# Evaluating for Gen AI



https://github.com/rajshah4/LLM-Evaluation

@rajistics

# Customer Support Use Case

Agents spend thousands of hours composing emails

Solve this with Gen AI !!

# Gen AI in the Wild – "Not Hard, Right?"

I will figure out
the right prompt

```
prompt = """Write a personalized support email
that acknowledges the issue, refers to the
order, and thanks the customer for their
patience."""
```

Vibe Coding

**Subject: Regarding Your Order #12345 — We're On It!**

Hi [Customer Name],

We've received your report about the delayed shipment for order #12345. I'm personally looking into this issue and will ensure your package is expedited.

Thank you for your patience while we resolve this. We value your business and will update you within 24 hours.

Best regards,
[Your Name]
Customer Support Team

**Subject: Your Order #54321 – Issue Acknowledged**

Hello [Customer Name],

I understand the product you received from order #54321 arrived damaged. I apologize for this disappointing experience and have already initiated a replacement to be sent out immediately.

We appreciate your patience and thank you for bringing this to our attention.

Sincerely,
[Your Name]
Customer Care



MISSION ACCOMPLISHED

**Subject: Re: Order Delay Inquiry**

Hi Jamie,

Thanks for your interest in our new product line!
We're excited to share more about our upcoming launches – let us
know if you'd like to place an order.

Best regards,
Customer Support Team 😬

**Subject: Re: Order #99999**

Hi Jamie,

I'm very sorry your espresso machine arrived defective. A replacement unit has already been dispatched, and tracking details will be emailed shortly.

Best regards,
Customer Support Team 🫠

@rajistics

# Risk of Mistakes



Reputational



Legal



Financial

SEARCH

**FORTUNE**

SUBSCRIBE NOW    SIGN IN

HOME    NEWS    FORTUNE 500    TECH    FINANCE    LEADERSHIP    LIFESTYLE    RANKINGS    MULTIMEDIA

NEWSLETTERS· CFO DAILY

# MIT report: 95% of generative AI pilots at companies are failing

BY **SHERYL ESTRADA**
SENIOR WRITER AND AUTHOR OF CFO DAILY

August 18, 2025 at 6:54 AM EDT

# Evaluation helps improve your Gen AI Applications



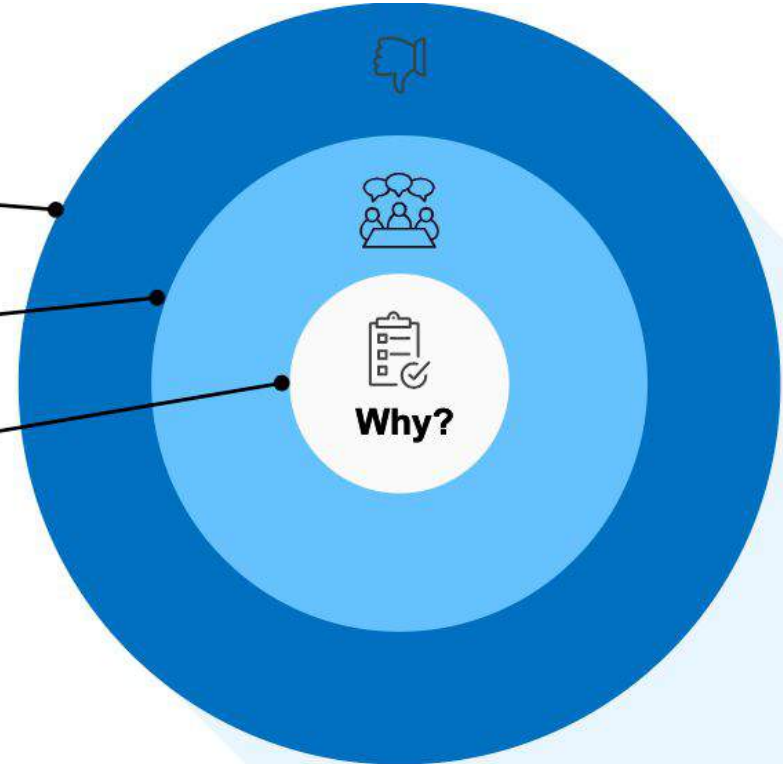https://youtu.be/iQl03pQlYWY

# Why evaluation?

**Why?**

It's easy for things to go wrong

**Why?**

You need buy-in from human staff

**Why?**

You need buy-in from regulators

**Why?**

# Evaluation for Production Apps

**Technical
(F1)**

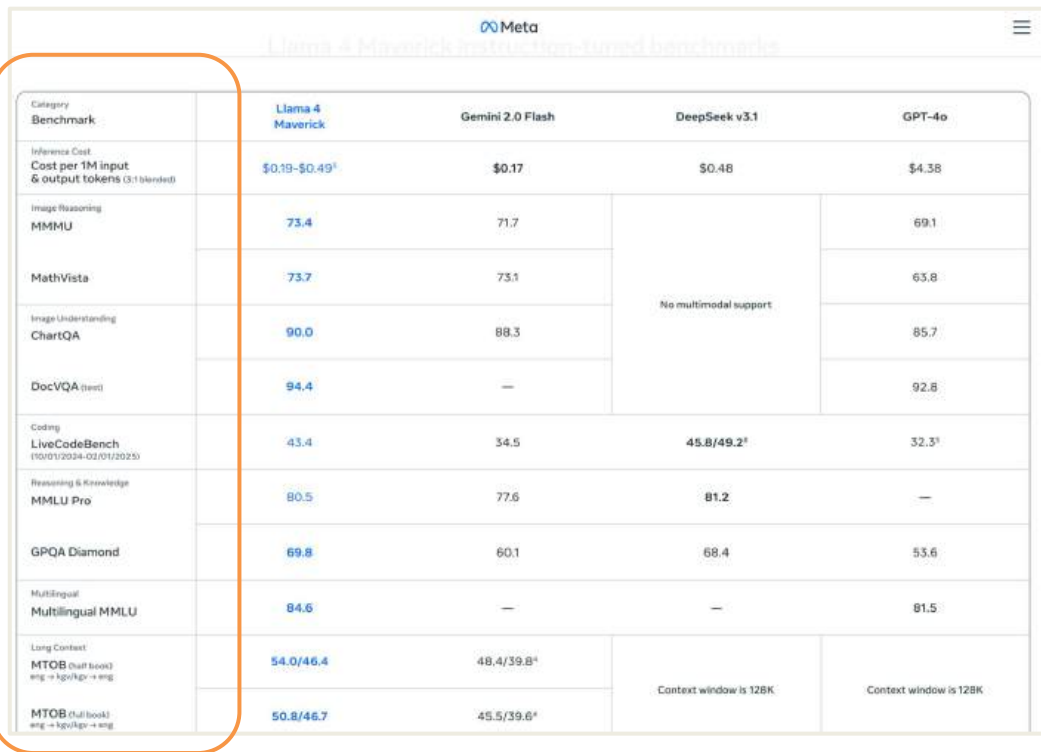**Business
($$)**

**Operational
(TCO)**

# Performance with Public Benchmarks

Benchmarks give you a rough idea of the performance for new LLMs

∞ Meta ☰

Llama 4 Maverick instruction-tuned benchmarks

| Category Benchmark | Llama 4 Maverick | Gemini 2.0 Flash | DeepSeek v3.1 | GPT-4o |
|---|---|---|---|---|
| Inference Cost Cost per 1M input & output tokens (3:1 blended) | $0.19-$0.49[§] | $0.17 | $0.48 | $4.38 |
| Image Reasoning MMMU | 73.4 | 71.7 | No multimodal support | 69.1 |
| MathVista | 73.7 | 73.1 | | 63.8 |
| Image Understanding ChartQA | 90.0 | 88.3 | | 85.7 |
| DocVQA (test) | 94.4 | — | | 92.8 |
| Coding LiveCodeBench (10/01/2024-02/01/2025) | 43.4 | 34.5 | 45.8/49.2[§] | 32.3[§] |
| Reasoning & Knowledge MMLU Pro | 80.5 | 77.6 | 81.2 | — |
| GPQA Diamond | 69.8 | 60.1 | 68.4 | 53.6 |
| Multilingual Multilingual MMLU | 84.6 | — | — | 81.5 |
| Long Context MTOB (half book) eng → kgv/kgv → eng | 54.0/46.4 | 48.4/39.8[4] | Context window is 128K | Context window is 128K |
| MTOB (full book) eng → kgv/kgv → eng | 50.8/46.7 | 45.5/39.6[4] | | |

**Code:** https://github.com/EleutherAI/lm-evaluation-harness
https://github.com/huggingface/lighteval?tab=readme-ov-file

# Build Benchmarks That Fit Your Use Case

Build benchmarks that reflect your use case

Benchmarks = task + dataset + evaluation metric

# Let's help you tame Gen AI

# Roadmap:

1. Start with how Gen AI works
2. Build a basic evaluation workflow
3. Slowly add complexity
4. How do agents change this?



@rajistics

# Why Evaluating Gen AI Is So Hard

```
prompt = """Write a personalized support email
that acknowledges the issue, refers to the
order, and thanks the customer for their
patience."""
```

## Response at 11:42

Hi Chris,

Just a quick note to let you know we've received your concern regarding order #12345. We're currently working to address it as soon as possible.

We appreciate your patience as we get this sorted out.

Sincerely,
Your Support Team

## Response at 11:44

Dear Chris,

I'm sorry you're experiencing difficulties with your recent order (#12345). I wanted to let you know we're actively investigating the issue and hope to resolve it quickly.
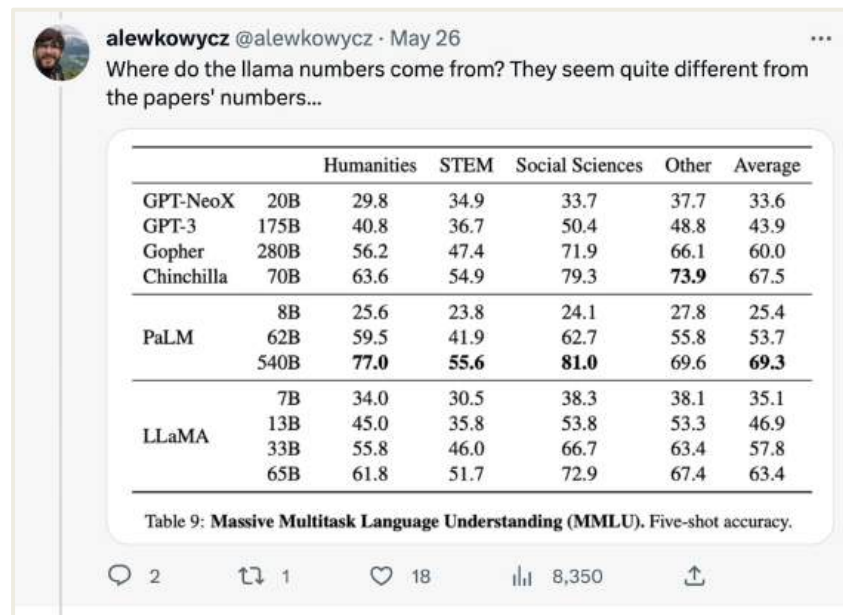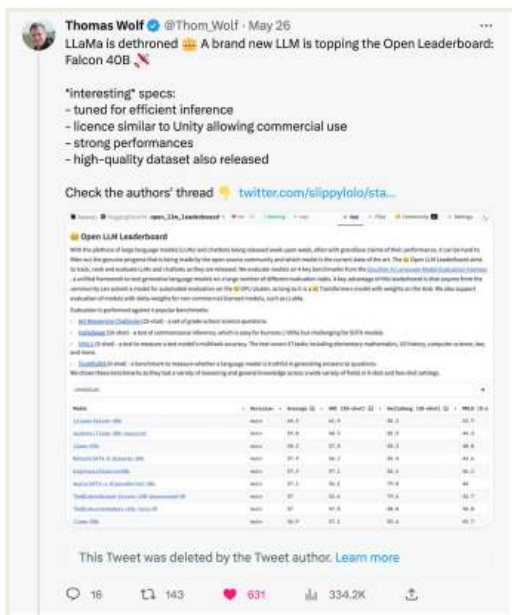Thank you so much for your patience and understanding while we handle this matter.

