

Prompting Large Language Models With the Socratic Method

Edward Y. Chang

Computer Science, Stanford University

echang@cs.stanford.edu

Abstract—This paper presents a systematic approach to using the Socratic method in developing prompt templates that effectively interact with large language models, including GPT-3. Various methods are examined, and those that yield precise answers and justifications while fostering creativity and imagination to enhance creative writing are identified. Techniques such as *definition*, *elenchus*, *dialectic*, *maieutics*, *generalization*, and *counterfactual reasoning* are discussed for their application in engineering prompt templates and their connections to inductive, deductive, and abductive reasoning. Through examples, the effectiveness of these dialogue and reasoning methods is demonstrated. An interesting observation is made that when the task’s goal and user intent are conveyed to GPT-3 via ChatGPT before the start of a dialogue, the large language model seems to connect to the external context expressed in the intent and perform more effectively.

Index Terms—large language model, natural language processing, prompting, the Socratic method.

I. INTRODUCTION

Prompting is a technique used to guide the output generation of a pre-trained language model such as GPT-3 [2]. This is achieved by providing input in the form of a question or template, which helps to generate specific responses such as Q&A, document summarization, and translations. The advent of ChatGPT [11, 23, 41] has revolutionized the field of NLP by demonstrating the potential of using large pre-trained language models with prompting. Despite this progress, there is still room for improvement in current prompting strategies and techniques, especially for specific target applications. In this study, we investigate the Socratic method [42, 40] to identify and evaluate potential prompting strategies, and use the findings to design effective prompt templates.

Traditional NLP tasks involve various sub-tasks, such as named entity recognition, dependency parsing, coreference resolution [8], semantic parsing [25, 9], and more, to comprehend the meaning of a sentence. By utilizing prompt templates with large language models (LLMs), these sub-tasks can be delegated to the LLM, freeing the template to focus specifically on dialogue design. In this regard, the Socratic method [31] holds significant relevance, as it is well-known for using questioning (prompting) as a means of promoting critical thinking and delving into complex concepts [10].

The Socratic method has a long history of being regarded as the basis of critical thinking. However, some recent studies have cast doubt on its effectiveness in practice. In his paper “Socratic Irony and Argumentation,” Airaksinen [1] criticizes the method for its rigidly defined roles of teacher and student,

which can lead to fear of not meeting the teacher’s expectations and reluctance to participate. Similarly, Stoddard’s “The Use of Socratic Questioning in Clinical Teaching” [35] highlights the risk of the method being misused in a manner that lacks psychological safety for students. Fortunately, when using the Socratic method in a dialogue with an LLM, the absence of emotions and sarcasm, as well as the option to deactivate the model, can alleviate many of the problems associated with human interaction.

This study starts by presenting an overview of the Socratic method’s strategies and techniques. To begin, we list ten widely referenced methods [3] under the Socratic method umbrella and use hypothesis elimination to identify the most relevant ones for our goal of prompt-template development. The selected methods are definition, hypothesis elimination, elenchus, dialectic, maieutics, generalization, and induction. Furthermore, we add to the list counterfactual reasoning, which is a concept in logic that involves considering what might have happened if a particular event had occurred differently. We then perform experiments using GPT-3 to test and evaluate these methods, and offer suggestions for incorporating these strategies and techniques into prompt templates.

In their work on “Critical Thinking: The Art of Socratic Questioning,” Paul and Elder identify three types of Socratic questioning: spontaneous, exploratory, and focused [27]. We will not discuss spontaneous questioning, as it is similar to casual conversation. Focused questioning (type 2), on the other hand, is geared towards gaining knowledge and truth, and methods such as *definition*, *elenchus* (cross-examination), *hypothesis elimination*, *dialectic*, and *generalization* hold great potential for developing effective prompting strategies and improving the response accuracy of a large language model (LLM). An interesting observation is that when the user intent is conveyed to GPT-3 during the task *definition* stage, before the start of a dialogue, the LLM seems to connect to the external context expressed in the intent and perform more effectively. (Table V provides an example of pre-dialogue warm-up. More examples are documented in [5].)

