# Introduction to Prompt Engineering with Large Language Models

Mobin Kheibary
Faculty of Electrical and Computer
Engineering
Urmia University
Urmia, Iran
Mobin.kh15@gmail.com

*Abstract*—**Large Language Models (LLMs), utilizing transformer architectures, have demonstrated remarkable capabilities in comprehending and generating text that closely mirrors human-like expressions. Their utility can be greatly enhanced with an engineering approach called Prompt Engineering. This technique manipulates the model's prompts to draw out more desired responses, leading to more interactive, intuitive, and efficient interactions with these models. This paper serves as an introduction to the concept of prompt engineering, discussing its necessity and shedding light on numerous techniques to effectively engineer prompts for LLMs, such as ChatGPT. The focus is on the implementation of these techniques, their advantages, and their potential implications for the future of AI interactions.**

*Keywords*—*Prompt Engineering, Large Language Models, Natural Language Processing, Human-AI Interaction, ChatGPT, Deep Learning.*

## I. INTRODUCTION (*HEADING 1*)

The advent of transformer-based models, such as GPT-3 and its derivatives, have revolutionized the AI landscape by demonstrating unprecedented capabilities in generating human-like text. These transformer models utilize a network architecture that is particularly conducive to modeling long-range dependencies in sequence data, forming the foundation of LLMs. Despite their potential, these models have not been fully exploited yet, and one way to do this is through Prompt Engineering.

Prompt Engineering is the practice of crafting and manipulating the input prompt to elicit more accurate and meaningful responses from LLMs. A well-crafted prompt can be the difference between an AI system understanding and appropriately responding to a user's request, and a system delivering incorrect or even nonsensical responses. This paper elucidates different techniques of prompt engineering, how they can be implemented, and their future implications for the field of human-AI interaction.

## II. PROMPT ENGINEERING TECHNIQUES

The art of prompt engineering has evolved with the primary aim of optimizing interactions with LLMs. Each technique caters to specific needs and is best suited for certain scenarios. They provide a structured approach to formulating prompts that guide the AI to produce the most appropriate response. In the realm of Natural Language Processing, effective dialogue with LLMs can be achieved by the artful application of prompts. Good prompts take advantage of the model's training data, teasing out the embedded knowledge and making it explicit. Prompt engineering helps craft such tailored cues to optimize the interaction with LLMs, creating an intricate mesh between the prompts and the LLM's knowledge base.

Prompts can be thought of as the bridge that connects the user's intent to the model's understanding. Hence, the quality of the bridge, i.e., the prompt, plays a crucial role in effective communication. The techniques described in this section enhance the quality of this bridge, resulting in responses that closely match the user's requirements.

### A. Prompt Priming

Prompt Priming involves providing some contextual information to the model before posing the actual question. This priming helps steer the model's responses in a particular direction, often leading to more specific and relevant responses. For instance, giving a brief background of the topic "Climate Change" before asking specific questions can result in responses that are more in line with the query.

Prompt priming can also be thought of as setting the stage before presenting the main act. By defining the context initially, the AI system is guided towards a specific theme, enhancing the chances of receiving a more contextually accurate response. However, the trick lies in providing just the right amount of context – too much might lead to biased responses, while too little may result in vagueness.

Moreover, priming prompts work best when dealing with complex or multi-layered topics. For simple or direct queries, these prompts may overcomplicate the process, leading to inefficiency. Therefore, understanding the situation and correctly choosing when to use priming prompts is as important as the prompt itself.

Consider a situation where we want to ask the LLM about the impact of climate change on agriculture. A priming prompt might look like this: "Climate change has far-reaching impacts on various sectors, one of which is agriculture. With rising temperatures and unpredictable rainfall patterns, how is climate change impacting the agricultural sector?"

### B. Shot Prompting

Shot Prompting is a technique that manipulates certain parameters such as 'temperature' and 'max tokens' to control the length and randomness of the model's response. A low 'temperature' value, for example, would generate a more deterministic and focused response, whereas a higher value would produce more diverse and creative output.

Shot prompting, especially when manipulating the 'temperature' and 'max tokens' parameters, provides control over the exploration-exploitation trade-off in the model's response. Lower temperature values tend to produce a narrow, focused output that stays close to the training data, whereas higher values lead to diverse, exploratory responses.

