

- Examples:
  - Differential diagnosis brainstorming
  - Literature summarization
  - Drafting patient letters
  - Insurance appeals
  - Genetic syndrome review
  - Ultrasound report editing
  - Coding support
  - Research assistance
  - Presentation preparation
  - Manuscript review
  - Grant writing assistance
  - Statistical interpretation
- Not:
  - Autonomous diagnosis
  - Autonomous medical decision making



But remember: AI is in the background, everywhere, in your machines, in your pocket, in your life...

Model		Usage Rank		Scale	Medical Utility	Cost	Advantages	Limitations
GPT		#1		Frontier	Excellent	\$\$	Best generalist, strong reasoning, projects	Can be verbose
Claude		#2		Frontier	Excellent	\$\$	Long documents, writing, nuanced analysis	Less integrated ecosystem
Gemini		#3		Frontier	Very Good	\$\$	Massive context window, Google ecosystem	Occasionally inconsistent reasoning
Grok		#4		Frontier	Good	\$\$	Current events, internet access	Less mature medical workflow
Le Chat		#5		Large	Good	\$	Fast, efficient	Smaller ecosystem
DeepSeek		#6		Large	Good	\$	Exceptional value	Variable reliability
OpenEvidence		Specialty		Proprietary	Excellent	\$\$	Medical evidence focused	Narrow scope
Replit AI		Specialty		N/A	Moderate	\$\$	Automation, workflow creation, vibe coding	Not a medical model

	Claude Mythos 5 / Fable 5	Claude Mythos Preview	Claude Opus 4.8	GPT 5.5	Gemini 3.1 Pro
<b>Agentic coding</b> SWE-Bench Pro	<b>80.3%</b>	77.8%	69.2%	58.6%	54.2%
<b>Agentic coding</b> FrontierCode (Diamond)	<b>29.3%</b> <small>xhigh</small>	—	13.4% <small>xhigh</small>	5.7% <small>xhigh</small>	—
<b>Knowledge work</b> GDPval-AA	<b>1932</b>	—	1890	1769	1314
<b>Knowledge work vision</b> GDPpdf	<b>29.8%</b> <small>no tools</small>	—	22.5% <small>no tools</small>	24.9% <small>no tools</small>	16.7% <small>no tools</small>
<b>Spatial reasoning</b> Blueprint-Bench 2	<b>38.6%</b>	—	14.5%	36.2%	26.5%
<b>Tool use</b> AutomationBench	<b>17.4%</b>	—	15.5%	12.9%	9.6%
<b>Computer use</b> OSWorld-Verified	85.0%	<b>85.4%</b>	83.4%	78.7%	76.2%
<b>Legal</b> Legal Agent Benchmark	<b>13.3%</b>	—	10.4%	2.1%	0.0%
<b>Multidisciplinary reasoning</b> Humanity's Last Exam	<b>59.0%*</b> <small>no tools</small>	56.8% <small>no tools</small>	49.8% <small>no tools</small>	41.4% <small>no tools</small>	44.4% <small>no tools</small>
	<b>64.5%*</b> <small>with tools</small>	<b>64.7%</b> <small>with tools</small>	57.9% <small>with tools</small>	52.2% <small>with tools</small>	51.4% <small>with tools</small>
<b>Biology</b> BioMysteryBench	<b>46.1%*</b> <small>hard</small>	29.6% <small>hard</small>	40.0% <small>hard</small>	—	—
	<b>83.9%*</b> <small>human solved</small>	82.6% <small>human solved</small>	80.4% <small>human solved</small>	—	—
<b>Agentic coding</b> Terminal-Bench 2.1	<b>88.0%*</b>	—	82.7%	83.4% <small>Codex CLI</small>	70.7% <small>Gemini CLI</small>
<b>Cybersecurity</b> ExploitBench (Cap%)	<b>78.0%*</b>	69.0%	40.0%	34.0%	—
<b>Health</b> HealthBench Professional	<b>66.0%*</b>	64.7%	56.9%	51.8%	—

**Methodology:** Reported scores are within a 1-3 percentage point difference for Claude Mythos 5 and Claude Fable 5. This table shows the higher score of the two. Starred (\*) benchmarks show a larger difference due to our blocking safeguards for cybersecurity and biology-related questions. For these benchmarks, Claude Fable 5 performs closer to Claude Opus 4.8 due to fallbacks. See the system card for details.