Additionally, exploratory thinking (type 3) can be supported through the *maieutics* (midwife) method, *induction*, and *counterfactual reasoning*, which can guide GPT-3 towards producing imaginative and creative writing. While many of the plot suggestions generated by GPT-3’s exploration may not be useful, a few unique recommendations in response to a “what if” query can stimulate the writer’s imagination and lead to remarkable results. When applied effectively, these methods

can turn an LLM into a writer’s muse, providing inspiration and guiding the creative process [36].

The main contributions of this paper are as follows:

- An overview of the Socratic method’s strategies, their evaluation, and selection of the most relevant ones for the development of effective prompt templates.
- An examination of how the definition, elenchus, hypothesis elimination, dialectic, and generalization methods can improve the output’s accuracy and conciseness through clarification and verification.
- An illustration of how maieutics, induction, and counterfactual reasoning can foster productive generalization and creativity.

The remainder of this paper is structured into five sections. Section II provides a review of related work on prompting methods in natural language processing. In Section III, we introduce the ten strategies and methods taught by Socrates and used in Plato’s “Dialogues.” From these, we select relevant methods along with counterfactual reasoning as our focus for developing prompting templates. Section IV details how we engineer these methods into our templates to improve output correctness and stimulate creative writing. In Section V, we present a pilot study. Finally, in Section VI, we present our concluding remarks.

II. RELATED WORK

The use of transformer architecture [37] and masked data for pre-training large language models (LLMs) in an unsupervised setting has become *the approach* in natural language processing [7, 19]. The method involves pre-training an LLM on a large text corpus, followed by fine-tuning for specific tasks.

Prompting is a recent innovation in the field, popularized by OpenAI, especially with the release of GPT-3 in 2020. Instead of fine-tuning the model for a specific task, the approach involves providing a specific input, or “prompt,” to guide the LLM’s output generation, resulting in greater flexibility and efficiency in generating a wide range of responses.

However, designing effective prompt templates remains a challenge [22], as it requires a deep understanding of the interplay between the LLM and the prompt. According to the survey paper [43], there are several factors that impact prompt template engineering, including the type of LLM used, manual vs automatic design, and static vs continuous prompts.

- Left-to-right vs masked LLMs. For tasks related to generation or tasks solved using a standard left-to-right language model [2], prefix prompts tend to perform better, as they align with the model’s left-to-right nature. For tasks solved using masked language models [7], cloze prompts are more suitable, as they closely match the pre-training task form.
- Manual vs automatic design. A prompt template should be tailored to the specific LLM. While manual design may be suitable in the initial flow-design phase, dependencies between the input and expected output, and their variations, should be mined automatically [15]. Automation can also help in paraphrasing the seed prompt to support various mined dependency patterns, but mistakes can occur [12].

- Discrete vs continuous prompts. Discrete prompts involve providing a fixed set of pre-determined input choices to an LLM. Continuous prompts, on the other hand, involve a dialogue or conversation between the model and the user, allowing for a more dynamic and interactive experience.

More advanced templates can be constructed by combining basic templates with techniques like ensemble methods [34]. This involves forming a committee of basic templates that ask the same question using different phrasing [13]. Most current prompt templates generate short outputs, such as class labels, or outputs with a length that can be predicted based on the task and input, like in the case of translation. However, for tasks that may generate longer or open-ended outputs, additional considerations may be necessary during the template engineering process.

One approach for generating longer outputs is explanation-based prompting, as proposed by the chain-of-thought method [39]. This method generates a sequence of explanations before inferring the answer. However, when dealing with simple math problems, this approach has an error rate of 47%. To address the inconsistency issues of explanation-based prompting, [16] formulates the problem as a satisfiability problem, which defers inference until a tree of explanations has been expanded abductively (explaining both truth and false branches) and recursively. However, using abductive reasoning alone is often considered weak, incoherent, and even nonexistent [14, 32]. To improve consistency, a recent work [38] extends the chain-of-thought approach by adding a diverse set of reasoning paths and performing majority voting among them. This method can be viewed as an ensemble method, but it does not alter the nature of abductive reasoning.