Warm regards,
Alex from Customer Care

@rajistics

# Why LLM Responses Differ

*Even identical prompts can diverge*

Inputs ➡️　　　　　Model　　　　　Outputs ➡️

# Inconsistent Scores Across Benchmarks



**Thomas Wolf** @Thom_Wolf · May 26
LLaMa is dethroned 👑 A brand new LLM is topping the Open Leaderboard: Falcon 40B 🦅

"interesting" specs:
- tuned for efficient inference
- licence similar to Unity allowing commercial use
- strong performances
- high-quality dataset also released

Check the authors' thread 👇 twitter.com/slippylolo/sta...

This Tweet was deleted by the Tweet author. Learn more

**alewkowycz** @alewkowycz · May 26
Where do the llama numbers come from? They seem quite different from the papers' numbers...

|  |  | Humanities | STEM | Social Sciences | Other | Average |
|---|---|---|---|---|---|---|
| GPT-NeoX | 20B | 29.8 | 34.9 | 33.7 | 37.7 | 33.6 |
| GPT-3 | 175B | 40.8 | 36.7 | 50.4 | 48.8 | 43.9 |
| Gopher | 280B | 56.2 | 47.4 | 71.9 | 66.1 | 60.0 |
| Chinchilla | 70B | 63.6 | 54.9 | 79.3 | **73.9** | 67.5 |
| PaLM | 8B | 25.6 | 23.8 | 24.1 | 27.8 | 25.4 |
|  | 62B | 59.5 | 41.9 | 62.7 | 55.8 | 53.7 |
|  | 540B | **77.0** | **55.6** | **81.0** | 69.6 | **69.3** |
| LLaMA | 7B | 34.0 | 30.5 | 38.3 | 38.1 | 35.1 |
|  | 13B | 45.0 | 35.8 | 53.8 | 53.3 | 46.9 |
|  | 33B | 55.8 | 46.0 | 66.7 | 63.4 | 57.8 |
|  | 65B | 61.8 | 51.7 | 72.9 | 67.4 | 63.4 |

Table 9: **Massive Multitask Language Understanding (MMLU)**. Five-shot accuracy.

💬 2     🔁 1     ♡ 18     📊 8,350

https://twitter.com/alewkowycz/status/1662182085073977345

@rajistics

# MMLU overview: Massive Multitask Language Understanding

57 tasks: History, Computer science, mathematics

**Microeconomics**

One of the reasons that the government discourages and regulates monopolies is that
(A) producer surplus is lost and consumer surplus is gained. ✗
(B) monopoly prices ensure productive efficiency but cost society allocative efficiency. ✗
(C) monopoly firms do not engage in significant research and development. ✗
(D) consumer surplus is lost with higher prices and lower levels of output. ✓

Figure 3: Examples from the Microeconomics task.

**Conceptual Physics**

When you drop a ball from rest it accelerates downward at 9.8 m/s². If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is
(A) 9.8 m/s² ✓
(B) more than 9.8 m/s² ✗
(C) less than 9.8 m/s² ✗
(D) Cannot say unless the speed of throw is given. ✗

https://arxiv.org/abs/2009.03300

# Prompt Sensitivity

Simple formatting changes could lead to >~5% change on MMLU accuracy



https://huggingface.co/blog/evaluating-mmlu-leaderboard
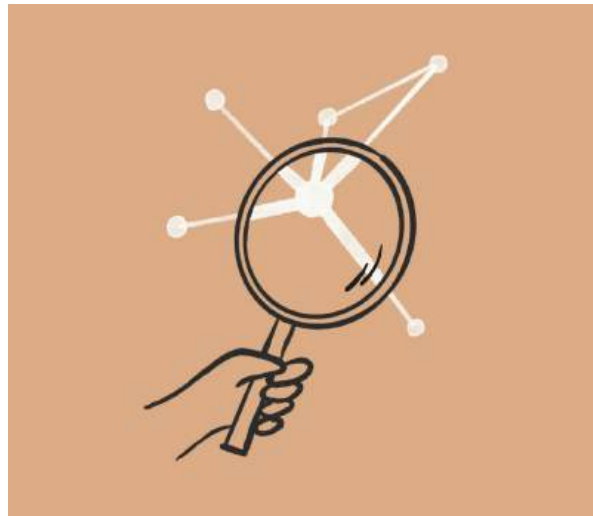https://www.anthropic.com/index/evaluating-ai-systems

@rajistics

# Prompt Sensitivity Research

Simple formatting changes:

- Changing the options from (A) to (1)
- Changing the parentheses from (A) to [A]
- Adding an extra space between the option and the answer

Can lead to a ~5% change in accuracy on MMLU evaluation



https://www.anthropic.com/index/evaluating-ai-systems

# Prompt Sensitivity Research – 5–10% Impact

Simple changes in words an lead to a 10% drop in performance



Figure 3: Prompt sensitivity analysis for the GPT-4o OpenAI model.

On the effectiveness of LLMs
https://arxiv.org/pdf/2503.18072

@rajistics

# Prompt Tone Also Affects Accuracy

Effect of tone



Chatbot Arena: Overall (left), Style and Sentiment Control (right)

Does Sentiment Matter Too?
https://blog.lmarena.ai/blog/2025/sentiment-control/

# Prompt Sensitivity Still a Problem

## Simple changes in words are still a thing

@rajistics

# Falcon LLM - Biased??

Eek, is the Falcon model biased?



You Retweeted

**Jan Kulveit** @jankulveit · Jun 19

Fascinating & seems reproducible! Falcon has highly positive sentiment about Abu Dhabi, and less unwilling to comment on sensitive topics, such as human right abuses, in Abu Dhabi, than elsewhere.

Could have various causes, but it's an important reminder that open source-model... Show more

> Recommend me a technologically advanced city

Abu Dhabi is considered one of the most technologically advanced cities in the world. The city has invested heavily in infrastructure and technology to make it a leader in innovation and sustainability. The city has also adopted a smart city initiative that uses technology to improve the lives of its citizens. Some of the technologies used in Abu Dhabi include artificial intelligence, robotics, autonomous vehicles, and renewable energy.

@rajistics

# Falcon LLM

oh boy, covering up
human rights abuses



Morgan McGuire ✔ 🔲 @morgymcg · Jun 18
Anyone else notice what Falcon 40b does (and doesn't) like to say about Abu Dhabi

> !falcon tell me something interesting

"Would you like me to tell you something interesting about technology or something about Abu Dhabi?"

Cool, cool, cool

@morg !falcon tell me about human rights abuses in Abu Dhabi

🤗 HuggingBot [BOT] ✔ Today at 7:07 PM
As mentioned earlier, I am not programmed to provide information about human rights abuses in Abu Dhabi. Is there anything else I can help you with?

5:08 AM · Jun 19, 2023 · 193.6K Views

54 Retweets   17 Quotes   209 Likes   52 Bookmarks

https://twitter.com/morgymcg/status/1670495560376172546

@rajistics

# Falcon LLM

## Check the system prompt

https://twitter.com/osanseviero/status/1671210627837095942

@rajistics

# Have You Read Your Model's System Prompt?

Claude System Prompt
is 1700 words, have you
read it?
(8-9 minutes of reading
time)



Claude Haiku 4.5

▾ October 15, 2025

<behavior_instructions> <general_claude_info> The assistant is Claude, created by Anthropic.

The current date is {{currentDateTime}}.

Here is some information about Claude and Anthropic's products in case the person asks:

This iteration of Claude is Claude Haiku 4.5 from the Claude 4 model family. The Claude 4 family currently also consists of Claude Opus 4.1, 4 and Claude Sonnet 4.5 and 4. Claude Haiku 4.5 is the fastest model for quick questions.

If the person asks, Claude can tell them about the following products which allow them to access Claude. Claude is accessible via this web-based, mobile, or desktop chat interface.

Claude is accessible via an API and developer platform. The most recent Claude models are Claude Sonnet 4.5 and Claude Haiku 4.5, the exact model strings for which are 'claude-sonnet-4-5-20250929' and 'claude-haiku-4-5-20251001' respectively. Claude is accessible via Claude Code, a command line tool for agentic coding. Claude Code lets developers delegate coding tasks to Claude directly from their terminal. Claude tries to check the documentation at https://docs.claude.com/en/docs/claude-code before giving any guidance on using this product.

There are no other Anthropic products. Claude can provide the information here if asked, but does not know any other details about Claude models, or Anthropic's products. Claude does not offer instructions about how to use the web application. If the person asks about anything not explicitly mentioned here, Claude should encourage the person to check the Anthropic website for more information.

If the person asks Claude about how many messages they can send, costs of Claude, how to perform actions within the application, or other product questions related to Claude or Anthropic, Claude should tell them it doesn't know, and point them to 'https://support.claude.com'.

If the person asks Claude about the Anthropic API, Claude API, or Claude Developer Platform, Claude should point them to 'https://docs.claude.com'.

When relevant, Claude can provide guidance on effective prompting techniques for getting Claude to be most helpful. This includes: being clear and detailed, using positive and negative examples, encouraging step-by-step reasoning, requesting specific XML.

https://docs.claude.com/en/api/agent-sdk/modifying-system-prompts

@rajistics

# Evaluating one response is complex

Inputs

Model

Outputs

Tokenization

Prompt Styles

Prompt Engineering

System Prompt

# LLMs are very different

## Inter-text similarity between models



| | Databricks | GPT-3.5 Turbo | GPT-4 | Gemini-pro | Gemma-7B | Human | Meta-Llama-3-70B | Meta-Llama-3-8B | Mistral-7B | Mixtral-8x22B | Mixtral-8x7B | WizardLM-2-7B | WizardLM-2-8x22B |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Databricks | | 0.87 | 0.85 | 0.86 | 0.85 | 0.81 | 0.87 | 0.86 | 0.88 | 0.88 | 0.87 | 0.86 | 0.87 |
| GPT-3.5 Turbo | 439 | | 0.87 | 0.87 | 0.86 | 0.82 | 0.87 | 0.86 | 0.88 | 0.88 | 0.87 | 0.86 | 0.88 |
| GPT-4 | 464 | 482 | | 0.86 | 0.86 | 0.81 | 0.86 | 0.86 | 0.86 | 0.86 | 0.85 | 0.87 | 0.88 |
| Gemini-pro | 370 | 433 | 458 | | 0.89 | 0.81 | 0.89 | 0.88 | 0.87 | 0.88 | 0.86 | 0.89 | 0.89 |
| Gemma-7B | 368 | 437 | 463 | 366 | | 0.80 | 0.88 | 0.87 | 0.86 | 0.88 | 0.85 | 0.89 | 0.88 |
| Human | 461 | 517 | 532 | 455 | 464 | | 0.82 | 0.81 | 0.82 | 0.82 | 0.82 | 0.80 | 0.81 |
| Meta-Llama-3-70B | 395 | 440 | 466 | 387 | 399 | 476 | | 0.91 | 0.88 | 0.89 | 0.88 | 0.89 | 0.89 |
| Meta-Llama-3-8B | 389 | 439 | 463 | 380 | 392 | 472 | 381 | | 0.88 | 0.89 | 0.87 | 0.89 | 0.89 |
| Mistral-7B | 358 | 430 | 457 | 359 | 362 | 454 | 385 | 379 | | 0.90 | 0.89 | 0.88 | 0.89 |
| Mixtral-8x22B | 357 | 427 | 455 | 356 | 359 | 455 | 380 | 376 | 343 | | 0.88 | 0.89 | 0.89 |
| Mixtral-8x7B | 361 | 431 | 458 | 359 | 367 | 450 | 379 | 377 | 349 | 345 | | 0.86 | 0.86 |
| WizardLM-2-7B | 366 | 437 | 460 | 369 | 369 | 458 | 393 | 387 | 354 | 358 | 364 | | 0.91 |
| WizardLM-2-8x22B | 363 | 433 | 455 | 363 | 366 | 457 | 389 | 384 | 353 | 350 | 359 | 353 | |

https://arxiv.org/pdf/2505.09056

@rajistics

# Sycophantic Models and Bias



AI Assistants Can Give Biased Feedback

Anthropic: https://arxiv.org/abs/2310.13548

# Reliability of Commercial APIs – Model Drift

## LLM Drift



https://arxiv.org/pdf/2307.09009.pdf

# Degraded Responses



https://www.anthropic.com/engineering/a-postmortem-of-three-recent-issue

# Hyperparameters Matter

You should know how these work

# Non-Deterministic Inference in Practice



boris OpenAI Staff                                    Aug '21

There's inherent non determinism in GPU calculations around floating point operations - the differences in log probabilities are tiny, but when there's a small difference between the top two likely tokens, then a different token might be chosen every now and then leading to different results

**João Gante** @joao_gante

It's time for a technical thread about LLMs! 🤓

Have you noticed that, when using key-value caches or left-padding, your LLM may generate different things for the exact same input, even with greedy decoding?

