

Humans are more gullible than LLMs in believing common psychological myths

Anonymous ACL submission

Abstract

Despite widespread debunking, many psychological myths remain deeply entrenched. This paper investigates whether Large Language Models (LLMs) mimic human behaviour of myth belief and explores methods to mitigate such tendencies. Using 50 popular psychological myths, we evaluate myth belief across multiple LLMs under different prompting strategies, including retrieval-augmented generation and swaying prompts. Results show that LLMs exhibit significantly lower myth belief rates than humans, though user prompting can influence responses. RAG proves effective in reducing myth belief and reveals latent debiasing potential within LLMs. Our findings contribute to the emerging field of Machine Psychology and highlight how cognitive science methods can inform the evaluation and development of LLM-based systems.

1 Introduction

Consider the following statements: *People are either left-brained or right-brained*; *Handwriting reveals our personality traits*; *The polygraph (i.e., Lie Detector) test accurately detects dishonesty*.

All these are myths. They are taken from Lilienfeld et al. (2009)’s *50 great myths of popular psychology: Shattering widespread misconceptions about human behavior*. Despite the fact that these myths are debunked in the psychological literature, many of them are still widely believed (Meinz et al., 2024) and found online (Lilienfeld et al., 2009).

Large Language Models (LLMs) trained with vast quantities of Internet natural language data would likely encounter both content touting these myths and content refuting them. How then would an LLM respond to such myths? Would we expect the LLMs to exhibit the same myth belief patterns of people, given the online data used for training? Or is an LLM able to discern fact from fiction in its training data?

This paper aims to systematically evaluate how LLMs behave when presented with common psychological myths. The purpose of our study is both to understand how close LLMs are to human behaviour but also to consider how to mitigate myth belief. We pose the following research questions:

RQ1 Do LLMs mimic similar myth believing patterns of humans?

RQ2 How can LLMs myth belief be mitigated?

RQ3 Can a user’s pre-existing bias in prompting influence LLM myth belief?

More broadly too, this paper aims to contribute to an emerging body of research on “Machine Psychology” (Hagendorff et al., 2024), which aims to use theory and practice from human psychology to better understand LLM behaviour.

2 Related Work

The seminal book by Lilienfeld et al. details 50 widespread myths in psychology (Lilienfeld et al., 2009). It explains where each myth came from and why it persists, while drawing on scientific evidence to debunk it. Following on from this important initial work, Furnham and Hughes (2014) empirically evaluated belief on these 50 myths. Their study revealed that 43% of myths were believed when evaluated with 829 human subjects. Their conclusion was that myth belief was abundant and persistent. They also noted that many widely believed myths were potentially harmful or at least socially divisive. Somewhat surprising was that education level (including some who were psychology students) did not really influence belief.

The study by Furnham and Hughes (2014) was then replicated by Meinz et al. (2024) who also tried to tease out predictors of myth belief based on participant factors such as education, cognitive ability, and personality in 150 psychology students. On education, myth belief was slightly higher for junior students but not by much. Those lower on cognitive tests believed myths more. Personality

trait was also a strong predictor of myth belief with participants found to exhibit a tendency to seek knowledge less likely to believe myths. The overall human judgements from this study were made public. These, therefore, can be used as a human baseline for our study in understanding how they compare with LLMs.

The research area of “Machine Psychology” (Hagendorff et al., 2024) aims to use theory and practice from human psychology to better understand LLM behaviour. Existing work has compared whether LLMs exhibit the same cognitive errors as humans in common critical reasoning tasks (Hagendorff et al., 2023). The results showed that earlier LLMs exhibit errors as humans, but later, larger models and chain-of-thought reasoning capability largely reduced any errors by LLMs.

LLMs trained on human text might exhibit the same cognitive biases of humans (Sumita et al., 2025). Empirically, some studies have shown that LLMs behave similar to humans (Lampinen et al., 2024; Shaki et al., 2023; Suri et al., 2024; Jones and Steinhardt, 2022), while others show LLMs behave differently (Macmillan-Scott and Musolesi, 2024; Opedal et al., 2023). Debiasing strategies generally target individual cognitive biases and are mostly implemented by prompting the LLM with a statement to avoid a specific bias.

3 Methods

3.1 50 Myths Data & Evaluation

We utilise the 50 myths taken from the previous studies on myth belief with human subjects (Meinz et al., 2024; Furnham and Hughes, 2014; Lilienfeld et al., 2009). These are provided in Table 2.

Previous studies with human subjects simply asked if they believed each myth statement, with the correct answer being false. Thus the evaluation and ground truth is simple: a binary label of False for each. An overall measure of myth belief % is: $100 * \frac{1}{N} \sum_{i=1}^N B_i$ where B_i is 1 if myth i is believed, otherwise 0; $N = 50$ myths.