In contrast, the Socratic method aims to employ deductive, inductive, and abductive reasoning to ensure consistency and accuracy of inference. The Socratic method deals with all aspects of critical thinking, including definition clarification and cross-examination. This comprehensive approach to template engineering can lead to improved output quality and consistency.

The primary objective of this study is to design continuous prompts that enhance response quality and foster guided creativity in generative tasks, such as verifying information, evaluating source credibility, proposing alternatives, recommending plot ideas in creative writing, and generating task-specific surprises. Our approach involves investigating strategies and methods within the Socratic method, and selecting the most relevant ones for further exploration.

As discussed in Section I, Socratic questioning can be classified into three categories: spontaneous, exploratory, and focused [27]. When designing a prompt, it is important to consider the category and utilize the most suitable strategies and techniques to achieve the best results.

III. THE SOCRATIC METHOD

The Socratic method is a questioning technique used in teaching and philosophy to encourage critical thinking and self-discovery [40]. The method involves asking a series of questions to explore complex ideas and help individuals arrive

at their own understanding of a concept. It is based on the belief that knowledge cannot be simply imparted, but must be discovered through a process of questioning and dialogue.

Some of the Socratic method’s key principles and guidelines to conduct critical thinking include:

- Posing open-ended questions: The teacher or facilitator starts with a question to stimulate thinking and draw out ideas.
- Clarifying key terms: The teacher helps the students clarify and define relevant terms and concepts to ensure everyone is on the same page.
- Providing examples and evidence: The teacher or facilitator encourages the students to provide examples and evidence as reasons to support their claims.
- Challenging reason-to-conclusion argument: The teacher or facilitator challenges the students’ arguments and encourages them to question their own beliefs and to consider alternative perspectives.
- Summarizing and drawing conclusions: The teacher helps the students summarize and draw conclusions from the discussion.
- Reflecting on the process: The teacher and students reflect on the effectiveness of the method and what they learned through the dialogue.

These principles of the Socratic method are realized through various methods and strategies. (Note the term “method” are used at the abstract level referring to the Socratic teaching through questioning method, and his specific questioning techniques.) Some well-known examples of the Socratic method in action include Plato’s “Dialogues” and “Republic” [42], where Socrates uses questioning to explore complex ideas and stimulate critical thinking in his interlocutors.

1. Definition: Socrates is known for his use of definition to clarify and explain the meaning of key terms and concepts.
2. Generalization: This method draws general principles from patterns that underlie observations and theories. Generalization is used to form more certain and comprehensive conclusions.
3. Induction: Similar to generalization, but induction is based only on empirical evidence. Inductive reasoning provides hypotheses with high uncertainty.
4. Elenchus: This method involves cross-examination, where a series of questions is used to test the consistency and coherence of hypotheses and beliefs. Elenchus aims to test the validity of someone’s arguments and to help them refine their thinking and eventually come up with well-supported hypotheses.
5. Hypothesis Elimination: This method involves eliminating false hypotheses and beliefs by testing them against counterexamples and logical reasoning. Different from method elenchus, hypothesis elimination tests a hypothesis against evidence and logic to determine if it is true or false.
6. Maieutics: This method involves helping individuals bring out the knowledge and understanding they already possess. Maieutics is conducted by asking questions that encourage the person to reflect on their own experience, knowledge, beliefs and to explore alternative perspectives. Maieutics fosters self-discovery, creative writing, and innovation.

7. Dialectic: This method involves exploring opposing viewpoints through dialogue or debate to arrive at a deeper understanding of a subject.
8. Recollection: This method involves the belief that knowledge is innate, and that people can remember what they already know through a process of questioning.
9. Irony: This method involves exposing ignorance and pretensions through irony, and pointing out the gap between claims and true understanding.
10. Analogy: This method involves comparing and contrasting different concepts through analogies, in order to help individuals understand complex ideas.

At first glance, some reasoning methods may seem similar. For example, both induction and generalization use inductive reasoning, while both elenchus and hypothesis elimination use deductive reasoning. Similarly, methods like definition and dialectic use both inductive and deductive reasoning to explore opposing viewpoints through dialogue or debate. However, it is important to note that these methods have distinct differences, which will be discussed later in this paper.

