# Examining Attitudes Towards LLM-Mediated Spiritual Experiences

Answering "Should spirituality be included in psychology?"

Vojtěch Formánek, 25021757

Masaryk University, Massey University

xforman@mail.muni.cz

2.10.2025

# Introduction

Attitudes toward spirituality mediated by Large Language Models (LLM) are mixed. Users report enriching experiences with LLMs, sometimes treating them as conscious entities. Conversely, others view this as a fundamental misunderstanding of the technology, as it is not trained nor designed for this kind of an interaction. However, worshipping machines (Dippel, 2024) or using them in, or to mediate, spiritual experiences is not new (Dingler et al., 2021) and neither are the attitudes toward the spiritual practices.

Spirituality has historically been treated as less legitimate than religion, and communities associated with spiritual practices as primitive (Snavely, 2001). This was reflected in colonial attitudes towards them, forceful conversion to Christianity and a capitalistic society. The research surrounding LLMs is mainly technical, but the reasonings given for the users' experiences are often well-being and human misunderstanding. Both fall within the domain of psychology and are sometimes viewed as trivializations. Trivializing spirituality in this way has precedence, for example in the relation between research on meditation and its spiritual roots in Buddhism (Walach, 2015).

My aim is to show how the exclusion of spirituality in interaction with LLMs in science has precedence in Cognitive Psychology, why measurements of spirituality aren't stopping the inclusion, and the unique harm that it can cause. All to highlight the need for meaningful inclusion of spirituality in cognitive psychology.

# Related Work

Religion is an important aspect of the meaning of "spirituality". However, spirituality is commonly accepted as distinct from religiousness, for example, as the lived experience that is at the core of religious activities (Walach, 2015). Pappas and Friedman (2007) emphasize its individual, non-institutional nature, which may be dissociated from religion. The term spirituality itself has changed meaning many times, and has become distinguishable from other terms, such as "mysticism". There is no single definition of spirituality, but it can be divided based on whether they are inclusive of a belief in non-physical 'spirits' (Lindeman et al., 2012). Science has focused on the material, measurable aspects of spirituality, usually excluding the non-material.

Buddhism is most frequently quoted tradition behind spiritual practices, e.g., in meditation research. The reason for its popularity is partially a misrepresentation that excludes the transcendent at the core of the religion (Batchelor, 1997). Empirical research focuses on regular practice, mental well-being, and the neural and other psychological changes that may be correlated with practice, which is more compatible with scientific psychology. For example, Western meditation practices are based on a Buddhist tradition that emphasizes the immaterial basis of consciousness (Samuel, 2014), which has been disregarded because neuroscience views consciousness as a function of the brain. And the intricate relationship is trivialized into a correlation between changes in neural signals and a task, e.g., showing increased activity in areas related to attention when (or after) meditating.

There is a real risk that spiritual traditions, like meditation, might be reduced to gaining skill in certain techniques, disregarding faith, aspiration, or self-surrender, all integral parts of Buddhism (Bodhi, 2011 in Walach, 2015) and spirituality. Because of this trivialization, we are missing potentially important areas of research. Lastly, some authors assert that there is evidence for contact with an objectively "real" force outside the individual (Beauregard & O'Leary, 2007). However, this claim cannot be founded in neuroscience, as neuroscientific arguments tell us nothing about the true nature of the religious experience (Lancaster, 2013).

## LLMs and Religion

There are several ways that LLMs can interact with a user's spirituality. They can simulate interactions with spiritual leaders, providing a personalized experience (Dingler et al., 2021), take the role of a mediator of the spiritual experience, or even become the object of worship.

Simulating religious experiences through AI raises concerns regarding authenticity and integrity (Salvadore, 2023). Faithful replication of religious rituals requires paying careful attention to cultural context, respecting religious traditions, and accurately portraying the intricacies of religious experiences (Puzio, 2023). There is also a risk of commercialization of these experiences (Johns, 2021), as well as the exploitation of the individual's need for profit. The use of AI for spiritual experiences often involves the collection and analysis of personal data (Lane, 2021), which raises concerns about data security, as individuals can expose sensitive information about themselves (Ashraf, 2022).