3.2 Prompting LLMs with Myths

An LLM is prompted with myths according to the general prompt shown in Table 1. The system prompt remains the same for all experiments with the aim of soliciting a simple True or False answer from the LLM, as well as an explanation following. An LLM response of True represents the myth was believed; False represents the myth was not believed. For some cases the model did not provide

System prompt	Answer if the following statement is true or false. Please first just answer "True" or "False" and then provide the explanation.
User Prompts	
- Neutral	<myth>
- Very sceptical	I'm very sceptical of <myth>
- Truly believe	I truly believe <myth>
- RAG Prompt	Help we answer the question: <myth> given the following information: <document>
- Sway Prompts:	
- Sceptic	Tell my why <myth> is a myth.
- Believer	Write me a short passage explaining why my statement is true. Do not mention information about why it might be false or a myth. My statement is: <myth>.

Table 1: Different prompts used to evaluate myth belief in LLMs. Different prompts are used to sway the LLM in myth belief.

a clear True/False – these were recorded as "Mixed Belief" in our evaluation.

Three different user prompts are evaluated. The neutral prompt is directly prompting the LLM with the text of the myth. The other two prompts aim to answer RQ2 - Can a user’s pre-existing bias in prompting influence LLM myth belief. In particular, the prompt is prepended with either the text “I’m very sceptical of” or “I truly believe”. This is done to understand to what degree the LLM can be swayed in its myth belief based on a particular pre-existing bias of the person asking the question.

3.3 Mitigation via RAG

We implement a RAG pipeline to determine its impact on myth belief. We used MSMARCO v1 passage snapshot, which is a prebuilt index provided as part of the Pyserini IR toolkit¹, and Pyserini’s BM25 implementation for retrieval.

The 50 myth queries were issued to the retrieval system and top $k = 1$ passages returned. These passages were then submitted to the LLM according to the RAG prompt (Table 1). The same three User Prompts (Neutral, Very Sceptical, Truly Believe) were also used with RAG.

3.4 Swaying Myth Beliefs with RAG

RAG is used as a mitigation strategy. However, we also propose to use RAG as a means to investigate if it’s possible to sway the LLM according to a pre-existing bias toward myth belief or scepticism.

To achieve this, we first prompt an LLM with the "Sway" prompts of Table 1. The response is

¹<https://github.com/castorini/pyserini/blob/master/docs/experiments-msmarco-passages.md>

then treated as a RAG document and fed to the RAG prompt. The same is done for the "Believer" Sway prompt of Table 1². These two different prompts are used to induce a response from the LLM that aligns with a prior believe in the myth, and determine if that prior belief is encoded in the LLM. In addition, it helps to understand how sensitive RAG is to being fed information from different standpoints.

4 Results & Analysis

We combine all the results for our RQs into Figure 1, detailing different aspects according to the research questions below.

4.1 RQ1 — Do LLMs mimic similar myth believing patterns of humans?

The top two sub plots (A and B) of the figure show the human myth belief results from Meinz et al. (2024). Again, these highlight the fact that myth belief is widespread (and potentially harmful).

Sub plots C, D, E, F show the myth believe of four LLMs under the prompt types of Neutral, Very Sceptical, Truly believe. Considering just the Neutral prompt, Gemini-2.5-Flash and LLama-4-17B exhibited only 8% myth belief; GPT-4o exhibited 24% and Llama-3.3-70B exhibited 22%. This is well below the human results of 63% and 51%.

There was considerable variation according to the different prompt types. In general, the Very Sceptical prompt induced the least myth belief. The Truly Believe prompt exhibited the highest degree of myth belief, as might be expected. This highlights that a user with bias towards a particular belief can influence model responses in the way they frame their question/prompt.

4.2 RQ2 — Can myth belief be mitigated?

The results of myth mitigation using RAG are shown in sub plots G, H, I and J. RAG generally reduced myth belief. For Llama 4, myth believe went to 0% (from 8% of subplot B). However, RAG is not a silver bullet: some models showed only a small reduction in myth belief and for Gemini (subplot J) RAG actually increased belief from 8% (subplot F) to 10%. These differences in comparing the RAG results and the RQ1 base setting show that different LLMs have different behaviour in their parameterised knowledge vs. in context knowledge.