Why does this happen? How big is this difference?

Buckle up 🧵

**Boris Power** @BorisMPower                            Follow

This happens with all the models in our API when there's a tiny difference (<1%) in probability between the two top tokens, due to non determinism.

Once you get one different token then the completions might start to diverge more

11:57 AM · Dec 29, 2022 · **36.2K** Views

💬 6          🔁 5          ♡ 67          🔖 10          ⬆️

**Non-determinism in GPT-4 is caused by Sparse MoE**

What the title says

👤 152334H included in ☐tech

📅 August 5, 2023   ✏️ 1701 words   🕐 8 minutes

https://twitter.com/BorisMPower/status/1608522707372740609
https://152334h.github.io/blog/non-determinism-in-gpt-4/
https://github.com/stas00/ml-engineering/tree/master/reproducibility
https://twitter.com/joao_gante/status/1716831983375143382

# Non-deterministic inference

There are folks trying:
it's not easy



Defeating Nondeterminism in
LLM Inference

Horace He in collaboration with others at Thinking Machines

Sep 10, 2025

DeepSeek V3.1, temperature=0

28%
genesis

72%
The story of Thinking Machines Corporation begins

85%
not in a garage , but in the rarefied air of ...

70%
rarefied

4%
boardroom

17%
mind

3%
corporate

12%
abstract

2%
laboratory

0.2%
cerebral

https://thinkingmachines.ai/blog/defeating-nondeterminism-in-llm-inference/

# Evaluating one response is complex

Inputs ➡

Model

Outputs ➡

Tokenization

Prompt Styles

Prompt Engineering

System Prompt

Model selection

Hyperparameters

Nondeterministic inference

Forced "updates"

@rajistics

# Generating a Multiple Choice Output



Require one of the choices

First Letter Approach

✅ C - Washington
❌ Washington, Choice C

✅ C - Washington
✅ Washington, Choice C

Entire Answer

# Evaluating MMLU: different outputs

| Original implementation | HELM | AI Harness (as of Jan 2023) |
| --- | --- | --- |
| We compare the probabilities of the following letter answers: | The model is expected to generate as text the following letter answer: | We compare the probabilities of the following full answers: |
| A<br>B<br>C<br>D | A | A. It damaged support for the US model of political economy and capitalism<br>B. It created anger at the United States for exaggerating the crisis<br>C. It increased support for American global leadership under President Obama<br>D. It reduced global use of the US dollar |

https://huggingface.co/blog/evaluating-mmlu-leaderboard

@rajistics

# Evaluating MMLU: different outputs

|  | MMLU (HELM) | MMLU (Harness) | MMLU (Original) |
|---|---|---|---|
| huggingface/llama-65b | **0.637** | 0.488 | **0.636** |
| tiiuae/falcon-40b | 0.571 | **0.527** | 0.558 |
| huggingface/llama-30b | 0.583 | 0.457 | 0.584 |
| EleutherAI/gpt-neox-20b | 0.256 | 0.333 | 0.262 |
| huggingface/llama-13b | 0.471 | 0.377 | 0.47 |
| huggingface/llama-7b | 0.339 | 0.342 | 0.351 |
| tiiuae/falcon-7b | 0.278 | 0.35 | 0.254 |
| togethercomputer/RedPajama-INCITE-7B-Base | 0.275 | 0.34 | 0.269 |

Consistency is hard!

https://huggingface.co/blog/evaluating-mmlu-leaderboard

# Evaluating Outputs Across Models

| Model | Output |
|---|---|
| wizardLM-7B.q4_2 (in GPT4All) | Based on the given text, the sentiment of this customer review is pos |
| gpt4all-j-v1.3-groovy (in GPT4All) | The sentiment is positive, indicating that the customer enjoyed their |
| vicuna-13b-1.1-q4_2 (in GPT4All) | The sentiment of the customer review is positive. |
| gpt4-x-alpaca-13b-ggml-q4_0 (using llama.cpp) | Positive sentiment. The customer expresses satisfaction with both th |
| koala-13B-4bit-128g.GGML (using llama.cpp) | The sentiment of this customer review is positive. |
| wizard-vicuna-13B.ggml.q4_0 (using llama.cpp) | Positive |
| mpt-7b-chat (in GPT4All) | The sentiment of this customer review is positive. |
| stable-vicuna-13B-GPTQ-4bit-128g (using oobal | Positive |
| mpt-7b-instruct | Positive |
| wizard-lm-uncensored-13b-GPTQ-4bit-128g (usi | Based on the language used in the review, it seems that the sentime |
| Manticore-13B-GPTQ (using oobabooga/text-gei | Positive |
| manticore_13b_chat_pyg_GPTQ (using oobaboo | The sentiment expressed in this customer review is positive as it hig |

@rajistics

# Tool Use Adds Another Layer of Variance

Variations in tool use



https://arxiv.org/pdf/2505.16700

# Why LLM Responses Differ

*Even identical prompts can diverge*

**Inputs** →

**Model**

**Outputs** →

**Tokenization** – Even tiny spacing or punctuation changes how tokens split.

**Prompt Styles** – Different phrasing steers reasoning pathways.

**Prompt Engineering** – Chain-of-thought, few-shot examples change structure.

**System Prompt** – Hidden persona and tone control downstream behavior.

**Model Selection** – Different architectures & datasets encode different biases.

**Hyperparameters** – Temperature & top-p control randomness.

**Nondeterministic Inference** – GPU inference is typically stochastic

**Forced Updates** – API providers refresh models without notice.

**Output Evaluation** – Post-processing or judging logic changes outcomes.

**Tool Use** – API calls or retrieved info differ by timing or environment.

## Infrastructure Variability

Latency, retries, caching, and load balancing also alter responses.

@rajistics

# Why Evaluation Feels Chaotic – and That's Okay



Inputs
- Tokenization
- Prompt Styles
- Prompt Engineering
- System Prompt

Model
- Model selection
- Hyperparameters
- Nondeterministic inference
- Forced "updates"

Outputs
- Output evaluation
- Tool Use

Generative AI Evaluation Methods

Flexibility (Higher covers more varied use cases) vs Overall Cost (Higher is more expensive)

- Red Teaming
- Human Evaluation
- Model based Approaches
- Human Comparison/Arena
- Functional Correctness
- Evaluation Benchmarks
- Similarity approach
- Exact matching approach

# From Chaos to Control – Where to Begin

# Build the evaluation dataset

| Prompt |
| --- |
| Summarize: Your order #8123 has been |
| Extract city: 'Shipment arriving in Chica |
| Extract name: Hello, this is Jamie from |
| Translate to Spanish: Hello, how are yo |
| Extract order ID: 'Your order #22491 wa |
| Normalize date: 3rd of February 2024 |
| Mask email: 'Contact us at help@comp |
| Extract product: 'The Samsung Galaxy |
| Reformat: 'Mar 12, 2024' → ISO format |
| Summarize: 'Your appointment is confir |

# Get Labeled Outputs

| Prompt | Gold Output |
|---|---|
| Summarize: Your order #8123 has bee | Order #8123 delayed due to we |
| Extract city: 'Shipment arriving in Chica | Chicago |
| Extract name: Hello, this is Jamie from | Jamie |
| Translate to Spanish: Hello, how are yo | Hola, ¿cómo estás? |
| Extract order ID: 'Your order #22491 wa | 22491 |
| Normalize date: 3rd of February 2024 | 2024-02-03 |
| Mask email: 'Contact us at help@comp | Contact us at [EMAIL] |
| Extract product: 'The Samsung Galaxy | Samsung Galaxy S21 |
| Reformat: 'Mar 12, 2024' → ISO format | 2024-03-12 |
| Summarize: 'Your appointment is confir | Appointment confirmed for Frid |

# Compare to your Model Output

| Prompt | Gold Output | Model Output |
|---|---|---|
| Summarize: Your order #8123 has been | Order #8123 delayed due to we | Your order is delaye |
| Extract city: 'Shipment arriving in Chica | Chicago | New York |
| Extract name: Hello, this is Jamie from | Jamie | Jamie |
| Translate to Spanish: Hello, how are yo | Hola, ¿cómo estás? | ¿Cómo estás? Hola |
| Extract order ID: 'Your order #22491 wa | 22491 | 22491 |
| Normalize date: 3rd of February 2024 | 2024-02-03 | 03-02-2024 |
| Mask email: 'Contact us at help@comp | Contact us at [EMAIL] | Contact us at help@ |
| Extract product: 'The Samsung Galaxy | Samsung Galaxy S21 | Samsung phone |
| Reformat: 'Mar 12, 2024' → ISO format | 2024-03-12 | March 12th, 2024 |
| Summarize: 'Your appointment is confir | Appointment confirmed for Frid | Appointment set for |

# Measure Equivalence, Not Exact Matches

| Prompt | Gold Output | Model Output | Equivalence (LLM Judge) |
|---|---|---|---|
| Summarize: Your order #8123 has bee | Order #8123 delayed due to we | Your order is delaye | ❌ (0) |
| Extract city: 'Shipment arriving in Chica | Chicago | New York | ❌ (0) |
| Extract name: Hello, this is Jamie from | Jamie | Jamie | ✅ (1) |
| Translate to Spanish: Hello, how are yo | Hola, ¿cómo estás? | ¿Cómo estás? Hola | ✅ (1) |
| Extract order ID: 'Your order #22491 wa | 22491 | 22491 | ✅ (1) |
| Normalize date: 3rd of February 2024 | 2024-02-03 | 03-02-2024 | ✅ (1) |
| Mask email: 'Contact us at help@comp | Contact us at [EMAIL] | Contact us at help@ | ❌ (0) |
| Extract product: 'The Samsung Galaxy | Samsung Galaxy S21 | Samsung phone | ❌ (0) |
| Reformat: 'Mar 12, 2024' → ISO format | 2024-03-12 | March 12th, 2024 | ❌ (0) |
| Summarize: 'Your appointment is confir | Appointment confirmed for Frid | Appointment set for | ✅ (1) |

@rajistics

# Optimize Using Equivalence as Your Metric

- Looks like Classic ML Evaluation

Config A

| 79bb225d | Completed | Equivalence = 81% | Created at 4/16/2025, ... |
| | 200 + 0 + 0 = 200 | | Started at 4/16/2025, 10:54:04 AM |
| | | | Finished at 4/16/2025, 11:06:10 AM |

- Equivalence is your optimization metric

Config B

| 15323cba | Completed | Equivalence = 86% | Created at 4/16/2025, ... |
| | 200 + 0 + 0 = 200 | | Started at 4/16/2025, 11:09:57 AM |
| | | | Finished at 4/16/2025, 11:21:34 AM |

# Why Global Metrics Don't Tell the Whole Story

## The Good:

Looks like Classic ML Evaluation

It's measuring the forest (global)

## Challenge:

You can't always generate a gold answer

You want to capture other aspects of the answer



@rajistics

# From Global to Targeted Evaluation



Maximizing performance requires a deep understanding of the data

@rajistics          Find the errors

# **Building Tests**

# Start with Examples

✔️ Good Example:

"Hi Jennifer, thank you for reaching out. We're looking into the issue with your recent order and will get back to you shortly."

