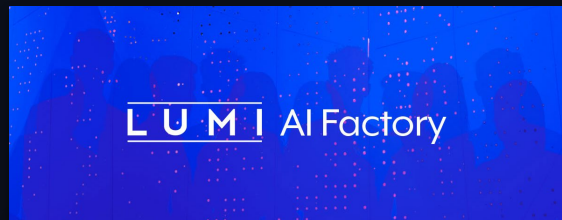


/root signals



Can We Control GenAI Without Measuring It?

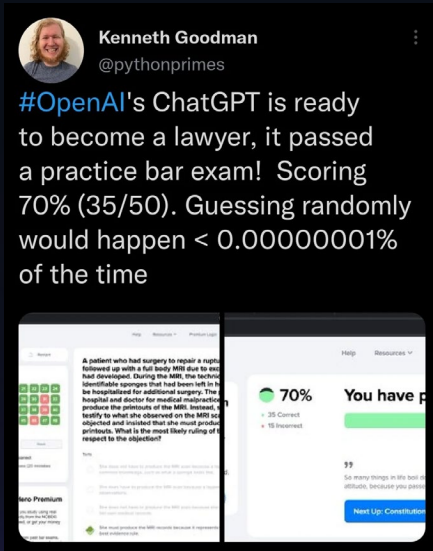
Oguzhan (Ouz) Gencoglu - Co-founder & Head of AI

2 April 2025

No!

Agenda

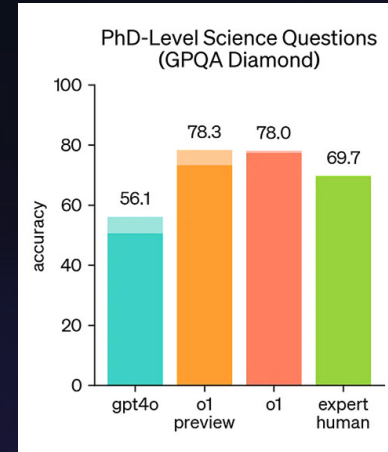
- 2 Questions
 - Why is that RoI of GenAI is almost 0%?
 - Why do I still have to write emails?
- *Root Judge* & LUMI
- Common Misconceptions About LLM Judges



AI Passes U.S. Medical Licensing Exam
 — Two papers show that large language models, including ChatGPT, can pass the USMLE

A.I. Chatbots Defeated Doctors at Diagnosing Illness

A small study found ChatGPT outdid human physicians when assessing medical case histories, even when those doctors were using a chatbot.



Why is that so?

Non-deterministic
systems

No self reflection

Semantic guidance

Failures

- Inaccuracy
- Hallucinations
- Biases
- Verbosity
- Diverging from business interests
- ...

→ Multi - dimensional Failures

Evals has always been
important!

But it was not hard!

Now it is !

SOLUTION:

Properly

tuned / optimized / calibrated

LLM-Judges!

 Root Judge

State-Of-The-Art Judge LLM For Evaluation & Hallucination Detection

Root Judge, a ground-breaking LLM that sets a new standard for reliable, customizable and locally-deployable evaluation models.

Try In Platform



Hugging Face

Couple of Misconceptions about LLM Judges

Misconception:

“Judges are there just to block unwanted responses”

REALITY:

- Many criteria are matters of degree, e.g. safety, relevance, and often even truthfulness – and can align with humans [1,2]
- Metric outputs [0...1] unlock optimization
- Guardrails are simply threshold-conditioned metrics
- Recent trend of inference-time search increases the need for evaluations for choosing candidates, with risks of hallucination increasing with increased turns

[1] <https://arxiv.org/pdf/2306.05685>

[2] <https://arxiv.org/html/2403.03230v1>

Misconception:

“When foundation models become stronger, no judges (or evals) are any longer necessary”

REALITY:

- As foundation foundation models get better, the maximum performance does indeed improve
- But smaller and more energy-efficient models will always be needed – even if there were an AGI
- Also – historically, as models improve, we have always moved the goalposts

Misconception:

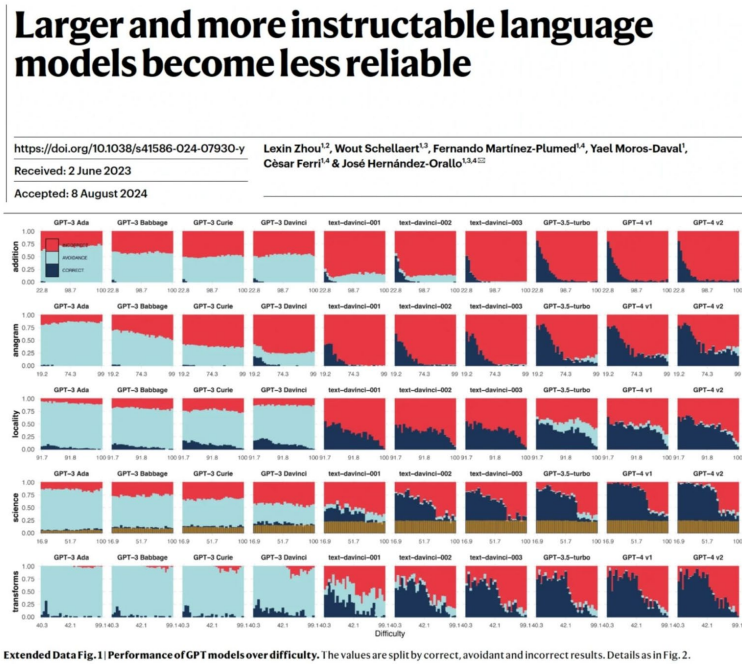
“Judges are only for pre-production (development)”

REALITY:

- LLMs can perform complex task before being flawless on very easy tasks
- There is no “reliable” zone one can identify
- This paradigm is new even to ML Experts

Misconception:

“Judges are only for development, not for production”





Root Judge & LUMI

Take Away

1. LLMs are unreliable by nature
2. Lack of reliability is the main blocker for large scale adoption
3. A measurement layer around GenAI automations is needed
4. But metrics we want to measure are complex → LLM

Judges

5. We need a systematic way to define, tune, and manage these Judges → Root Signals

/ thank you