Another key aspect of shot prompting is controlling the length of the response with 'max tokens'. This can be particularly useful when seeking brief responses or limiting

verbose explanations. However, it must be used judiciously, as limiting the output excessively might result in incomplete or nonsensical responses.

Suppose we want a brief and specific answer to a question about a historical event. The 'max tokens' parameter can be set to a small value to ensure a short response. A shot prompt with a low temperature could be: "Briefly, what was the significance of the Battle of Waterloo?"

### C. Chain of Thought Prompting

This technique involves using a series of related prompts to guide the model's responses along a desired path, akin to steering a conversation by asking connected questions. This allows for more intricate and complex dialogues, leading to a more engaging conversation with the AI model.

Chain of thought prompting is akin to having a conversation with the model. Each prompt is a link that connects to the next, creating a chain of coherent thoughts. It gives the user greater control over the conversation's trajectory and can be particularly useful when dealing with complex or abstract topics.

This technique also improves the interpretability of the model's responses by keeping the conversation grounded and related. However, chain of thought prompting requires careful management of the conversation's context, as losing track can lead to irrelevant responses.

If we're discussing a complex topic like the philosophical implications of artificial intelligence, we can use a series of related prompts. For example:

"What are the primary ethical considerations surrounding artificial intelligence?"

"How do these ethical considerations impact the development and deployment of AI technologies?"

"What measures can be taken to mitigate these ethical concerns?"

### D. Tabular Format Prompting

Presenting information in a tabular format can aid in eliciting responses that require comparisons or relative assessments. For example, presenting features of different smartphones in a table before asking the AI to recommend the best one, allows the model to process structured data effectively.

Tabular format prompting works effectively when dealing with structured information or comparison-based queries. The structured format can help the model grasp the correlation between different data points better, enabling it to generate more informed responses.

However, the challenge with tabular format prompting lies in structuring the data effectively. The model's understanding is directly related to how well the data is structured in the table. Incorrect or confusing table structures may lead to incorrect interpretations by the model.

Instead of asking "What are the differences between an apple and an orange?", you could structure the prompt like this:

"Compare an apple and an orange in terms of color and taste:

Color: Apples are usually red, while oranges are, well, orange.

Taste: Apples are sweet and sometimes slightly tart, oranges are usually a balance of sweet and sour."

The model would then generate a comparison based on these provided details.

### E. Ask Before Answering Prompting

In this approach, the model is prompted to ask clarifying questions before providing an answer. This technique encourages the AI model to clarify ambiguous queries, leading to more accurate and nuanced responses, and closely mirrors how human communication often occurs.

By prompting the model to seek clarifications, we encourage it to explore the query more deeply. This strategy can often unearth underlying complexities in the query, leading to a more fitting response. This technique closely emulates the human practice of seeking clarifications when faced with ambiguous queries.

However, overuse of this technique can lead to unnecessary back-and-forths that might frustrate the user. Therefore, it's crucial to strike a balance between seeking clarifications and providing direct responses.

If a user asks the model, "What's the best programming language?" The model, using the Ask Before Answering technique, might respond with: "The choice of a programming language can depend on many factors. Could you please specify if you're interested in web development, data science, game development, or another area?"

### F. Perspective Prompting

This technique involves posing the prompt from a specific perspective, such as a particular profession, to guide the model's responses. For example, asking a question from a doctor's perspective might yield different responses compared to asking from a historian's perspective.

Perspective prompting is a powerful tool that can lend a unique flavor to the model's responses. By adopting a particular perspective, the model can generate responses that mirror the thinking of that perspective. This can be particularly useful in scenarios requiring subject-specific insights.

However, this technique's effectiveness largely depends on how well the model has been trained on the chosen perspective's corpus. Inadequate or biased training data can lead to less than satisfactory responses.

Suppose we want to understand how a botanist would see a forest. The prompt could be: "As a botanist, how would you describe the importance of a healthy forest ecosystem?"

### G. Constructive Critic Prompting

In this technique, the model is guided to critique or evaluate a particular idea or concept. This can be useful for brainstorming sessions or idea generation where the AI can provide novel insights or point out potential flaws in a plan.

This technique can be used as a creative or critical thinking tool. By prompting the model to critique or evaluate a concept

or idea, we can gain new insights or identify potential pitfalls. This technique essentially leverages the model's vast knowledge base to think creatively or critically.