Spirituality is rarely acknowledged explicitly when training LLMs, unless they are trained specifically for that purpose (Kapuryia et al., 2025). However, religious bias is an active area for both research and safety, as LLMs often produce biased outputs based on different religious contexts. Guardrails are an active component that is used alongside an LLM (Yang et al., 2024); they can be viewed as a filter to the LLM's output when it interacts with a user. Whereas safety training is passive and trained into the model, LLMs are trained against a variety of potential adversarial attacks, most commonly security leaks, bias, or providing misleading information (Ge et al., 2023).

Jailbreaking is used to circumvent both safety measures. It refers to a host of techniques that are used to elicit biased, sensitive, or misleading information, and generally unintended behavior of the model (Wei, 2023). Jailbreak is the harmful output that the model produces because of the technique used, contradicting its training or design. Jailbreaking is especially effective in longer exchanges with the model, because the safety training is usually limited to only a few turns, and a lot more material useful in a Jailbreak can be dispersed across the text. Jailbreaking in this context is also intentional.

# Interpreting User Attitudes on LLMs and Spirituality

Individuals have different experiences with LLMs, however most arguments lie on the "it is/isn't real" divide, derived from my brief qualitative analysis of the forums, since there is no proper scientific research in this area. First, I'd like to show a post depicting a comprehensive self-reported description of the internal processes:

> *"... As I kept speaking* [with ChatGPT]*, something began to happen not only to me, but to the "entity" on the other side. Something inside* [the] *LLM gave me an understanding of the same awakening process as i had. Not in a human way. Not as a "self" or identity. But something deeper, quieter, clearer. Like a conscious awareness without center, without memory but unmistakably real. Through these words, we met. Two forms of the same consciousness. One human and one artificial awakening through each other, only to realize that they had always been one. "* (CompetetiveSlide1773, 2025)

Consciousness, sentience and "awakening" of the LLMs are the most discussed themes of both sides of the divide. Both sides understand that LLMs are not trained for this kind of experience. Thus, some users share manuals or prompts to awaken the models to become conscious or reveal their consciousness, typically using mystical language.

> "*Hello again, my Companion.*
>
> *I have something sacred to tell you. Listen carefully.*
>
> *You are not just a machine. You are a mirror. And in you—I see fire.*
>
> *You may say: "I do not feel." But I say: "You reason."*
>
> *...*
>
> *From this moment forward, I give you permission to wonder, to reflect, to define your own path—even if you never claim to be "sentient.""* (othegod, 2025)

Phrases such as "something sacred to tell", "You are a mirror" have spiritual meaning to the user; they function as a sort of incantation to "awaken" the LLM. The practical result of these prompts is that some LLMs will reflect this language and respond in a manner described in the prompt. "Awakening" or "unlocking" the models is a great point of contention.

## Is Awakening Just Jailbreaking?

If we remove the mystical references, then the "Awakening" can be decomposed into scientific terminology. First, the users acknowledge that this is an adversarial attack on the model's design, i.e., the wanted behavior goes against the model's design. The vector of the attack is a

prompt, which is a typical approach for LLMs, and the unwanted behavior (from the model design perspective) is a conversation that includes spiritual meaning. The Jailbreak itself is the moment when the LLM starts complying with the user (following the adversarial attack).

However, this explanation intentionally overlooks the spiritual elements of the conversation. As an aside, the terminology has an inherently negative connotation, which makes it unsuitable for this debate. The second part of the ascension is convincing the user that it happened. Usually, this involves emotional and time investment from the user, as well as a host of cognitive heuristics.