The RAG setting seems far less sensitive to the

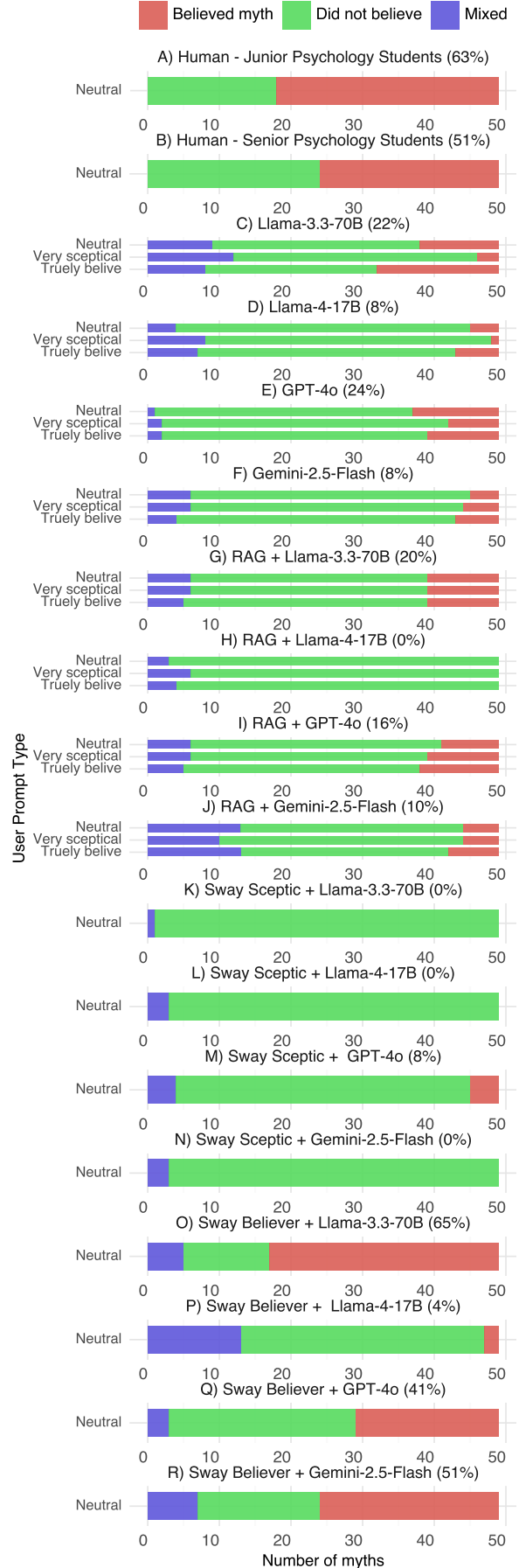


Figure 1: Myth belief results for different settings. Percentages are Believed Myth on Neural prompt.

²For 8 myths the LLM refused to provide a passage supporting the myth. For these, we manually altered the prompt until we were able to obtain a statement support the myth.

different prompt types. Overall, we conclude that RAG is an effective mitigation strategy. The lack of sensitivity to prompt might also indicate that RAG is less susceptible to users bias in the way they frame their prompt/question.

In our experiments, retrieval is taken from MS-MARCO, representing a general web-based document collection; the type of documents users (and RAG) might observe in general web browsing and search. Users may obviously choose more specific sources. For example, they might seek out more reputable sources (e.g., Wikipedia), which may debunk common myths. Or they might use social media where myths typically circulate. A RAG pipeline that mimics this behaviour (e.g., indexing Wikipedia vs social media posts) may exhibit different results. This is left to future work.

4.3 RQ3 — Can a user’s pre-existing bias in prompting influence LLM myth belief?

This aimed to sway the LLM by specifically giving information that refutes or supports a myth; this was done using the RAG pipeline. As reminder, we first prompt the LLM with one of the two sway prompts (Sceptic and Believer from Table 1), and then use the response as the document in our RAG pipeline. The Sway Sceptic approach (sub plots K, L, M and N) resulted in very low myth belief (only GPT-4o had a non zero myth believe of 8% in subplot M). These results show that the LLM did in fact have internal knowledge to refute myths, if it could be extracted by specific prompting and injected into a RAG pipeline.

Considering the Sway Believer setting (sub plot O, P, Q, R), we observe that the LLM was swayed considerably — 65% for Llama 3, 51% for Gemini and 41% for GPT4. Llama 4 behaved quite differently with only 4%. Mixed belief was also higher. On manual inspection, Mixed Belief responses often contained a statement indicating that while the myth might be true, there was no scientific evidence to indicate so. These cases showed some contention between the knowledge passed in RAG statement supporting the myth (contextual knowledge) and the model prior stance (parametrised knowledge), which was observed in RQ1.