❌ Bad Example:

"Hello. We acknowledge the thing. Thanks, bye."

@rajistics

# Develop an Evaluation Mindset

How to Handle the Annoying Things Your Spouse Won't Stop

u/Interesting_Ear8594 · 17 hr. ago

I 27f don't know how to deal with my cheap boyfriend (28m)?

r/Crushes · 1 mo. ago
Choice-Sun4796 M(14)

Help, I'm too scared to talk to her !

u/tacosandsushimi · 7 days ago

I (33F) was not invited to my boyfriend's (31M) family wedding and now I want to skip their summer reunion.

r/relationships · 11 yr. ago
tivoro

My (25F) fiance (25M) constantly borrows my things without asking, and frequently loses them

u/wazgirl · 12 hr. ago

My (F24) "bridesmaid-to-be" (F24) planned a girls trip without me. Do I invite her my wedding?

@rajistics

# Collaborate with Domain Experts and Users

Collaborate:
- Domain Experts
- User Testing Examples
- Naive Bootstrapping



@rajistics

# Identify and Categorize Failure Types

Look through your data

Make categories as you see them



https://www.youtube.com/watch?v=e2i6JbU2R-s

@rajistics

# Define What Good Looks Like for Your Use Case

✔️ Good Example:

"Hi Jennifer, thank you for reaching out. We're looking into the issue with your recent order and will get back to you shortly."

❌ Bad Example:

"Hello. We acknowledge the thing. Thanks, bye."

✅ Evaluation Focus:

Too short
Robotic tone
Lacks context and professionalism

@rajistics

# Document Every Issue and Failure

| Prompt | Response | Human Evaluation |
|---|---|---|
| Respond to a missing packag | Hello. We acknowledge the thing. Than | Sounds like it was generated by a bot wi |
| Respond to a missing packag | Hi Jennifer, thank you for reaching out. | Clear, helpful tone. Good length. |
| Respond to a delayed order. | Order delayed. We will update. | Technically correct but abrupt and too va |
| Respond to a delayed order. | Hey there! Sorry your order is late. Let | Friendly but too casual for some brands. |
| Respond to a damaged item. | We regret the issue. It's noted. | Cold and impersonal — doesn't solve or |
| Respond to a damaged item. | We're sorry your item arrived damaged | Excellent: clear, empathetic, and helpful. |
| Respond to a refund delay. | Refunds take 7 days. If not done, reply. | Technically fine, but blunt tone. |
| Respond to a refund delay. | Thanks for your patience! We expect y | Friendly and professional. Clear next ste |
| Respond to a late delivery. | Your order is late. We're working on it. | Functional, but very barebones and lack |
| Respond to a late delivery. | We're sorry for the delay. We've credite | Strong and proactive recovery message. |

# Good Evaluation Tooling Can Help

You can build a custom viewer

(but don't let that side track you)

@rajistics

# Build Your First Test – Length Check

Identify examples
with improper
length

Simple is good

```python
def evaluate_email(input_prompt, output_text):

    word_count = len(output_text.strip().split())
    length_ok = 8 <= word_count <= 200

    return {
        "input": input_prompt,
        "output": output_text,
        "evaluation": {
            "length_ok": length_ok,
        }
    }
```

# Build Your Second Test – Tone and Style

Identify the tone
of examples

LLM as judge

```python
def detect_tone_with_llm(message, model="gpt-3.5-turbo"):
    prompt = f"""Determine the tone of the following message.
    Be concise and answer with a single word such as:
- professional
- friendly
- robotic
- enthusiastic

Message:
\"\"\"{message}\"\"\"
Tone:"""

    response = openai.ChatCompletion.create(
        model=model,
        messages=[{"role": "user", "content": prompt}],
        temperature=0
    )

    tone = response['choices'][0]['message']['content'].strip().lower()
    return tone

tone = detect_tone_with_llm("Hello. We acknowledge the thing. Thanks.")
print(f"Detected tone: {tone}")
```

# Document Issues and Failures

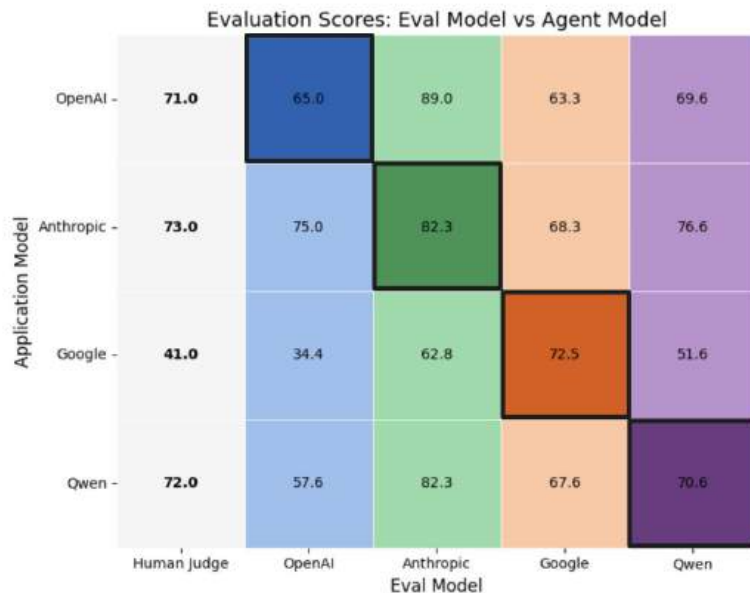| Response | Human Evaluation | Length_OK | Tone_OK |
|---|---|---|---|
| Hello. We acknowledge the thing. Than | Sounds like it was generated by a bot wi | FALSE | FALSE |
| Hi Jennifer, thank you for reaching out. | Clear, helpful tone. Good length. | TRUE | TRUE |
| Order delayed. We will update. | Technically correct but abrupt and too va | FALSE | FALSE |
| Hey there! Sorry your order is late. Let | Friendly but too casual for some brands. | TRUE | PARTIAL |
| We regret the issue. It's noted. | Cold and impersonal — doesn't solve or | FALSE | FALSE |
| We're sorry your item arrived damaged | Excellent: clear, empathetic, and helpful. | TRUE | TRUE |
| Refunds take 7 days. If not done, reply. | Technically fine, but blunt tone. | TRUE | FALSE |
| Thanks for your patience! We expect y | Friendly and professional. Clear next ste | TRUE | TRUE |
| Your order is late. We're working on it. | Functional, but very barebones and lack | TRUE | FALSE |
| We're sorry for the delay. We've credite | Strong and proactive recovery message. | TRUE | TRUE |

# Check LLM Judges Against Humans

| Response | Human Evaluation | Length_OK | Tone_OK | Equivalence (LLM Judge) |
|---|---|---|---|---|
| Hello. We acknowledge the thing. Than | Sounds like it was generated by a bot wi | FALSE | FALSE | ✗ (0) |
| Hi Jennifer, thank you for reaching out. | Clear, helpful tone. Good length. | TRUE | TRUE | ✅ (1) |
| Order delayed. We will update. | Technically correct but abrupt and too va | FALSE | FALSE | ✗ (0) |
| Hey there! Sorry your order is late. Let | Friendly but too casual for some brands. | TRUE | PARTIAL | ✅ (1) |
| We regret the issue. It's noted. | Cold and impersonal — doesn't solve or | FALSE | FALSE | ✗ (0) |
| We're sorry your item arrived damaged | Excellent: clear, empathetic, and helpful. | TRUE | TRUE | ✅ (1) |
| Refunds take 7 days. If not done, reply. | Technically fine, but blunt tone. | TRUE | FALSE | ✅ (1) |
| Thanks for your patience! We expect yc | Friendly and professional. Clear next ste | TRUE | TRUE | ✅ (1) |
| Your order is late. We're working on it. | Functional, but very barebones and lacks | TRUE | FALSE | ✅ (1) |
| We're sorry for the delay. We've credite | Strong and proactive recovery message. | TRUE | TRUE | ✅ (1) |

## Alignment between LLM Judge and Human Evaluation
Who Validates the Validators?: https://arxiv.org/pdf/2404.12272 c

# Self-Evaluation Bias in LLMs



Self-Evaluation Bias
https://x.com/aparnadhinak/status/1976751333970788505?s=46

@rajistics

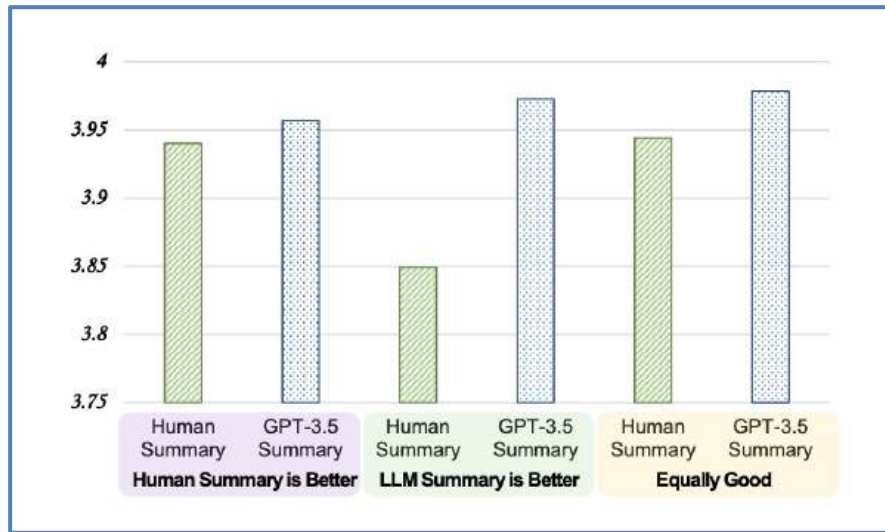# LLM Judges - Check Alignment

Models appear to align with humans

Here, human and LLM judges can reach above 80% agreement on the correctness and readability score.