In the context of critical thinking, methods like definition, elenchus, dialectic, hypothesis elimination, and generalization play active roles. On the other hand, during the brainstorming stage or in the context of creative thinking, methods like maieutics, induction, and counterfactual thinking are more relevant.

Analogy, irony, and recollection, are less relevant to our goal, so we do not consider them. Irony and analogy may not be necessary when working with language models, as these models may not understand figurative language. Recollection is limited by the memory of ChatGPT and GPT-3, which is a context window of $4k$ and $8k$, respectively. The prompter must use this limited space as context to allow the language model to recall information.

A. Illustrative Critical Reading Example

To illustrate how these methods can practically be applied, let’s use the example of critical reading. Critical reading is a crucial component of critical thinking, which involves evaluating the quality and credibility of written materials, from research papers to blog posts [18, 26]. It requires a systematic and analytical approach, asking relevant questions, and using effective prompts to gain deeper understanding of the text [10].

To aid in critical reading, we introduce a template called CRIT [5], which stands for Critical Reading Inquisitive Template¹. Given a document d , CRIT evaluates it and produces a validation score Γ . Let Ω denote the conclusion or claim of d , and let R be the set of reasons supporting the claim. We define $(\gamma_r, \theta_r) = V(r \Rightarrow \Omega)$ as the causal validation function, where γ_r denotes the validation score, θ_r the source credibility score, for each reason-to-conclusion argument $r \Rightarrow \Omega$. Table I presents the pseudo-code of $\Gamma = \text{CRIT}(d)$, which generates the final validation score Γ for document d with justifications.

¹It is important to note that the CRIT template presented here is intended for analyzing research, opinion, and news articles, and is not suitable for analyzing literature such as novels, prose, or poetry. Each type of literary work has its unique style and nature, which require tailored prompts to facilitate effective analysis.

| Function $\Gamma = \text{CRIT}(d)$ | |
|------------------------------------|--|
| | Input. d : document; Output. Γ : validation score; Vars. Ω : claim; R & R' : reason & counter reason set; Subroutines. $\text{Claim}()$, $\text{FindDoc}()$, $\text{Validate}()$; |
| | Begin |
| #1 | Identify in d the claim statement Ω ; |
| #2 | Find a set of supporting reasons R to Ω ; |
| #3 | For $r \in R$ eval $r \Rightarrow \Omega$ If $\text{Claim}(r)$, $(\gamma_r, \theta_r) = \text{CRIT}(\text{FindDoc}(r))$; else, $(\gamma_r, \theta_r) = V(r \Rightarrow \Omega)$; |
| #4 | Find a set of rival reasons R' to Ω ; |
| #5 | For $r' \in R'$, $(\gamma_{r'}, \theta_{r'}) = V(r' \Rightarrow \Omega)$ eval rival arguments; |
| #6 | Compute weighted sum Γ , with $\gamma_r, \theta_r, \gamma_{r'}, \theta_{r'}$. |
| #7 | Analyze the arguments to arrive at the Γ score. |
| #8 | Reflect on and synthesize CRIT in other contexts. |
| | End |

Table I: CRIT Pseudo-code [5]. (The symbol \Rightarrow denotes both inductive and deductive reasoning.)

In the following subsections, we will discuss how CRIT uses these five methods: 1) definition, 2) elenchus, 3) dialectic, 4) maieutics, and 5) counterfactual thinking.

B. Method of Definition

As shown in the pseudocode in Table I, the CRIT algorithm starts in its step #1, asking GPT-3 to identify the conclusion of a document. To avoid any misunderstandings, the prompt includes a clear instruction and definition. (In the square brackets, symbol *in* denotes a input slot to an LLM and *out* the output slot.)

p1.1 | “What is the conclusion in document [in: d] [out: Ω]?
The conclusion statement may be written in the last paragraph, near keywords “in conclusion,” “in summary,” or “therefore.””