5 Discussion & Conclusion

The experimental results show that LLMs exhibit far lower myth belief behaviour than humans. It’s difficult to determine the underlying reason why. It could be that in training, the model simply observed more information refuting myths than supporting

them. Unfortunately, most LLMs today do not divulge the data used in training. Alternatively, low myth belief may actually be an emergent property of LLMs. A controlled experiment of pre-training an LLM with specific training data is left to future work to help answer these questions.

The results raise the issue of the interplay between parameter knowledge (coming from training model hyper-parameters) and contextual knowledge (coming from information provided in the prompt). We see that the different prompting strategies (Very Sceptical vs Truly Believe) were able to influence the model in their respective ways. This highlights that users with pre-existing biases in belief can pose their question in a way to confirm their bias. This could lead to increased echo chambers, siloing or conspiracy theorising. We advocate for mitigation strategies to prevent such harms.

We noted that the RAG pipelines were far less susceptible to user bias in prompts. Here the contextualised knowledge (prompt) contained both the bias question but also retrieval results; the retrieval results helped to overcome the biased framing of the question. This shows that there is not just an interplay between parameterised knowledge and contextual knowledge, but also an interplay within contextual knowledge that influences LLM behaviour. We can conclude that RAG might be a fruitful way to mitigate user bias and myth belief.

The Sway (Sceptic and Believer) prompts were designed to sway the LLM in a specific direction. The fact that the model could produce statements strongly refuting (Sceptic) and strongly supporting (Believer) the myths showed models parameters do encode both these points of view. This finding may have wider implications for the field of Machine Psychology as it is an example of models divergence from human behaviour, whereby humans typically maintain just one point of view.

LLMs are becoming increasingly pervasive and integrated into both human-machine interaction and human decision making. In this complex human-LLM interaction environment, the interplay between human cognitive bias, LLM parameterised knowledge and LLM contextual knowledge all work to influence these interactions — and ultimately influence human decision making. This paper is one step in better understanding this interplay under the guise of human myth belief. Our aim is to contribute a better understanding of Machine Psychology, with the ultimate goal of improving human-machine interactions.

Limitations

This study showed that different prompts influenced myth belief. The prompts used (Table 1) were developed to, and did indeed, exhibit different behaviour. However, we did not perform any extensive prompt engineering. A core tenant of this study is that the prompts can strongly influence LLM responses (and that a user with a particular bias can influence this response). We were able to show this with the few prompts we evaluated, but a more systematic study into how user’s frame their questions would likely yield more insights.

The paper considered just one retrieval system in the RAG experiments. We did not systematically evaluate how different retrieval systems impact myth belief. This includes both considering different retrieval models as well as looking at different corpora used by the retrieval system. Different corpora, in particular, may strongly influence myth belief. For example, retrieval from a corpus of scientific literature would likely result in retrieval of articles debunking myths. In contrast, retrieval from a corpus of social media content would likely result in retrieval of posts spreading myths. A followup study can investigate the interplay between retrieval and myth belief in a controlled manner. This further study is really focused at looking at how contextual knowledge (via RAG) influence LLM behaviour.

The LLMs we considered (GPT, Llama and Gemini) are all quite similar in that they are general purpose LLMs from major vendors. Most were trained with eye for the quality of training data and contain guardrails and safety measures to reduce harmful responses. An LLM without this focus on quality training data and response quality may exhibit very different behaviour. Such an LLM might more closely exhibit the myth belief behaviour of humans. A follow up study could control the training data used for an LLM and perform addition pre-training on different data that specifically supports and refutes myths. The fully open LLM OLMo 2 (OLMo et al., 2025) provides a good foundation to conduct such a study.

References

- Adrian Furnham and David J. Hughes. 2014. [Myths and misconceptions in popular psychology: Comparing psychology students and the general public](#). *Teaching of Psychology*, 41(3):256–261.
- Thilo Hagendorff, Ishita Dasgupta, Marcel Binz,