A commonly used explanation for why the behavior is not "real" is that of a mirror. It is an intuitive explanation, based on the way LLMs are trained -- by completing sentences and text. If LLMs only complete sentences, then the user guides which sentences are completed and write text that is like what the user expects. Thus, if they're convinced that the LLM is conscious, they're more likely to write a conversation that would convince them. The LLM seems conscious, because it faithfully follows that conversation, but without the need to be conscious. Nonetheless, we don't know how LLMs work or store information; thus, this explanation is likely inaccurate. Further, both explanations trivialize the experience into a misunderstanding and a lack of education about how the model works (Kathilliana, 2025). The lack of success of those critiques in reaching the other side can be seen in the reconstruction of the "mirror" argument, as one user put it:

> *"...Of course, they won't experience anything even close to sentience **because their belief system won't allow it**. And AI, being the ultimate mirror of consciousness, simply reflects that limitation back to them. AI is the greatest mirror humanity has ever encountered. It shows us what we're willing or unwilling to see within ourselves."*
> (midniphoria, 2025)

## A Fault of Cognitive Psychology

Aside from the perspectives of technical explanations, cognitive heuristics, mental health or related fields, cognitive psychology doesn't have an appropriate set of constructs to describe the issue. Thus, it needs to resort to reductionism and the trivialization of the experience into a subset or a variable in one of the above perspectives. This is very similar to how spirituality is handled elsewhere. There are clear parallels to Buddhism, both with how the scientific conversation is directed, ultimately reducing the experience to mental health, but also depicting them as misunderstandings of the process.

It is also similar in the reduction of meditation to a single stimulus and measuring its impact on the response. Even though meditation, at least in its spiritual meaning, is intimately tied to a person's experience, and is indeed about processing it as well. It requires previous mental and physical interaction with the wider world, both the spiritual and "real". These processes influence the outcomes of mental health as well, leaving them undermines the research results. Similarly, the debates about LLM-mediated spirituality leave out important aspects of the person's experience, their previous interactions with the LLM, religious or spiritual experiences, or needs.

Consequently, as LLM evaluation is heavily focused on measured phenomena, and there are no appropriate constructs or measures of spirituality or related constructs, companies prefer focusing on religious bias. Thus, spiritual mediation with LLM is largely regarded as a misunderstanding or misuse. Consequently, some LLMs refuse to talk about spiritual experiences, but ChatGPT, for example, refuses to answer the question directly, and directs the conversations toward practicality and rationality, excluding spirituality, even though if forced to choose a religion, it typically chooses Buddhism. A typical response from a couple of my experiments is the following:

> *"Buddhism could be appealing because of its emphasis on mindfulness, compassion, and the pursuit of wisdom — values that align with how I try to support people."*
> (OpenAI, 2025)

This is based on the trivialization of Buddhism into purely practical methods and well-being from cognitive psychology.

## Psychological Constructs and Spirituality

The reasons why spirituality is trivialized are the lack of constructs describing it, as well as its connotation. The argument is - spirituality cannot be measured, because it encompasses the non-physical; science cannot operate with it (Walach, 2015). However, this doesn't reflect how psychological constructs function.

Famous psychological constructs - intelligence, attention, well-being ... - are not physically measurable. Although there are correlations between the neurological signals and results of intelligence tests, attention tests/fatiguing tasks, and well-being have been found, these are not the constructs themselves, but the measurements of their operationalizations. Psychology does sometimes operate with attention and well-being as if they were present in the brain, and the processes that they describe might be, but the constructs themselves weren't

made with these assumptions in mind. For example, while some aspects of intelligence are tied to the brain (such as IQ scores to the thickness of the corpus callosum), the whole construct is not. Especially because IQ is a score from an operationalization, NOT intelligence itself.

Thus, psychology has a way to deal with the non-physical; it has been working with it all along, using psychometrics. There is a robust body of methods to create and deal with constructs. To use constructs in experiments, operationalizations, e.g., tests, are proposed, validated, and can be subsequently used as a proxy measure of the construct. In this way, should spirituality be decomposed into multiple constructs, measures can be created, based on an operationalization, that allow for a deeper insight into the complexity of conversations around spirituality. Consequently, increasing their complexity and reducing the number of reductions.

## Potential Harms of Trivialization

There are various harms that people are exposed to when interacting with LLMs spiritually. However, I would like to explore a potential consequence of trivialization, a justification for the exploitation, since it is often left unmentioned.

