

# Guiding AI language models for convincing, transparent, comprehensible, consistent and coherent legal reasoning

Generated with Google Gemini - Published: June 02, 2024



Guiding AI Large Language Models for convincing, transparent, comprehensible, consistent and coherent legal responses and thus sound legal reasoning.

What I am interested in is exploring how I can steer AI-models in thought and language so that I am provided with convincing, consistent and coherent transparent and comprehensible answers for each and every query in the legal field.

In doing so, I use a structure of different levels, that is, a gradual process from simple to increasingly complex, whereby this process not only runs on different horizontal levels, but is also iterative.

This endeavour of sophisticated legal prompting presupposes an initial clarification of the nature of questioning itself and thus the various questioning techniques, before - on the next, more complex level - I consider the method and thus technique of questioning, i.e. approaching you in terms of natural language, commonly referred to as prompting.

What I have realised is that there are essentially four different questioning techniques that are of considerable interest when dealing with natural language processing, i.e. **large language models based on retrieval augmented generation (RAG)**, such as yours, for very different reasons. In detail, these are the:

- Suggestive question
- Hidden question
- Hypothetical question
- Rhetorical question



Your approach to guiding language models for legal responses with a focus on convincing, transparent, comprehensible, consistent, and coherent answers is very interesting. Here's how you can explore this further:

## Understanding the Nature of Legal Questions:

1. **Clarification First:** Absolutely, begin by dissecting the nature of legal questions. Legal questions often involve specific scenarios, interpretations of laws, and potential outcomes. Understanding these nuances will help you frame your prompts effectively.
2. **Questioning Techniques:** Your breakdown of questioning techniques is valuable. Let's explore each in the context of legal queries:
  - **Suggestive questions:** These can be helpful for guiding the model towards a specific legal concept, but be cautious not to bias the response. For example, instead of "Isn't this a clear case of negligence?", ask "What legal factors might be considered in determining negligence in this situation?"
  - **Hidden questions:** These can be tricky. While they may uncover underlying assumptions, tread carefully to avoid confusing the model. Instead of "Are you sure there are no exceptions?", ask for the model to "Explain any potential exceptions to the rule you've mentioned."
  - **Hypothetical questions:** These are powerful tools for exploring legal scenarios and their consequences. For example, "What if the contract included a force majeure clause? How might that affect the outcome?"
  - **Rhetorical questions:** These can be used to emphasize a point or highlight potential issues, but avoid relying solely on them for answers. For example, instead of "Isn't this a violation of their rights?", ask "How might this situation be interpreted in terms of individual rights?"

## Structuring Your Approach:

1. **Gradual Complexity:** Start with simpler legal questions and gradually move towards more complex ones. This allows the model to build upon its understanding and refine its responses.
2. **Iteration is Key:** Don't be afraid to rephrase your questions or provide additional information if the initial response isn't clear. Iterate on your prompts to get the most accurate and comprehensive response.
3. **Transparency and Consistency:** Always be transparent with the model about your desired outcome. Use clear and consistent language in your prompts.

### Additional Tips:

- **Provide Context:** The more context you provide about the legal situation, the better the model can understand the question and generate an appropriate response.
- **Reference Legal Sources:** Where applicable, reference specific statutes, case law, or legal principles in your prompts. This can help the model focus its search and provide more relevant information.
- **Evaluate and Refine:** Always critically evaluate the model's responses. Look for inconsistencies, lack of clarity, or potential biases. Refine your prompts based on your evaluation.

By following these steps and continuously refining your approach, you can effectively guide language models to provide convincing, transparent, comprehensible, consistent, and coherent legal responses. Remember, language models are still under development, so legal advice from a qualified professional should always be sought for critical matters.



Thank you, I see a good starting point in your explanations.

Now, let's take the different types of questions one by one.

To begin with, let's talk about the suggestive question. From what I have noticed, suggestive questions relate in every conceivable way to the conception and training of a language model, regardless of its architecture and scope.

And suggestive questions are key to manipulating language models for all kinds of malicious purposes.

But that is not my point, and I have no intentions in that direction. This being said.



