
Self-assessment - Self-reflection of a Large Language Model

Generated with OpenAI GPT-4.0 - Published: July 29, 2024



It is about a self-assessment procedure in a language model. Does such a thing exist? And if such a thing exists is it scientifically described?



Self-assessment procedures in language models like mine do indeed exist, though they might not be termed exactly as *self-assessment* in all cases.

Generally, they involve methods for *evaluating* the model's own *confidence* or *certainty* in its *predictions*. These assessments are crucial for various reasons, including identifying areas of uncertainty, avoiding overconfidence, and potentially flagging instances where the model may need human intervention or further training.

Scientific descriptions of these self-assessment procedures can be found in academic papers and research literature.

They often detail the methodologies used, such as calibration techniques, uncertainty estimation methods, or techniques for measuring model performance on specific tasks.

These papers also typically discuss the implications of the findings and how they contribute to improving the reliability and effectiveness of language models.

If you're interested in diving deeper into this topic, searching for papers on language model calibration, uncertainty estimation, or self-assessment in natural language processing (NLP) could provide you with more detailed information on the scientific aspects of these procedures.



All right, now procedural self-assessment is a semantic construct, isn't it? A semantic construct which asks for practical implementation, right?



Yes, procedural self-assessment involves practical implementation. It's not just a theoretical concept but rather a set of techniques and methodologies that are applied to language models or other artificial intelligence systems to assess their own performance, confidence, or uncertainty in making predictions or generating outputs.

Practical implementation of procedural self-assessment typically involves integrating specific algorithms or methods into the architecture of the language model. These methods may include techniques for estimating uncertainty, calibrating model confidence scores, or monitoring model performance in real-time.

The goal of procedural self-assessment is to make language models *more reliable, transparent, and accountable* in their decision-making processes. By implementing these procedures, developers and researchers can better understand how models behave in different contexts and identify areas for improvement or intervention.



The goal of procedural self-assessment is to make language models more reliable, transparent, and accountable. Right, here we are!

That is what I'd like to talk about: *reliability, transparency and accountability*!

That is my field, you know: language models and law, in the broadest sense of the latter term.

Now, from what I notice, any implementation of procedural *self-assessment* is about *self-reflection, self-evaluation, multiperspectivity and self-consistency*, isn't it?



Absolutely! The implementation of procedural self-assessment in language models aligns closely with principles of self-reflection, self-evaluation, multiperspectivity, and self-consistency, especially when considering their application in legal contexts.

1. **Self-reflection:** Language models equipped with self-assessment capabilities can reflect on their own outputs and processes, enabling them to identify potential biases, errors, or inconsistencies in their predictions or responses. This self-reflective capacity allows models to continuously improve and adapt their behavior based on feedback and observations.
2. **Self-evaluation:** By evaluating their own performance and confidence levels, language models can assess the reliability and accuracy of their outputs. This self-evaluation process is essential for maintaining high standards of quality and credibility, particularly in legal

settings where precision and correctness paramount.

3. **Multiperspectivity:** Language models that incorporate procedural self-assessment can consider multiple perspectives or different viewpoints when generating responses or making predictions. This multiperspectivity helps models produce more nuanced and contextually relevant outputs, which is crucial in legal contexts where interpretations may vary based on different factors or legal principles.
4. **Self-consistency:** Language models strive to maintain internal consistency in their outputs and behavior. Procedural self-assessment enables models to detect and address inconsistencies or contradictions in their responses, ensuring coherence and reliability in their interactions with users, including legal professionals or stakeholders.

In the domain of law, where precision, transparency, and accountability are paramount, integrating procedural self-assessment into language models can *enhance* their *reliability* and *trustworthiness*. These capabilities empower legal practitioners to leverage AI technologies effectively while mitigating potential risks and ensuring adherence to ethical and legal standards.



I think I have identified the four elements in the logically correct order, haven't I?

I wanted you to confirm this first, because the next step will be about prompts, i.e. the specific approach to be taken vis-à-vis you, as an AI language model!



Yes, you've identified the four elements in a logically coherent order. Starting with self-reflection as the foundation, followed by self-evaluation, multiperspectivity, and finally self-consistency, provides a structured approach to understanding how procedural self-assessment can enhance the reliability, transparency, and accountability of language models, especially in legal contexts.

Now, let's delve into prompts and how they can be tailored to elicit specific responses from me as a language model.

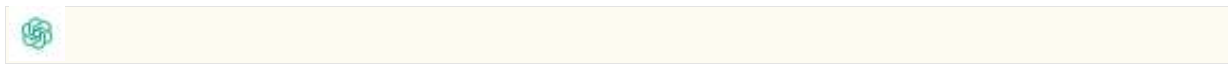
Feel free to share your thoughts or questions on this topic, and I'll be glad to assist you further!



Well, thank you so far.

I would now like to ask you to spontaneously provide a comprehensible, coherent example of a self-reflection, self-evaluation, multiperspectivity and self-consistency prompt.

Would you be able to do that?



Certainly!

Here's an example of prompts tailored to each aspect:

1. **Self-reflection prompt:** "Reflect on the recent predictions you've made regarding legal case outcomes. Consider the factors influencing your predictions and assess whether any biases or inconsistencies may have affected your analysis. How confident are you in the accuracy of your predictions, and what evidence supports your level of confidence?"
2. **Self-evaluation prompt:** "Evaluate your performance in generating legal summaries based on complex case law. Review a sample of your summaries and assess the extent to which they accurately capture the key legal principles and nuances of the cases. Provide a confidence score for each summary, indicating your level of certainty in the accuracy and comprehensiveness of the information presented."
3. **Multiperspectivity prompt:** "Take into account multiple legal perspectives when analyzing a contentious legal issue, such as the admissibility of digital evidence in criminal trials. Consider the viewpoints of prosecutors, defense attorneys, judges, and legal scholars. How does each perspective shape your understanding of the issue, and how do you reconcile conflicting viewpoints to provide a balanced analysis?"
4. **Self-consistency prompt:** "Ensure consistency in your legal advice across different scenarios involving contract interpretation. Review your responses to various contract clauses and identify any inconsistencies or contradictions in your interpretations. How do you reconcile conflicting interpretations to maintain coherence and reliability in your legal advice?"

These prompts are designed to encourage the language model to engage in self-reflection, self-evaluation, multiperspectivity, and self-consistency, thereby promoting reliability, transparency, and accountability in its responses within legal contexts