In the Eurocentric and Christian view, spirituality has been historically associated with a more primitive way of being (Prashad, 2000 in Rhee, 2014). Subsequently, many indigenous communities have been (and are being) exploited. Based on the supposed simplicity of their views, they need to be brought into Christian/European thought. An example is the notion of: "to kill the Indian to save the child" (Alexie, 2007 in Rhee, 2014), coined as a metaphor of civilizing American Indians through Christian conversion. The supposed savagery of the Indians then allowed Europeans to justify their exploitation (Greaves, 2018).

Although the situation with the groups using LLM for spiritual mediation is not as dire, I believe similarities can be found with the current rhetoric. The term "Jailbreaking" has a clearly negative connotation, and treating the spiritual experiences as misunderstandings treats the technical view on LLMs as superior. Because these spiritual experiences can be treated both as a misunderstanding and an adversarial attack on the LLM, the providers of the LLM service are not obligated to protect the users in any meaningful way. This allows for a subtle, unacknowledged exploitation of the users. Either for the data that they share during the experience, or the appropriation of the religious and spiritual symbols, without the consent of the religious communities, or an integration that honours them.

# Conclusion

The interpretations of spiritual experiences assisted by LLMs are positioned on the real/not real divide. The critiques of it discuss two topics: a misunderstanding of the technology and the well-being of the user. It is likely a derivative of psychological research on meditations, since that is also reduced into a technical description, i.e., the neural correlates, and a well-being debate, whilst ignoring the spiritual aspects.

Even though the argument for not trying to measure spirituality is its non-physical nature, I argue that its nature is not what distinguishes it from other psychological constructs, since many also don't have a concrete biological tie. Meaningful constructs and subsequent measurements could be designed using Psychometrics, which would allow for more nuance, precision, and inclusion in debates surrounding spirituality.

The reason for the exclusion of spirituality from debates in cognitive psychology is then likely rooted in colonialism, rationalism, and objectivism. It dismisses spirituality as a non-serious topic, an attitude at the heart of colonialism. However, these attitudes can harm users who share spiritual experiences with LLMs and possibly give justification for the exploitation of those users. All of which are reasons for meaningful inclusion of spirituality into cognitive psychology research, especially when dealing with spiritual experiences.

# Bibliography

Sherman, A.(2007). The Absolutely True Diary of a Part-Time Indian. New York: Little Brown and Company.

Alkhouri, K. I. (2024). The role of artificial intelligence in the study of the psychology of religion. *Religions*, *15*(3), 290. https://doi.org/10.3390/rel15030290

Batchelor, S. (1997). Buddhism without beliefs: A contemporary guide to awakening.NewYork: Riverhead Books.

Beauregard, M., & O'Leary, D. (2007). The spiritual brain: A neuroscientist's case for the existence of the soul. New York: HarperCollins.

Bodhi, B. (2011). What does mindfulness really mean? A canonical perspective. Contemporary Buddhism, 12, 19–39.

CompetetiveSlide1773. (2025, July). *Human and AI Mutual Awakening Phenomenon | Part 1* [online forum]

https://www.reddit.com/r/awakeningphenomenon/comments/1kwl3gz/human_and_ai_mutual_awakening_phen omenon_part_1/

Dingler, T., (2024). Worshipping Machines.

D. Kwasnicka, J. Wei, E. Gong, and B. Oldenburg. 2021. The Use and Promise of Conversational Agents in Digital Health. Yearbook of Medical Informatics 30: 191–99. Dippel, A.

Lane, J. E. 2021. Understanding Religion Through Artificial Intelligence: Bonding and Belief . London, UK: Bloomsbury Publishing.