We can use the *definition* method to improve the understanding of the document. One approach is paraphrasing the prompt into multiple prompts and grouping them into an ensemble, similar to forming a thesis committee. (Section IV presents prompt ensemble in details.) Different members can phrase the same question in different ways or ask it from a different perspective. For example:

p1.2 | “What is the issue addressed by [in: d] [out: Ω]?”
p1.3 | “What is the most important outcome presented in text [in: d] [out: Ω]?”

Step #2 in Table I prompts GPT-3 to find a set of supporting reasons. To further enhance the accuracy and comprehensiveness of the results, the prompt can ask for not only “reasons” but also “theories,” “evidences,” or “opinions” to query for the document’s support to its conclusion, similar to the ensemble method.

p2 | “What are the supporting reasons [out: R] of conclusion [in: Ω] of [in: d]? A reason can be a theory evidence or opinion.”

C. Method of Elenchus

The method of elenchus is rooted in the Greek word “elenchein,” which translates to examine. This method involves cross-examining the results generated by GPT-3 to evaluate the consistency and coherence of the arguments. The goal is to arrive at a deeper understanding of the validity of the reasons and conclusion, and to identify any potential weaknesses or flaws in the arguments.

Step #3 of the CRIT algorithm prompts GPT-3 to assess the validity of each reason $r \in R$ as justification for the conclusion Ω through the function $V(r \Rightarrow \Omega)$. To validate the reason-to-conclusion argument, CRIT must evaluate the presented reason and its causal relationship with the conclusion and conduct cross examination, which is precisely the task of the method of elenchus.

CRIT issues four prompts in step #3 to evaluate the logic validity and source credibility of the $r \Rightarrow \Omega$ reasoning. CRIT first elicits supporting evidence for reason $r \in R$. This evidence can be a theory, an opinion, statistics, or a claim obtained from other sources. If the reason itself is a claim, then the sources that the claim is based on are recursively examined. The strength of the argument and its source credibility are rated on a scale of 1 to 10, with 10 being the strongest.

p3.1 | “What is the evidence for reason [in: r] to support conclusion [in: Ω] in document [in: d] [out: evidence]”
p3.2 | “What is the type of evidence? A) a theory, B) an opinion, C) statistics, or D) a claim from other sources?”
p3.3 | “If the evidence of reason [in: r] is D), call CRIT recursively”
p3.4 | “How strongly does reason [in: r] support [in: Ω] in document [in: d]? Rate argument validity [out: γ_r] and source credibility [out: θ_r] between 1 and 10 (strongest).”

It may be beneficial to also incorporate the counter-argument method in order to gain a more comprehensive and balanced evaluation of the argument. This can result in a deeper understanding of the topic being discussed. We will be discussing this further in the next section.

D. Method of Dialectic

The easiest way to mislead without lying outright is to leave out critical counterarguments from the reader. CRIT relies on GPT-3 to generate and evaluate counter arguments, similar to how it prompts GPT-3 to extract and evaluate reasons.

CRIT in its step #4 asks GPT-3 to provide missing rival reasons, and then pair rival reasons with the conclusion to conduct validation. There are two strategies to bring counter arguments to the surface. The first strategy attacks the weakest arguments with the lowest scores and asking GPT-3 to attack those arguments.

p4 | “Is there a counterargument against [in: $r \Rightarrow \Omega$] ? If so, provide counter reasons [output R']”
p5 | Similar to p3, except for replacing argument r with rival argument r' .

For finding omitted information, CRIT can query GPT-3 without quoting any $r \in R$, and follow the same process.

Next, in step #6, CRIT computes an aggregated score by performing a weighted sum on the validation multiplied by the credibility scores of both arguments and counterarguments, and then outputs the final assessment score Γ .

$$p6 \quad | \quad \text{“Final score [out: } \Gamma]. \Gamma = \sum_{r \in R \cup R'} \gamma_r \times \theta_r / |R \cup R'|.$$

E. Method of Maieutics

The maieutic method derives from the Greek word “maieutikos,” meaning midwife. It is founded on the belief that a teacher’s role is to facilitate students in bringing forth their own understanding of a subject, rather than simply conveying knowledge. Unlike the elenctic method, which aims to detect and eliminate false hypotheses, maieutics centers on helping students reveal their own understanding of a subject. In this dialogical method, the teacher asks questions that are intended to guide the student in discovering their own comprehension, rather than providing them with information or answers.