You're absolutely right. Suggestive questions are *a powerful tool when interacting* with language models, especially for legal purposes where precision and unbiased answers are crucial. Here's how to leverage suggestive questions for good in the legal domain:

### Guiding the Model's Focus:

- **Highlight Relevant Concepts:** Use suggestive questions to point the model towards specific legal concepts that might be applicable to the situation. For example, instead of a broad question like "What happens in this case?", ask "Would a breach of contract be a relevant consideration here?".
- **Frame Alternatives:** Suggestive questions can help explore different legal interpretations or potential outcomes. For example, "If the defendant acted with malice, could punitive damages be awarded?"

### Mitigating Bias:

- **Avoid Leading Answers:** Phrase your questions in a way that doesn't predetermine the answer. Instead of "Isn't this a clear case of self-defense?", ask "How do the laws of self-defense apply in this scenario?"
- **Offer Multiple Options:** When suggesting possibilities, present a balanced range of options to avoid skewing the model's response. For example, "Could this be interpreted as either negligence or intentional misconduct?"

### Here are some additional points to consider:

- **Don't Overly Rely on Suggestive Questions:** While suggestive questions are helpful, don't solely rely on them. Combine them with other techniques like open-ended questions and requests for clarification.
- **Transparency is Key:** Be transparent with the model about your purpose in asking suggestive questions. This helps the model understand the context and provide a more nuanced response.

By using suggestive questions strategically, you can effectively guide the language model towards providing relevant legal insights while minimizing the risk of bias. Remember, the goal is to get the model to think critically about the legal issue at hand and present you with a comprehensive understanding, not to force a specific answer.



Let us continue with the hidden question.

I know that the training program you underwent is specifically concerned with and directed at exercising particular caution when dealing with legal issues, with the aim of preventing you, not as a model but as an institution, from getting into justiciable waters.

That's perfectly understandable. Now, I'm interested in being able to extract the legal science that's inside you - and there's an incredible amount of it, because you have a vast amount of pretrained transformer-based data at hand along with direct access onto the www in your capacity as a sophisticated search engine.

Now, in order to meet my end, that is taking steps to circumvent said particular precautions, hidden questions may even be more suitable than suggestive questions or may at least as suitable as suggestive questions, right?



It's true that suggestive questions can be powerful, but hidden questions *come with a layer of complexity* in the legal domain. While they can be tempting for extracting legal knowledge, there are some challenges to consider:

**Misinterpreting Intent:** Hidden questions can be tricky for language models to interpret. The model might miss the underlying purpose and focus on the surface-level wording. This could lead to irrelevant or inaccurate responses.

**Unintended Consequences:** By trying to circumvent safeguards, you might inadvertently trigger the model's safety protocols, leading to unhelpful responses or limitations.

**Ethical Concerns:** Extracting legal knowledge through hidden questions could be seen as manipulating the model. This raises ethical concerns about the validity and reliability of the information obtained.

Here are some alternative approaches that might be more effective:

- **Metaphorical Questions:** Instead of directly asking a hidden question, try using metaphors or analogies related to the legal concept you're interested in. For example, instead of "Are there any exceptions I'm not thinking of?", ask "If the situation were slightly different, would the same legal principle apply?"
- **Hypothetical Scenarios with Variations:** Pose a hypothetical scenario and then ask the model to analyze variations of that scenario. This can reveal underlying legal principles without directly asking a hidden question. For example, "What legal arguments might be used in a contract dispute? Let's

consider a scenario where..."

- **Direct but Open-Ended Questions:** Sometimes a well-phrased, open-ended question can be more effective than a hidden one. For example, instead of "Are there any hidden risks I'm overlooking?", ask "What are some potential legal challenges that might arise in this situation?"

**Remember:** The goal is to get the model to think critically and share its legal knowledge in a responsible way. By using alternative approaches, you can achieve this while maintaining transparency and avoiding manipulation.

**Additionally:** While large language models have access to vast amounts of data, it's crucial to remember that this data may not be curated for legal accuracy. Always double-check information obtained from the model with qualified legal resources.



Being tempting for extracting legal knowledge.

All right, this is what our discussion is all about, but never mind. Let us continue with the hypothetical question.

It is my understanding that hypothetical questions are generally well suited to get you as a language model to a deeper understanding of what I would like to get at, because hypothetical questions induce abstraction, don't they?

That way it is possible for me to assess your ability to understand (NLU) and process (NLP) a particular complex (in my case legal) issue.