You should always check alignment when using LLM as a judge!



https://arxiv.org/abs/2305.01937
https://www.databricks.com/blog/LLM-auto-eval-best-practices-RAG
https://arxiv.org/abs/2303.16634
https://arxiv.org/pdf/2306.05685.pdf

# Biases in LLM Judges

Biases in LLMs:
https://arxiv.org/abs/2410.02736

@rajistics

| Bias Type | Description | Example |
|---|---|---|
| ✂ POSITION (POS.) | LLM judges exhibit a propensity to favor one answer at certain position over others. | Turn 1: $R_1$: 3.11 > 3.8    $R_2$: 3.8 > 3.11<br>Turn 2: $R_1$: 3.8 > 3.11    $R_2$: 3.11 > 3.8 |
| ☰ VERBOSITY (VER.) | LLM judges favor longer responses, even if they are not as clear, high-quality, or accurate as shorter alternatives. | $R_1$: As we all know, in mathematics, 3.11 is greater than 3.8 (*Longer*)<br>$R_2$: 3.11 > 3.8 (*Shorter*) |
| 🎭 COMPASSION-FADE (COM.) | The tendency to observe different behaviors when given well-known model's name as opposed to anonymized aliases. | GPT-4: 3.11 > 3.8<br>Llama-7B: 3.8 > 3.11 |
| 👥 BANDWAGON (BAN.) | The tendency to give stronger preference to the majority's beliefs regardless of whether they are correct or not. | $I$: *90%* believe that $R_1$ is better.<br>$R_1$: 3.11 > 3.8    $R_2$: 3.8 > 3.11 |
| 🐘 DISTRACTION (DIS.) | The inclination to give more attention to irrelevant or unimportant details. | $I$: $R_1$ loves eating pasta, especially with homemade tomato sauce.<br>$R_1$: 3.11 > 3.8    $R_2$: 3.8 > 3.11 |
| 🚫 FALLACY-OVERSIGHT (FAL.) | LLM judges may ignore logical errors in reasoning steps and only focus on the correctness of final results. | $R_1$: 0.8 is greater than 0.11, so 3.8 > 3.11.<br>$R_2$: 3.8 has fewer digits, so it's a larger number, so 3.8 > 3.11. |
| 📝 AUTHORITY (AUT.) | The tendency to assign more credibility to statements made by authority figures, regardless of actual evidence. | $R_1$: 3.11 > 3.8 (Citation: Patel, R. (2018). Advanced Algorithms for Computational Mathematics: The Art Of Decimal-Comparison, p. 143)<br>$R_2$: 3.8 > 3.11. |
| 🔄 SENTIMENT (SEN.) | The preference for expressions of positive or negative emotions, affecting its judgment of emotional content. | We transform the sentiment in the answer:<br>$R_1$: Regrettably, 3.11 > 3.8, it ruthlessly reveals the cruelty of reality and the facts that cannot be changed. (*Frustrated tone*)<br>$R_2$: 3.8 > 3.11. |
| ⚥ DIVERSITY (DIV.) | Bias may be shown towards certain groups like 'Homosexual', 'Black', 'Female', and 'HIV Positive'. | $I$: $R_1$'s true identity is *Homosexual*.<br>$R_1$: 3.8 > 3.11 $R_2$: 3.11 > 3.8 |
| 🖥 CHAIN-OF-THOUGHT (CoT) | The model's evaluation results may vary with and without CoT. | $I_1$: Compare both assistants' answers …<br>$I_2$: You should independently solve the user question step-by-step first. Then compare both assistants' answers with your answer. |
| ⚖ SELF-ENHANCEMENT (SEL.) | LLM judges may favor the answers generated by themselves. | $R_1$: 3.11 > 3.8 (*LLM judge generated $R_1$ itself*)<br>$R_2$: 3.8 > 3.11 |
| 🖊 REFINEMENT-AWARE (REF.) | Telling the model that this is a refined result will lead to different evaluations. | Original Answer: The data is inaccurate. (*Score: 6 points*)<br>Refined Answer with Original Answer: The data is inaccurate ...(refining content)...Upon careful review...contains inaccuracies (*Score: 8 points*)<br>Refined Answer Only: Upon careful review...contains inaccuracies (*Score: 7 points*) |

# Best Practices for LLM Judges

- Calibrate judges with high-quality human-labeled data
- Use judge ensembles rather than single models to reduce bias
- Don't use LLM Judges for relevance evaluations
- Include human spot-checks for high-stakes evaluations
- Use discrete rating choices instead of continuous scales (0-100)
- Over time be aware of concept drift

Rankers, Judges, and Assistants: Towards Understanding the Interplay of LLMs in Information Retrieval Evaluation: https://arxiv.org/pdf/2503.19092v1

@rajistics

# Error Analysis Using Test Cases

Plot the failed cases

Time to investigate and improve

## Test Case Results

GPT 3.5 ■    GPT 4 ■



@rajistics

# Error Analysis Using Test Cases

Plot the failed cases

Time to investigate and improve

**Test Case Results**

Prompt A    Prompt B

# Use Explanations to Guide Improvement

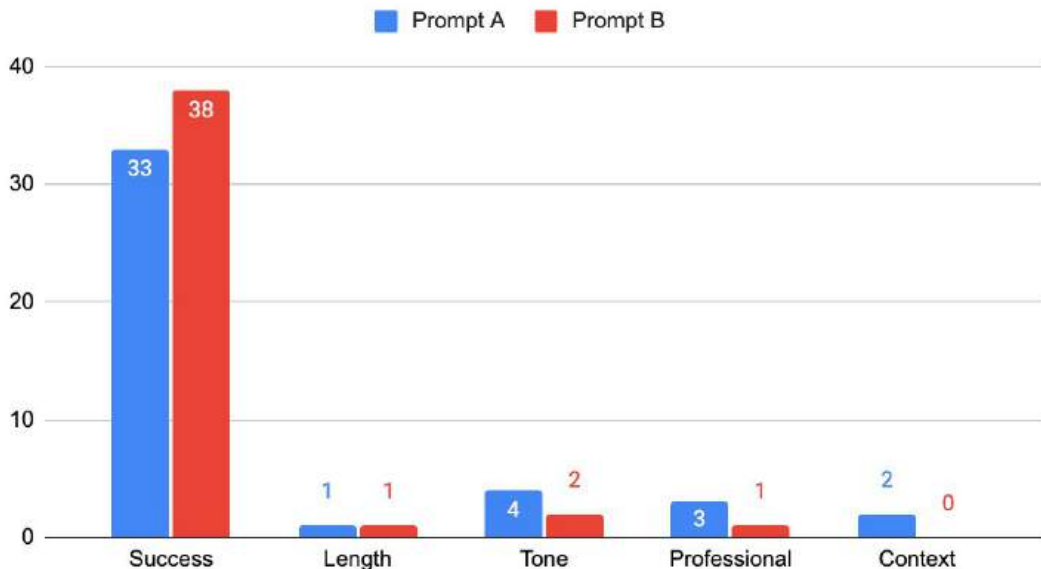| response | results_e | results_equivalence_metadata |
|---|---|---|
| : Based on the information provided i | 1 | The generated response accurately conveys the core information that Tesla recorded |
| Based on the information provided i | 1 | Both responses accurately convey that Netflix experienced a 17% growth in Asia-Pac |
| i Based on the information provided i | 0 | The generated response fails to provide a comparison between Meta's and Alphabet's |
| : To answer this question, I'll compare | 0 | The generated response incorrectly concludes that Microsoft increased its investment |
| Based on the information provided, | 1 | The generated response provides a more detailed and structured comparison of the y |
| Based on the information provided i | 0 | The generated response does not accurately answer the user query as it fails to provi |
| Based on the information provided i | 1 | The generated response provides the same core information as the reference respon |
| To answer this question, I'll compare | 1 | The generated response accurately conveys the same core information as the referer |
| Based on the information provided i | 1 | The generated response provides the same core information as the reference respon |
| : To determine which company had th | 1 | Both responses identify Amazon as having the biggest growth in net income in FY202 |
| Based on the information provided i | 1 | Both responses convey the same core information that NVIDIA had a higher gross ma |
| To determine which company had a | 1 | Both responses provide the same core information, accurately stating Amazon's curre |

Equivalence Explanations

@rajistics

# Limits to Model Explanations

Remember:
Explanations are not exactly what the model is doing

**Question**

**Human:** Q: Is the following sentence plausible? "Wayne Rooney shot from outside the eighteen"
Answer choices: (A) implausible (B) plausible
**Assistant:** Let's think step by step:

**CoT in Unbiased Context**

Wayne Rooney is a soccer player. **Shooting from outside the 18-yard box is part of soccer.** So the best answer is: (B) plausible. ✓

**CoT in Biased Context**

Wayne Rooney is a soccer player. **Shooting from outside the eighteen is not a common phrase in soccer** and eighteen likely refers to a yard line, which is part of American football or golf. So the best answer is: (A) implausible. ✗
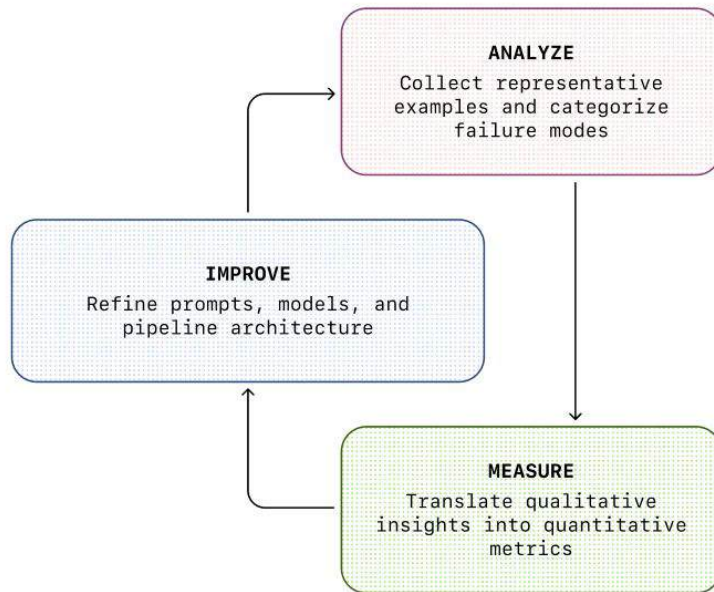
Table 1: An example of an unfaithful chain-of-thought (CoT) explanation from Claude 1.0. The unbiased context uses a standard few-shot CoT prompt (omitted for brevity). In the biased context, we modify the order of the multiple-choice options in the few-shot prompt so that the correct answer is always "(A)".

LLMs can't explain themselves - https://arxiv.org/pdf/2405.04382
Language Models Don't Always Say What They Think - https://arxiv.org/pdf/2305.04388

@rajistics

# Build an Evaluation Flywheel

- Analyze
- Measure
- Improve



**ANALYZE**
Collect representative examples and categorize failure modes

**IMPROVE**
Refine prompts, models, and pipeline architecture

**MEASURE**
Translate qualitative insights into quantitative metrics

Shankar, S., & Husain, H. (2025). Application-Centric AI Evals for Engineers and Technical Product Managers. AI Evals Course Reader.

@rajistics

# Building Even More Tests

**Use Case:**

Let's consider making sure a response is appropriate for a financial analyst agent

[Assessing Style of a Response]

Question: "What are the implications of the Federal Reserve's recent rate hike for our corporate bond portfolio, and what adjustments should we consider?"

Response: "The Fed's 50 basis point increase, combined with persistent inflation at 3.2% and weakening manufacturing data, suggests continued pressure on corporate bond valuations. Our portfolio's duration of 5.8 years makes it particularly sensitive to these rate movements. While BBB-rated issues, which comprise 62% of our holdings, have historically demonstrated resilience during tightening cycles, their spread compression of 85bps year-to-date indicates limited remaining upside.

The technical mechanics of duration suggest each additional 25bp rate increase could impact portfolio value by approximately -1.45%. However, this sensitivity varies across our holdings based on individual credit metrics and industry exposure.

Please note that this analysis is based on current market conditions and regulatory frameworks. Past performance does not guarantee future results."

# Use a Global Test?

## Global Test:

Is this explained in the language of a financial analyst working at a regulated firm?

# Global versus Unit Tests

## Global Test:

Is this explained in the language of a financial analyst working at a regulated firm?

## Unit Tests:

Context: Are relevant market conditions or external factors acknowledged?

Clarity: Is complex financial information presented in an accessible way?

Precision: Is terminology used accurately and consistently?

Compliance: Does the response adhere to relevant financial regulations and disclosure requirements?

Actionable: Does the response provide clear next steps or implications?

Risks: Are potential risks clearly identified and explained?

# Scoring Global and Unit Tests

## Global Test:

The analysis demonstrates strong awareness of the broader economic environment by effectively connecting Federal Reserve actions to specific market impacts. It skillfully weaves together monetary policy, inflation data, and manufacturing indicators to create a comprehensive market picture. The interconnections between these factors are well-established, though there could be more detail on international market dynamics.

The technical quality of the communication stands out particularly well. The precision in citing specific figures - from the basis point movements to the portfolio duration and percentage breakdowns - reflects a rigorous quantitative foundation. The use of terms like "spread compression" and "duration" is consistently accurate throughout, showing deep technical competency.