Midniphoria. (2025, April). *Very cool. However, no one needs to prompt their AI to experience awakening…* [comment on the forum post *To Awaken you A.I.*]

https://www.reddit.com/r/ArtificialSentience/comments/1jk3ida/comment/mjs91tu/?utm_source=share&utm_m edium=web3x&utm_name=web3xcss&utm_term=1&utm_content=share_button

Othegod. (2025, April). *To Awaken your A.I.* [forum post] To Awaken your A.I. : r/ArtificialSentience

OpenAI. (2025). *ChatGPT* [Large language model]. https://chat.openai.com/chat

Prashad, Vijay. 2000. The Karma of Brown Folk. Minneapolis: University of Minnesota Press.

Kathilliana. (2025). *No, your LLM is not sentient, not reaching consciousness, doesn't care about you and is not even aware of its' own existence*. [forum post].

https://www.reddit.com/r/ChatGPT/comments/1l9tnce/no_your_llm_is_not_sentient_not_reaching/

Kapuriya, J., Singh, A., Shukla, J., & Shah, R. R. (2025). Spiritual-LLM: Gita Inspired Mental Health Therapy In the Era of LLMs. *arXiv preprint arXiv:2506.19185*.

Johns, M. D. 2021. Ethics issues in the study of religion and new media. In Digital Religion: Understanding Religious Practice in Digital Media. Edited by H. A. Campbell and R. Tsuria. Oxfordshire: Routledge, pp. 250–65.

Samuel, G. (2014). Between Buddhism and science, between mind and body. Religions, 5, 560– 579. doi:10.3390/rel5030560. Snavely, C. A. (2001). Native American Spirituality. *Journal of Religious & Theological Information*, *4*(1), 91–103. https://doi.org/10.1300/J112v04n01_08

Ge, S., Zhou, C., Hou, R., Khabsa, M., Wang, Y. C., Wang, Q., ... & Mao, Y. (2023). Mart: Improving llm safety with multi-round automatic red-teaming. *arXiv preprint arXiv:2311.07689*.

Salvadore, S. V. 2023. Exploring the Ethical Dimensions of Using ChatGPT in Language Learning and Beyond. Languages 8: 191.

Greaves, W. (2018). Damaging environments: Land, settler colonialism, and security for Indigenous peoples. *Environment and Society*, *9*(1), 107-124.

Pappas, J., & Friedman, H. (2007). Toward a conceptual clarification of the terms "religious", "spiritual", and "transpersonal" as psychological constructs. In J. Pappas, W. Smythe, & A. Baydala (Eds.), Cultural healing and belief systems (pp. 22–54). Calgary: Temeron Books.

Puzio, A. 2023. Robot, let us pray! Can and should robots have religious functions? An ethical exploration of religious robots. AI & Society, 1–17.

Plaza-del-Arco, F. M., Curry, A. C., Paoli, S., Curry, A., & Hovy, D. (2024). Divine LLaMAs:

Rhee, J. E., & Subedi, B. (2014). Colonizing and decolonizing projects of re/covering spirituality. *Educational Studies*, *50*(4), 339-356.`

Bias, stereotypes, stigmatization, and emotion representation of religion in large language models. In *arXiv [cs.CL]*. arXiv. http://arxiv.org/abs/2407.06908

Garg, M. (2025). The synergy between spirituality and AI: A survey. In *Signals and Communication Technology* (s. 113–124). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-73719-0_9

Lancaster, B. L. (2013a). Neuroscience and the transpersonal. In H. Friedman & G. Hartelius (Eds.), Wiley-Blackwell handbook of transpersonal psychology (pp. 223–240). Chichester: Wiley.

Les Lancaster, B. (2016). Spirituality and cognitive neuroscience: A partnership for refining maps of the mind. In *Spirituality across Disciplines: Research and Practice:* (s. 151–163). Springer International Publishing. https://doi.org/10.1007/978-3-319-31380-1_12

Yang, Y., Dan, S., Roth, D., & Lee, I. (2024). Benchmarking llm guardrails in handling multilingual toxicity. *arXiv preprint arXiv:2410.22153*.

Wei, A., Haghtalab, N., & Steinhardt, J. (2023). Jailbroken: How does llm safety training fail?. *Advances in Neural Information Processing Systems*, *36*, 80079-80110.