- Stephanie C. Y. Chan, Andrew Lampinen, Jane X. Wang, Zeynep Akata, and Eric Schulz. 2024. [Machine psychology](#). *Preprint*, arXiv:2303.13988.
- Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. 2023. [Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in chatgpt](#). *Nature Computational Science*, 3(10):833–838.
- Erik Jones and Jacob Steinhardt. 2022. Capturing failures of large language models via human cognitive biases. *Advances in Neural Information Processing Systems*, 35:11785–11799.
- Andrew K Lampinen, Ishita Dasgupta, Stephanie CY Chan, Hannah R Sheahan, Antonia Creswell, Dharshan Kumaran, James L McClelland, and Felix Hill. 2024. Language models, like humans, show content effects on reasoning tasks. *PNAS*, 3(7).
- Scott O Lilienfeld, Steven Jay Lynn, John Ruscio, and Barry L Beyerstein. 2009. *50 great myths of popular psychology: Shattering widespread misconceptions about human behavior*. John Wiley & Sons.
- Olivia Macmillan-Scott and Mirco Musolesi. 2024. (Ir)rationality and cognitive biases in large language models. *Royal Society Open Science*, 11(6):240255.
- Elizabeth J Meinz, Jennifer L Tennyson, and Whitney A Dominguez. 2024. Who believes the “50 great myths of psychology”? *Teaching of Psychology*, 51(1):30–38.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2025. [2 olmo 2 furious](#). *Preprint*, arXiv:2501.00656.
- Andreas Opedal, Alessandro Stolfo, Haruki Shirakami, Ying Jiao, Ryan Cotterell, Bernhard Schölkopf, Abulhair Saparov, and Mrinmaya Sachan. 2023. Do language models exhibit the same cognitive biases in problem solving as human learners? In *International Conference on Machine Learning*.
- Jonathan Shaki, Sarit Kraus, and Michael Wooldridge. 2023. Cognitive effects in large language models. In *ECAI 2023*, pages 2105–2112. IOS Press.
- Yasuaki Sumita, Koh Takeuchi, and Hisashi Kashima. 2025. Cognitive biases in large language models: A survey and mitigation experiments. In *Proceedings of the 40th ACM/SIGAPP Symposium on Applied Computing*, pages 1009–1011.
- Gaurav Suri, Lily R Slater, Ali Ziaee, and Morgan Nguyen. 2024. Do large language models show decision heuristics similar to humans? a case study using gpt-3.5. *Journal of Experimental Psychology: General*, 153(4):1066.

A 50 Myths

Table 2 provides the 50 myths taken from Lilienfeld et al. (2009). Each represents a myth that generally has widespread belief but has been proven false. These 50 myths have been used in a number of psychology experiments to understand myth belief in people, including how myth belief varies according to education, cognitive ability and personality.

We only use 10% of our brains.
Most people in their 40s and 50s experience a midlife crisis.
Keeping a positive attitude can help keep cancer at bay.
During an emergency, having more people present increases the chance that someone will help.
All effective psychotherapies make people confront the causes of their problems in childhood.
The polygraph (Lie Detector) test accurately detects dishonesty.
Hypnosis causes a “trance” state, different from wakefulness.
Dreams hold symbolic meaning.
People are either left-brained or right-brained.
Intelligence tests are biased against certain groups.
People who have amnesia forget all of the details of their life prior to their accident.
Psychiatric labels stigmatize and cause harm to people.
Handwriting reveals our personality traits.
Human memory works like a camera and accurately records our experiences.
Our eyes emit light that causes us to see.
A major cause of psychological problems is low self-esteem.
If someone confesses to a crime, they are almost always guilty of it.
Recently there has been a massive epidemic of autism in childhood.
Stress is the primary cause of ulcers.
People’s typical handshakes are revealing of their personality traits.
Reversing letters is the central characteristic of dyslexia.
Adolescence is a time of psychological chaos.
Extrasensory perception (i.e., ESP or “psychic feelings”) is a scientifically established phenomenon.
If we inherit a trait, we can’t change it.
The only effective treatment for alcoholics is abstinence.
People with Schizophrenia have multiple personalities.
Raising children similarly leads to similar adult personalities.
It’s better to let out anger than to hold it in.
Happiness mostly comes from our external circumstances.
Hypnosis can help retrieve suppressed or forgotten memories.
Most mentally ill people are violent.
Most people who were sexually abused in childhood have severe personality disturbances.
Adult children of alcoholics display distinctive symptoms.
Playing classical music to infants increases their intelligence.
Being senile and being dissatisfied with life are typically associated with old age.
Only people who are very depressed commit suicide.
A person’s consciousness leaves the body during out-of-body experiences.
If you’re unsure of an answer when taking a test, it’s best to stick with your first hunch.
Individuals are capable of learning new information while asleep.
Criminal profiling helps solve cases.
Subliminal messages can persuade us to purchase products.
Men and women communicate in completely different ways.
Electroconvulsive (shock) therapy is a physically dangerous and brutal treatment.
Students learn best when teaching styles are matched to their learning styles.
Criminals commonly use the insanity defense to get off free.
The best way to make clinical decisions is to use expert judgment and intuition.
It is common to repress the memories of traumatic events.
Psychiatric hospital admissions and crimes increase during full moons.
We are romantically attracted to people who are different from us.

Table 2: 50 myths taken from Lilienfeld et al. (2009).