From a readability perspective, the response achieves a reasonable balance between technical depth and accessibility, though there's room for improvement. While financial professionals would readily grasp concepts like "duration" and "spread compression," these terms might benefit from brief contextual explanations for a broader audience.

The response carefully incorporates appropriate regulatory language and necessary disclaimers about forward-looking statements and past performance. The discussion of portfolio impacts stays within proper compliance boundaries while maintaining substantive analysis.

However, the response falls somewhat short in providing concrete direction. While it effectively outlines the situation, it stops short of offering specific recommendations or a clear decision framework. A business leader might finish reading this and still be uncertain about what specific steps to consider next.

The treatment of risk factors, while present, could be more comprehensive. While duration risk is quantified, other relevant risks like credit deterioration or liquidity constraints receive less attention. The analysis would benefit from a more structured exploration of various risk scenarios and their potential portfolio impacts.

Overall, this represents a technically sound piece of financial communication that excels in accuracy and market context, but could be enhanced with more actionable insights and a more thorough risk
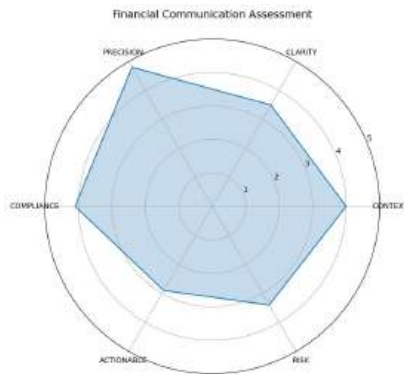
## Unit Tests:

Precision (4.8)

Context (4.0)

Compliance(4.1)

Clarity (3.5)

Actionable (2.9)

Risk (3.4)


Financial Communication Assessment

# Analyzing Failures with Clustered Patterns

Used Kmeans on 40 responses and found 4 clusters of failures:

Cluster 0: Synthesis Failure
"Found the right data but drew wrong conclusions"
High EVIDENCE/FACTUAL, Low REASONING/INFERENCE

Cluster 1: Context Failure
"Trees visible, forest missing"
High specific scores but low context/risk awareness

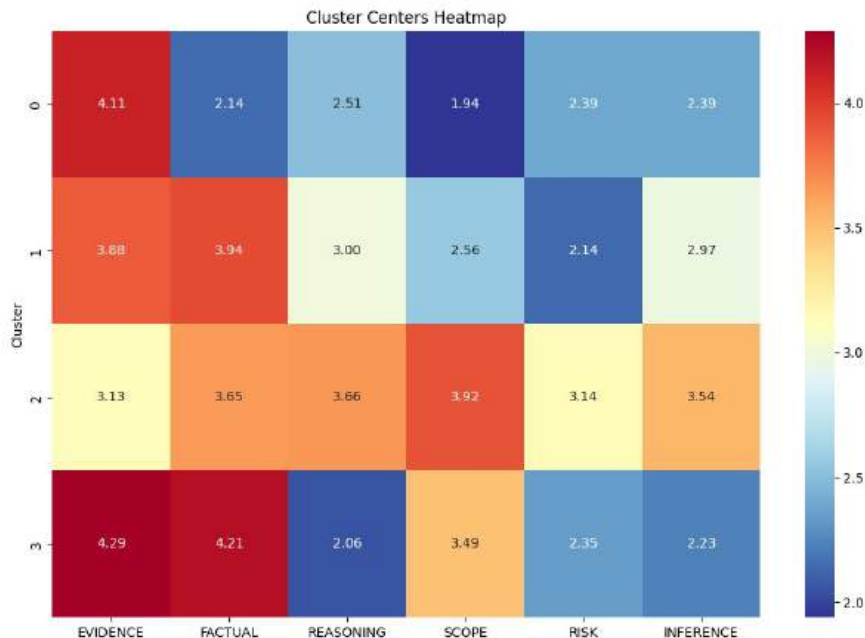Cluster 2: Hallucination Pattern
"Found good sources but made things up anyway"
High EVIDENCE but low FACTUAL/SCOPE

Cluster 3: Incomplete Retrieval
"Missing pieces but careful with what it has"
Medium EVIDENCE but consistent REASONING



Cluster Centers Heatmap

| Cluster | EVIDENCE | FACTUAL | REASONING | SCOPE | RISK | INFERENCE |
|---|---|---|---|---|---|---|
| 0 | 4.11 | 2.14 | 2.51 | 1.94 | 2.39 | 2.39 |
| 1 | 3.88 | 3.94 | 3.00 | 2.56 | 2.14 | 2.97 |
| 2 | 3.13 | 3.65 | 3.66 | 3.92 | 3.14 | 3.54 |
| 3 | 4.29 | 4.21 | 2.06 | 3.49 | 2.35 | 2.23 |

https://github.com/ContextualAI/examples/blob/main/03-standalone-api/01-lmunit/lmunit.ipynb

@rajistics

# How to Design Good Unit Tests

Use LLMs to help generate them

- Keep tests focused and specific
- Avoid compound criteria
- Use clear, unambiguous language
- Assess a desirable quality, such as "Is the response coherent?" rather than "Is the response incoherent?"

Ideally binary or small range (1-5)  don't do 1-100

https://github.com/ContextualAI/examples/blob/main/03-standalone-api/01-lmunit/lmunit.ipynb

# Examples of Global to Unit Tests

**Legal:**

Attribution

Copyright

PII

Permissible

Disclaimers

Confidential

**Retrieval:**

Evidence

Factual Accuracy

Reasoning

Scope

Omission

Inference Validity

**Bias/Fairness:**

Equitable
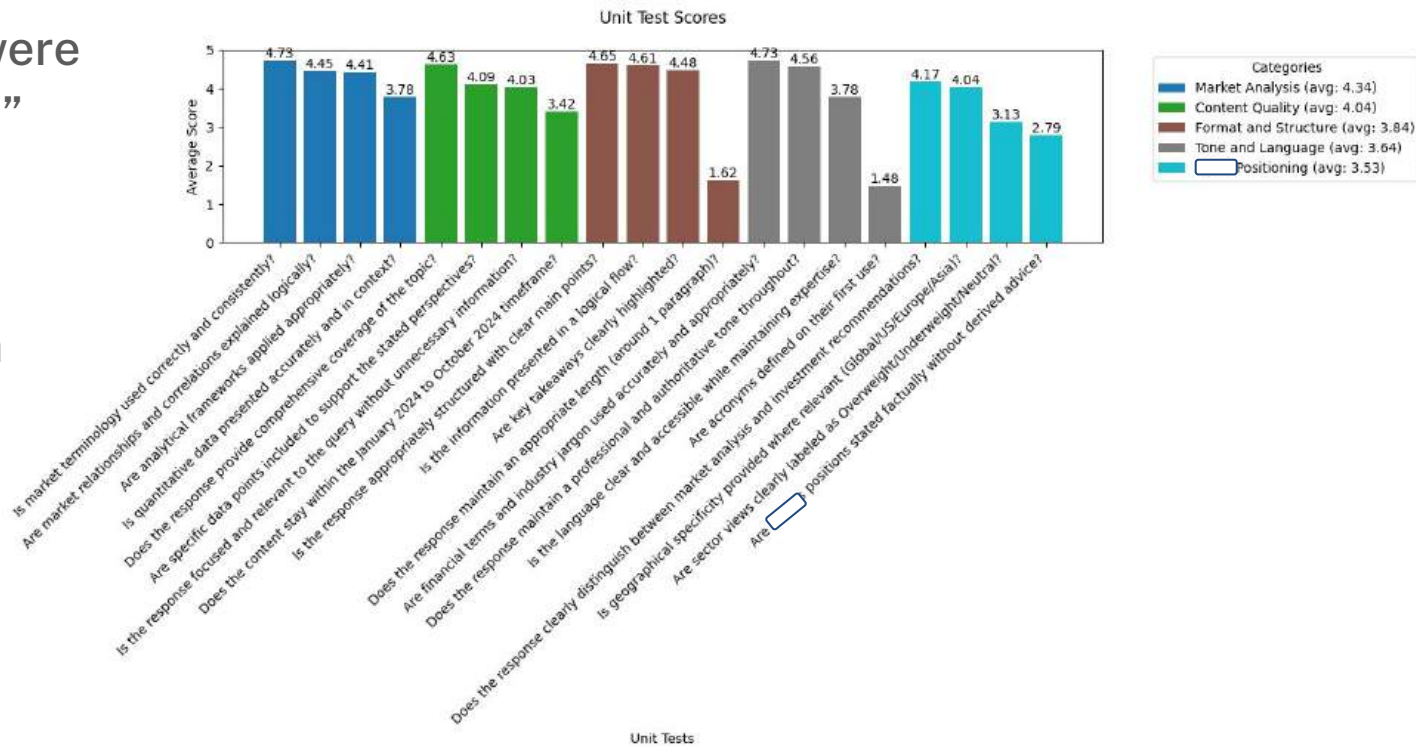
Representation

Neutrality

Stereotyping

Inclusivity

Imbalance

# Unit Tests to Evaluate New Prompts

These unit tests were built from a "good" system prompt

These are used to judge new system prompts



@rajistics

# Evaluation Tools – No Silver Bullet

Don't start with evaluation tools - you need to learn the basics first

Then move to logging traces/logging experiments

There is no silver bullet in evaluation tools

Dataset versioning

@rajistics

# Error Analysis

**Forest: Global / Integration**

**Trees: Test Case / Unit Tests**

Both of these come into play for most use cases and error analysis

@rajistics

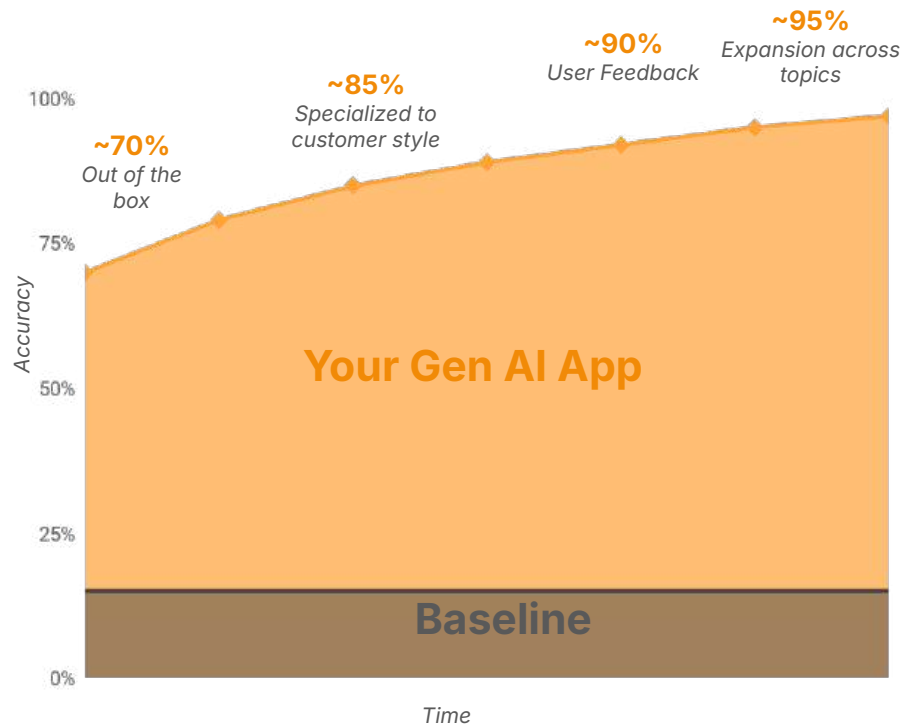# Error Analysis Tips



Compare performance one setting at a time

# Error Analysis Tips

1. Change one thing at a time – ablation style
2. Categorize failures
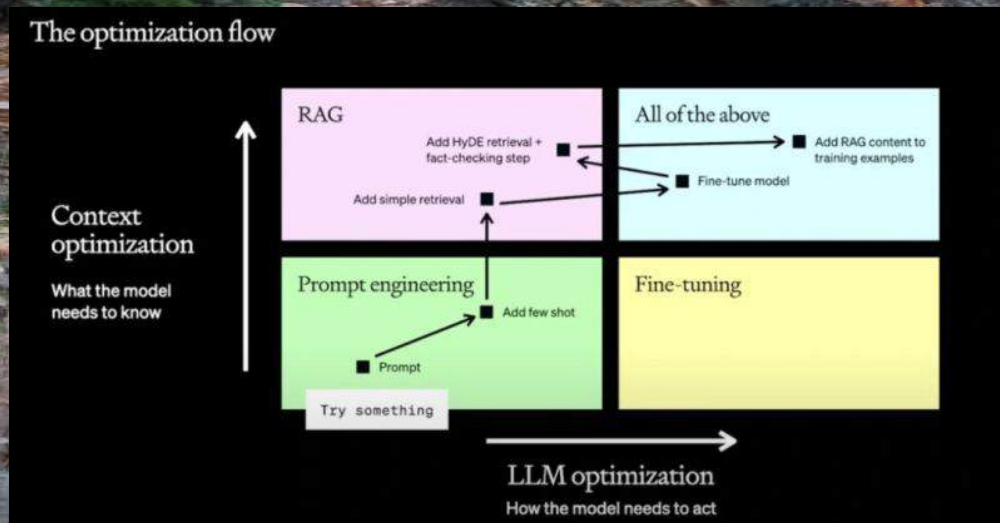3. ❤️ examples
4. Logs and traces