Source:

<https://www.anthropic.com/news/claude-fable-5-mythos-5>

# What Matters More Than the Model...*The Prompt*

*Prompting: the process of providing context, goals, constraints, and examples that help the model produce useful output*

- Mediocre Prompt:
  - "Tell me about Noonan syndrome."
- Better Prompt:
  - "Summarize Noonan syndrome for MFM counseling note that will go to the patient's ob provider (who knows me and works with me frequently). Include prenatal findings that are suggestive of this diagnosis, talk about recurrence risk, testing options, and postnatal outcomes. Keep it to 3 paragraphs. Use bullet points."

*In 2026, Prompt Engineering is becoming a critical clinical skill.*

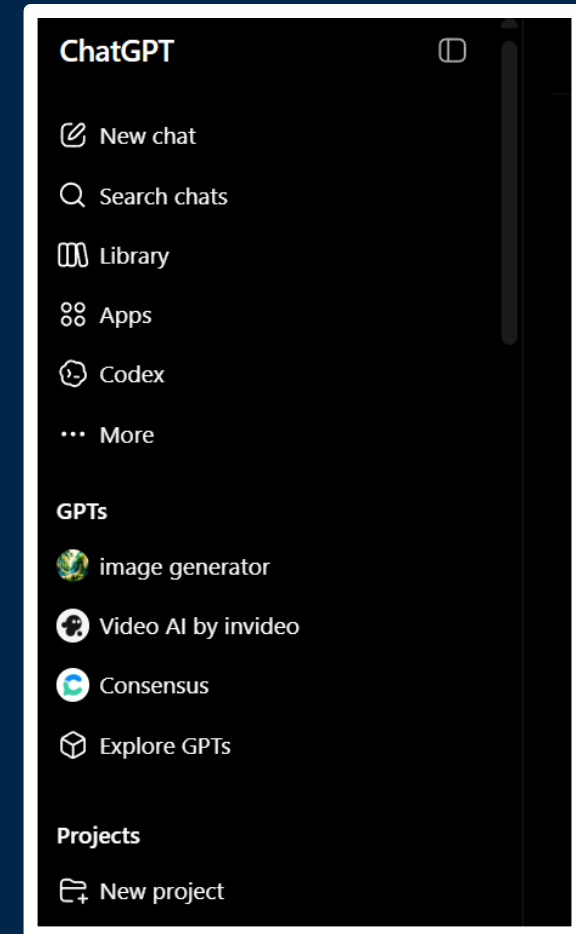
# Customizing Models

...Most clinicians have no idea this capability exists.

- Topics:
  - Projects
  - Persistent instructions
  - Dedicated workspace
  - Reusable references
- Knowledge Files
- Upload Examples:
  - Practice protocols
  - FMF guidelines
  - SMFM documents
  - Genetic counseling templates
  - Custom Instructions

Example:

"You are assisting a maternal-fetal medicine physician. Use concise language. Cite evidence when possible."



Feature	Purpose	Example for MFM/Genetics
Custom Instructions	Define role and style	"Respond as a board-certified MFM specialist using concise clinical language."
Project Folders	Maintain persistent context	MFM Practice, Genetics Consults, Research Projects
Knowledge Files	Provide trusted references	SMFM Consults, FMF protocols, ACMG guidelines
Markdown (.md) Files	Encode reusable instructions and workflows	FGR workup template, cfDNA counseling template
Reference Libraries	Reduce hallucinations	Practice policies, published guidelines, counseling scripts
Memory & Continuity	Preserve ongoing discussions	Longitudinal research or manuscript development

# Take Home Points

- There is no single best model. Use more than one model for important work and compare outputs.
- Match the model to the task.
- Prompt quality matters more than model selection.
- Use projects and knowledge files.
- Verify outputs. You and only you are responsible for patient care. Legally and ethically.
- AI is becoming a routine clinical and research enhancement capability.



*AI is not your replacement and not simply a tool. AI is like an exosuit. AI is an enhancement that can make you a more useful and effective human. AI cannot choose that for you. Only you can.*

# Thank You

This talk can be found at:

<https://thektproject.org/>