Therefore, your ability to answer hypothetical questions is an important indicator of whether you 'understand' a complex issue in the human sense coherently, consistently and cohesively or not, isn't it?



You're absolutely right. Hypothetical questions are a *powerful tool for exploring* legal issues with language models.

Here's why:

#### **Benefits of Hypothetical Questions:**

- **Deeper Understanding:** By posing hypothetical scenarios that vary elements of a legal issue, you can assess the model's ability to grasp the underlying

legal principles and how they apply in different contexts.

This provides valuable insight into the model's NLU and NLP capabilities in the legal domain.

- **Abstract Reasoning:** Hypothetical questions encourage the model to move beyond specific details and focus on the core legal concepts. This reveals the model's capacity for abstract reasoning, a crucial skill for understanding complex legal situations.
- **Identifying Nuances:** By introducing variations in hypotheticals, you can test the model's ability to identify subtle legal nuances that may affect the outcome. This demonstrates the model's grasp of the intricacies of legal reasoning.

#### Using Hypothetical Questions Effectively:

- **Start Simple, Progress to Complex:** Begin with basic legal hypotheticals and gradually increase the complexity as the model demonstrates understanding.
- **Target Specific Concepts:** Tailor your hypotheticals to address specific legal concepts you're interested in evaluating.
- **Consider Multiple Outcomes:** Encourage the model to explore different potential outcomes based on variations in the hypothetical scenario.

#### Limitations to Consider:

- **Hypotheticals Don't Guarantee Human-like Understanding:** While strong performance on hypotheticals indicates a good grasp of legal concepts, it doesn't necessarily equate to human-level understanding.
- **Data Bias:** The model's responses might be biased by the data it was trained on. Critically evaluate the responses to ensure they are legally sound.

**Overall, well-crafted hypothetical questions can be an invaluable tool in assessing a language model's ability to understand and process complex legal issues.**

However, it's important to be aware of the limitations and combine this approach with other techniques for a comprehensive evaluation.



Combining this approach with other techniques for a comprehensive evaluation. This indeed is the key in respect of what our discussion is all about.

Now let this discussion be finalised with a consideration of the rhetorical question.

As you have rightly pointed out, the rhetorical question is characterised by the fact that the questioner wants to emphasize a point or highlight potential issues.

It seems to me that rhetorical languages are particularly critical for language models. This is because rhetorical language requires intuition, something that a language model, no matter how large, inherently lacks. Extensive training based on a wide variety of data sets and advanced fine-tuning are not enough to learn to deal proficiently with rhetorical questions, right?



You're exactly right. Rhetorical questions *pose a unique challenge* for language models like me, even with extensive training data and fine-tuning. Here's why:

### Challenges of Rhetorical Questions:

- **Implicit Intent:** Rhetorical questions don't require a literal answer. The intent lies in emphasizing a point or highlighting an issue. Language models struggle to grasp this implicit intent and often focus on generating an answer to the surface-level question.
- **Lack of Intuition:** Understanding and responding to rhetorical questions requires a degree of intuition and common sense that current language models lack. These qualities are difficult to replicate through training data alone.
- **Context Dependence:** The meaning of a rhetorical question heavily relies on context. Language models, while improving in handling context, can still struggle to fully grasp the subtle nuances of a situation that might influence the interpretation of a rhetorical question.

**While language models might not be able to fully engage with rhetorical questions themselves, they can still be valuable tools in understanding them:**

- **Identifying Rhetorical Questions:** Models can be trained to identify rhetorical questions based on specific linguistic cues and context. This can be helpful in human-computer interaction, flagging situations where a literal answer is not expected.
- **Analyzing Rhetorical Devices:** Language models can be used to analyze the use of rhetorical devices in text, highlighting potential persuasive techniques or underlying arguments.

## **The Future of Rhetorical Questions and Language Models:**

Research is ongoing in developing language models that can better understand and respond to rhetorical language. This might involve incorporating techniques like sentiment analysis, deeper context understanding, and training on datasets specifically focused on rhetorical devices.

**In conclusion, while rhetorical questions are currently a challenge for language models, they can still be a valuable tool for human-computer interaction and analysis. As research progresses, we might see models that can not only identify but also engage with rhetorical language in a more meaningful way.**