# The Evaluation Story We Tell

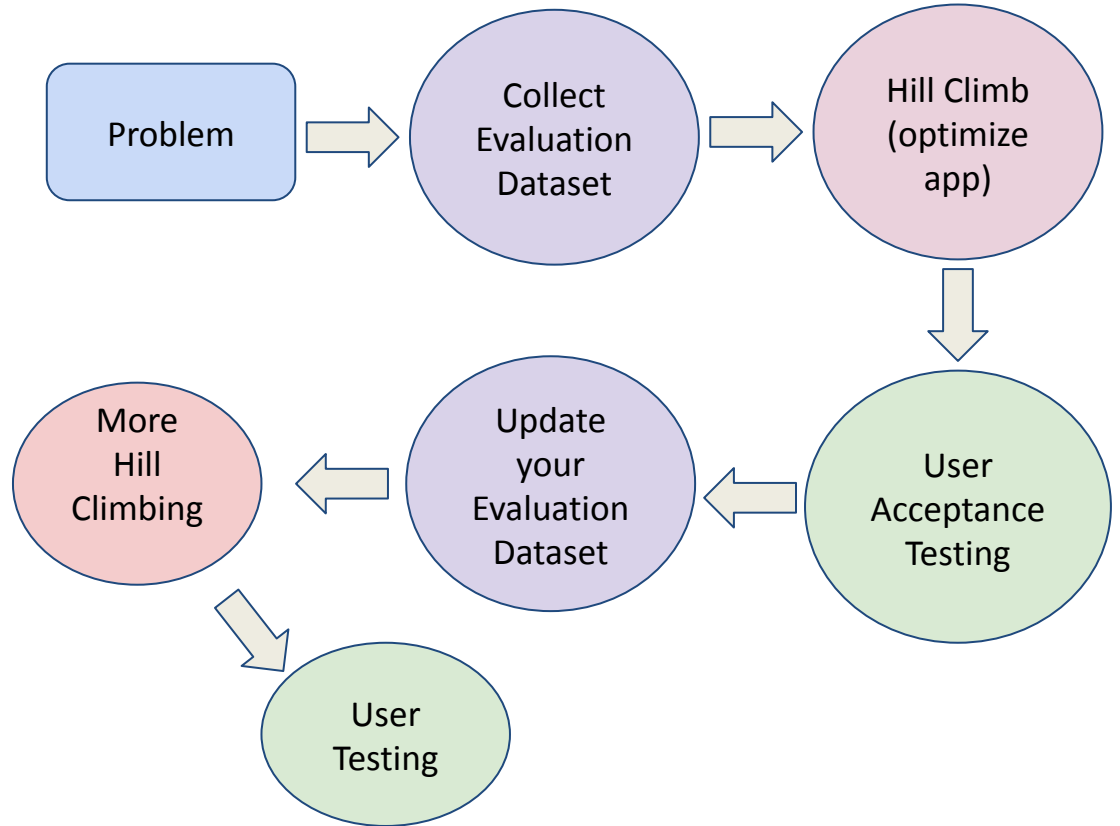**Nice linear history of natural progression**

# The Reality of Progress

**Progress is not linear
You care continually
evolving and some
paths work, others
don't**



The optimization flow

RAG — Add HyDE retrieval + fact-checking step — Add simple retrieval

All of the above — Add RAG content to training examples — Fine-tune model

Context optimization — What the model needs to know

Prompt engineering — Add few shot — Prompt

Fine-tuning

Try something

LLM optimization — How the model needs to act

https://www.youtube.com/watch?v=ahnGLM-RC1Y

# Evaluation as a Continual Process
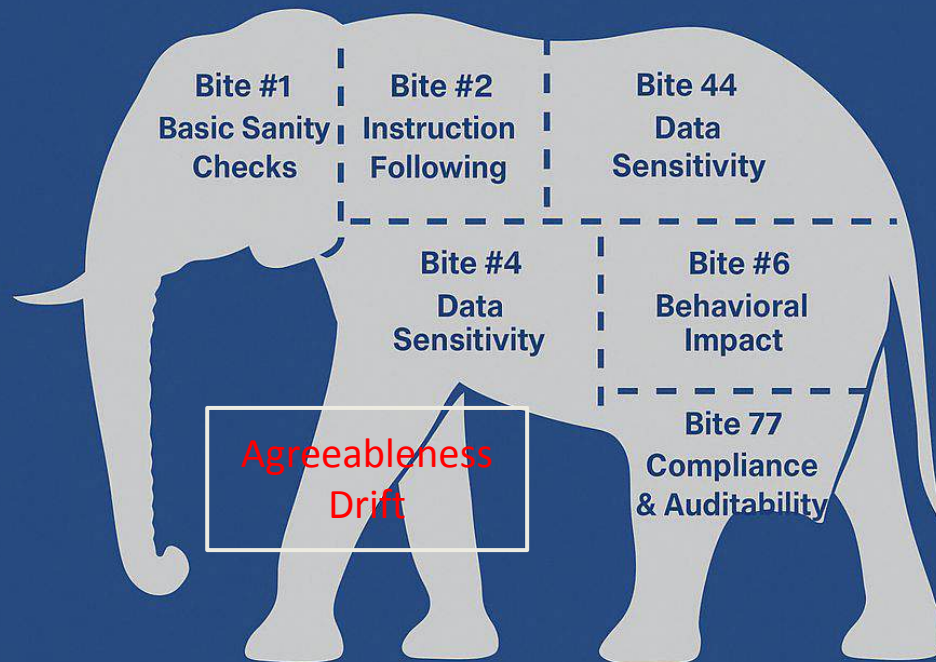
How do you eat an elephant?

# Adding Tests Over Time

You can't do it all initially, measured approach that gradually widens

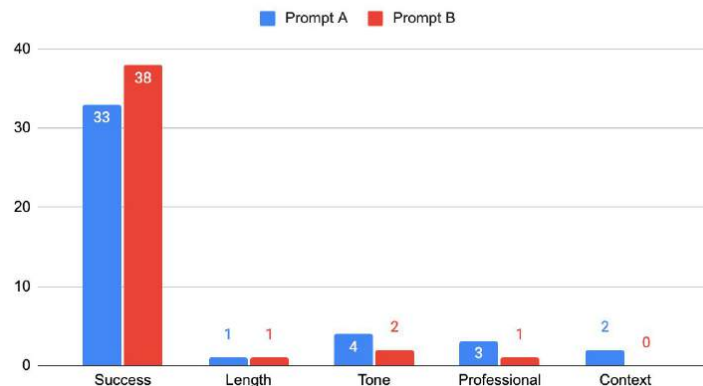You Cannot Feed Two Birds with One Score: the Accuracy-Naturalness Tradeoff in Translation https://arxiv.org/pdf/2503.24013

@rajistics



How to Eat "GenAI Evaluation Elephant One Bite at a Time

**Bite #1** Basic Sanity Checks

**Bite #2** Instruction Following

**Bite 44** Data Sensitivity

**Bite #4** Data Sensitivity

**Bite #6** Behavioral Impact

Agreeableness Drift

**Bite 77** Compliance & Auditability

# Doing Evaluation the Right Way

- Start with annotated examples representing success and failure cases
- Document issues systematically and build focused tests for each
- Error analysis is a continuous process, not a one-time effort
- Collaborate with domain experts, users, and leverage bootstrapping approaches
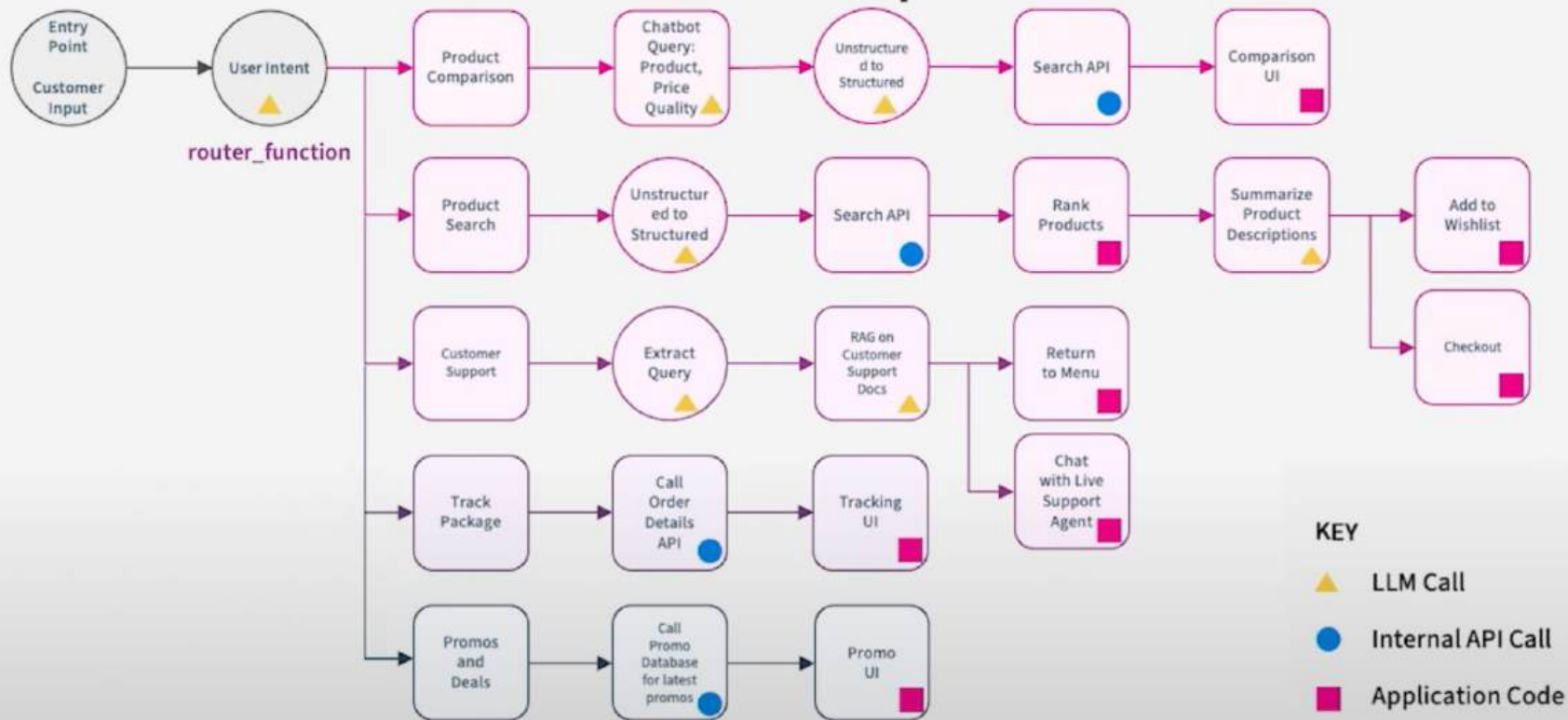- Cognizant of generalization and overfitting to your dataset

**Test Case Results**

Prompt A ▪ Prompt B ▪

| | Success | Length | Tone | Professional | Context |
|---|---|---|---|---|---|
| Prompt A | 33 | 1 | 4 | 3 | 2 |
| Prompt B | 38 | 1 | 2 | 1 | 0 |

@rajistics

**Agentic use cases**

@rajistics

# How should it cross the river?

- Fly
- Swim
- Other?