However, the effectiveness of this technique largely depends on the quality of the critique provided by the model. This in turn depends on the quality of the training data and how well the model has been trained in critical thinking tasks.

If we're discussing a new business idea, we could ask the model to critique it. For example: "Consider a startup idea focusing on delivering homemade meals. What potential challenges and opportunities do you foresee with this idea?"

## III. EXTERNAL TOOLS AND LIBRARIES

A plethora of external tools, Integrated Development Environments (IDEs), and libraries can greatly assist in effective prompt engineering. These tools provide the environment for crafting, testing, and refining prompts, and can highlight potential issues in real-time. Examples include Python libraries for text processing, data analytics tools for assessing model responses, and various IDEs that support AI development.

The significance of external tools and libraries extends beyond just crafting and refining prompts. These tools can also be used to analyze the effectiveness of different prompts, providing valuable insights that can help improve future prompts. Moreover, these tools can also assist in handling more complex tasks, like managing multiple simultaneous interactions or handling asynchronous responses.

Additionally, tools and libraries can assist in tracking the model's performance over time. This can be used to identify trends or anomalies, which can then be used to improve the model's training or adjust the prompting strategies.

## IV. FUTURE OF PROMPT ENGINEERING

As AI continues to evolve, the scope of prompt engineering is projected to grow in tandem. As models increase in complexity and understanding, so too will the prompts that we can feed into these models. Improved techniques will enhance the capabilities of LLMs, enabling them to understand more nuanced prompts and provide better responses. However, this progression might encounter challenges such as increased complexity of prompts and the need for a better understanding of contextual information.

The growth of prompt engineering is closely tied to advancements in AI and Natural Language Processing. As these fields evolve, we can expect to see newer and more advanced prompting techniques. However, this advancement will also bring along challenges, like managing the complexity of the prompts or ensuring that the prompts don't lead the model into generating inappropriate or harmful responses.

As the frontier of AI continues to push forward, one can expect the challenges in prompt engineering to keep pace. We are likely to witness the evolution of automated prompt generation and refinement tools, greater integration of AI systems with external databases for dynamic prompting, and the development of regulatory frameworks to guide prompt creation.

## V. CONCLUSION

Prompt engineering presents an effective way to optimize interactions with LLMs. It introduces a variety of techniques to enhance the capabilities of these models, allowing them to generate more accurate, detailed, and contextually relevant responses. As AI progresses and the demand for more human-like interaction grows, prompt engineering will undoubtedly play a key role in the future of human-AI interaction. The rapid development of LLMs and their widespread applications makes the field of prompt engineering an exciting area of research, promising significant advancements in the way we interact with AI systems.

In summary, prompt engineering offers a potent toolset to refine and customize interactions with Large Language Models. By carefully designing prompts with techniques such as priming, shot prompting, and chain of thought prompting, one can guide these AI systems towards more accurate, relevant, and context-specific responses. Tabular format prompting and the use of perspective-based or critique prompts further push the boundaries of what these models can accomplish in terms of complex tasks and assessments. The usefulness of these techniques, coupled with the incorporation of external tools and libraries, promises a more versatile, interactive, and intuitive Human-AI interaction paradigm.

As we move into the future, the evolution of these techniques will likely shape the development and application of LLMs. Just as LLMs have revolutionized text generation and Natural Language Processing, advancements in prompt engineering can lead to the creation of more refined AI interfaces, enabling more nuanced interactions and more precise outputs. As such, understanding and harnessing these techniques becomes pivotal for those looking to harness the power of AI in the most effective way. By bringing together the theory and practice of prompt engineering, we hope this work contributes to a deeper understanding of the field, promoting the exploration and adoption of these techniques in various AI-related pursuits.

### REFERENCES

[1]  T. B. Brown, B. Mann, N. Ryder, et al., "Language Models are Few-Shot Learners," in Proc. of the 37th International Conference on Machine Learning, vol. 119, pp. 1410-1423, July 2020.

[2]  V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," in Proc. of the 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing, pp. 1-6, December 2019.

[3]  D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2017.

[4]  A. Vaswani, N. Shazeer, N. Parmar, et al., "Attention is all you need," in Advances in Neural Information Processing Systems, vol. 30, pp. 5998-6008, December 2017.

[5]  S. Ruder, "An overview of gradient descent optimization algorithms," arXiv preprint arXiv:1609.04747, 2016.