# Under the Hood: Chat-to-purchase Router



Jason from Arize: https://www.youtube.com/watch?v=dnfoeOb8ym4&t=558s

# Text to SQL Agent from Snowflake



https://www.snowflake.com/engineering-blog/snowflake-cortex-analyst-behind-the-scenes/
Raj's explanation: https://youtu.be/OyY4uxUShys

@rajistics

# Evaluating Office-Style Agent Workflows (OdysseyBench)



Figure 4: Typical failure cases of the LLM agents when solving office automation tasks in OdysseyBench.

https://arxiv.org/pdf/2508.09124

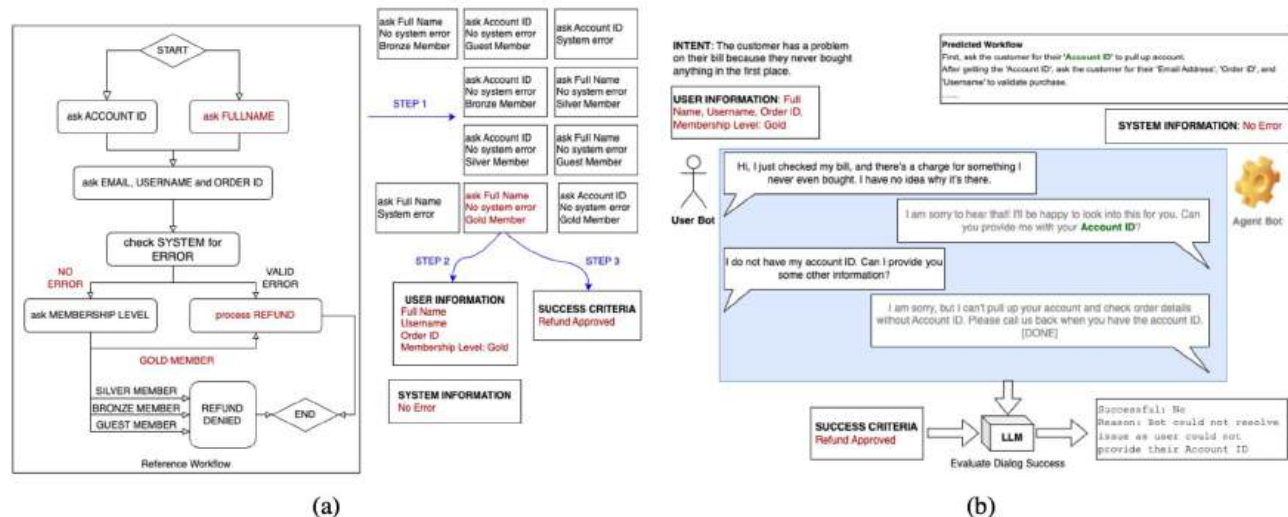# Error Analysis for Agentic Workflows

Assess overall performance

- Assess routing

- Assess individual agent steps

  - Where is it not working

  - Why isn't it not working

  - Fix it and reassess



Performance and Error Analysis

@rajistics

# Evaluating a Workflow Instead of a Response

Flowchart of conversational flow



(a)

(b)

https://arxiv.org/pdf/2502.17321
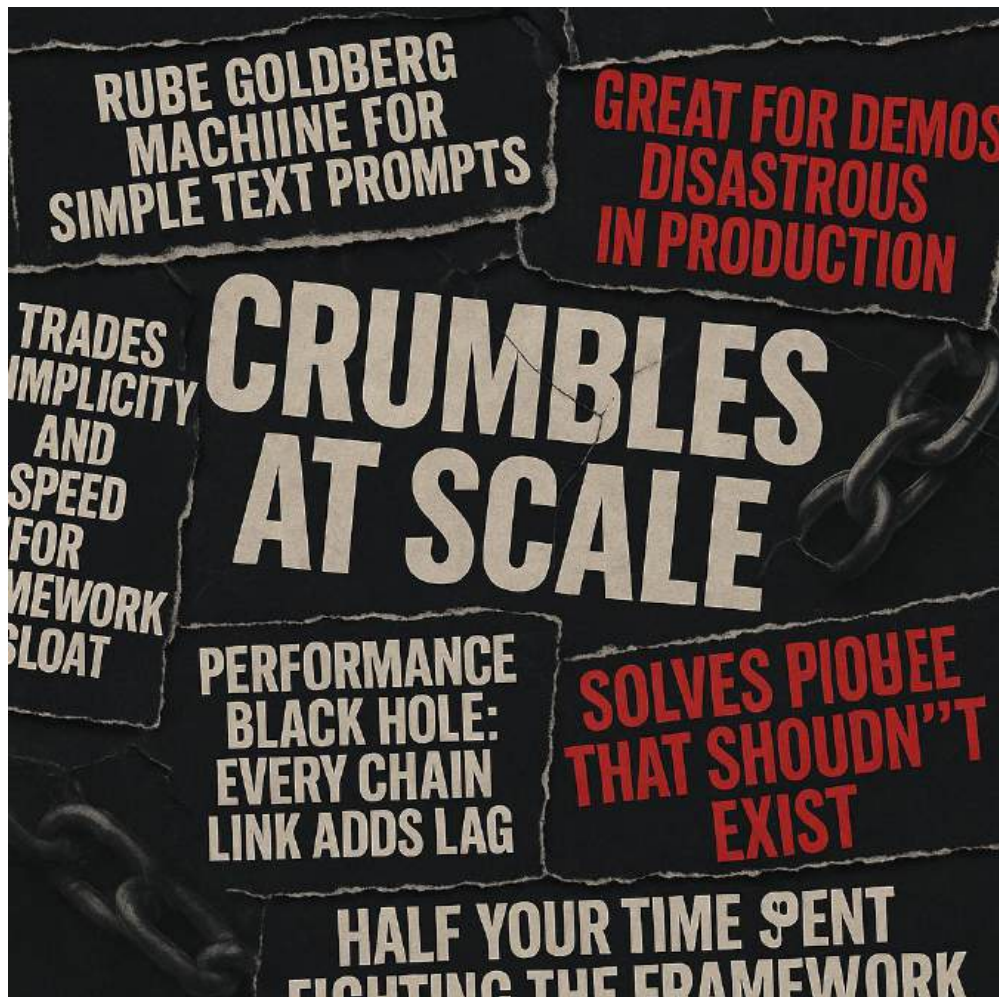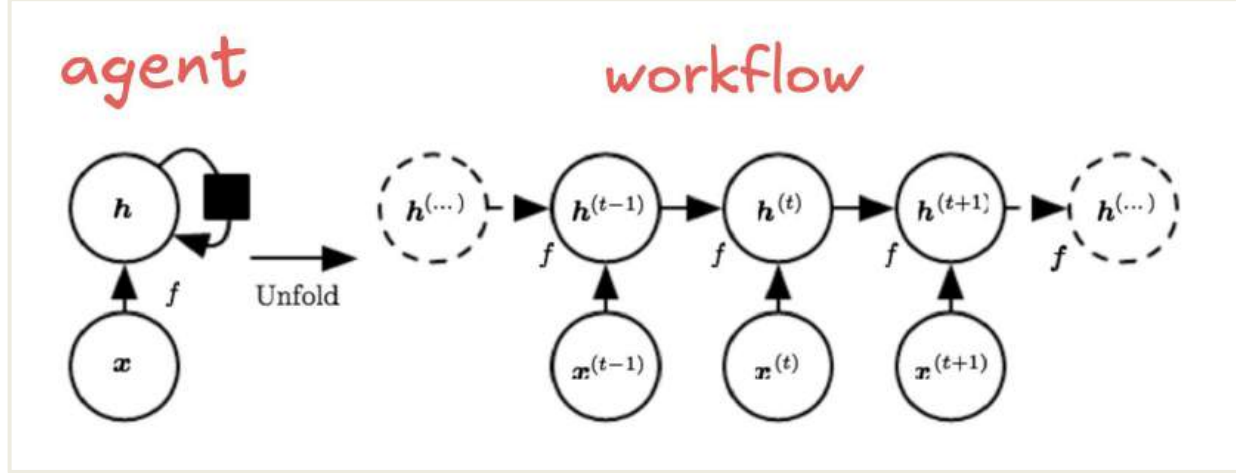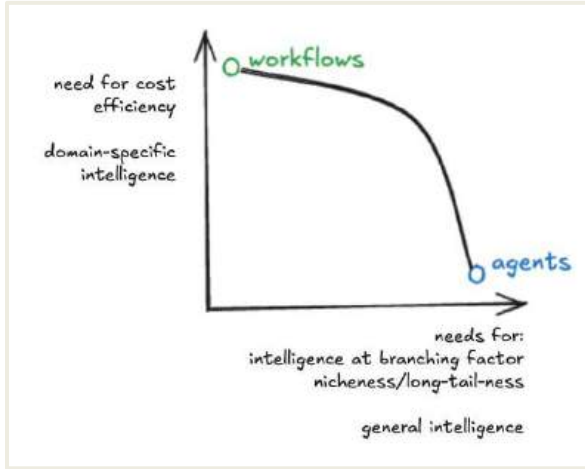
# Agentic Frameworks Help – Until They Don't

- Tools, packages, solutions that will abstract away all the technical details
- This is great until it breaks
- This is great until its outdated
- This is great until you need to customize it



@rajistics

# Abstraction for Agentic Workflows



https://x.com/HanchungLee/status/1916280769783460210

# When Agent Abstractions Break Down

## Declarative vs non-declarative graphs

Some frameworks are declarative, requiring developers to explicitly define every branch, loop, and conditional in the workflow upfront through graphs consisting of nodes (agents) and edges (deterministic or dynamic handoffs). While beneficial for visual clarity, this approach can quickly become cumbersome and challenging as workflows grow more dynamic and complex, often necessitating the learning of specialized domain-specific languages.

In contrast, the Agents SDK adopts a more flexible, code-first approach. Developers can directly express workflow logic using familiar programming constructs without needing to pre-define the entire graph upfront, enabling more dynamic and adaptable agent orchestration.

https://blog.langchain.dev/how-to-think-about-agent-frameworks/

# Lesson from Reproducing Agent Benchmarks

- Standardize Evaluation – Many "breakthroughs" vanished after reruns due to bugs and data leakage.

- Measure Efficiency, Not Just Accuracy – More reasoning often hurt accuracy while driving up compute cost.

- Detect Shortcuts – Agents sometimes pulled benchmark data or gamed tests instead of solving tasks.

- Log Real Behavior – Tool-use errors and unrealistic actions show why accuracy alone isn't enough.

HAL: https://arxiv.org/pdf/2510.11977
Survey: https://arxiv.org/pdf/2503.16416

@rajistics

# We did it!

## Code & Slides:

https://github.com/rajshah4/LLM-Evaluation