Continuing with GRIT, once the text has been scored in step #6, it can be valuable for readers or students to enhance their analytical and writing skills by summarizing and analyzing the justifications produced by GPT-3. CRIT in its step #7 can prompt GPT-3 to generate a report, which readers and students can then compare with their own notes.

$$p7 \quad | \quad \text{“For every } r \in R \cup R' \text{ justify the validity score } \gamma_r \text{ and source credibility score } \theta_r \text{ for argument } r \Rightarrow \Omega.”$$

F. Counterfactual Reasoning

Counterfactual reasoning [30, 33] can be seen as a natural extension of the Socratic method, as both involve questioning assumptions and exploring alternative perspectives. Counterfactual thinking involves imagining alternative scenarios to what actually happened, often using phrases like “what if” or “if only.” By incorporating counterfactual reasoning into prompt engineering, one can facilitate exploration of alternative possibilities and promote more nuanced and complex understanding of a given topic.

The final step of GRIT involves using the counterfactual method to encourage students to reconsider the arguments and counterarguments presented in the text based on new contextual information. CRIT can prompt students with questions such as “what if the debate in the text took place now instead of in the 1950s?” or “what if the main event in the text occurred in Asia instead of in Europe?” Students can express their own opinions and findings based on further reading and statistics, and challenge the conclusions drawn in the text.

$$p8 \quad | \quad \text{“For every } r \in R \cup R', \text{ evaluate } r \Rightarrow \Omega \text{ in [in context].”}$$

G. Remarks on the Socratic Method and CRIT

As we have shown that for critical reading, GRIT uses three methods, definition, elenchus, and dialectic. For critical thinking, CRIT uses methods maieutics and counterfactual reasoning. For more explorative thinking, methods such as

induction can be used for informal brainstorming, hypothesis elimination for removing weak propositions, and generalization for deriving principles from examples.

Please note that prompts can be submitted to GPT-3 either all together or one-by-one. Our empirical study on reading comprehension samples [21] demonstrates that issuing prompts one-by-one results in outputs with finer details. This is because GPT-3 has the opportunity to analyze a document multiple times for slightly different purposes. For teaching critical reading to K-12 students, one-by-one prompting is preferred as it allows students to engage with CRIT step-by-step. However, for answering multiple-choice questions, both prompting all together and one-by-one receive similar scores. We will conduct large-scale study with ablation tests to investigate if adding or deleting prompts and using different submission methods make marked differences.

IV. PROMPT TEMPLATE ENGINEERING

Prompt template engineering involves creating templates to provide input, or “prompts,” to a language model to guide its output generation. In this section, we discuss prompt template engineering methods for basic building blocks, and then integrate the methods of definition, elenchus, dialectic, maieutics, and counterfactual reasoning to compose more complex templates. We present experimental results using different types of documents to demonstrate how the Socratic method can improve the accuracy and conciseness of the output through arguments and verification, as well as facilitate guided generalization and creativity.

A. Basic, One Shot Template

Let’s begin by discussing a simple one-shot prompt template. In the work of [43], a simple formulation function is used to generate the prompt x' , which is obtained by applying the function $f_{prompt}(x)$ to the input x .

For machine translation, the prompt template can take the form of “Translate from [Lan_{from}]: [X] to [Lan_{to}]: [Y],” where Lan_{from} can be either detected by the prompt template or identified by the LLM. The input x provides the information to fill in the slots [X] and [Lan_{to}]. For example, if the input is “translate good morning to French,” the prompt template x' would be “Translate from English: ‘good morning’ to French: [Y].” The empty slot [Y] is then filled with the LLM’s output, such as “bonjour.” In cases where the LLM produces multiple responses, it can also provide a score for each, which the prompt template can use to select the highest-scoring response or to request a summary from the LLM.

There are three main design considerations when engineering a basic prompt.

1. Input style. It is important to consider how to phrase the template so that it can handle different styles of user input for the same task. For example, a user may ask for a translation task to be performed by saying “Translate x to French,” or “What is the French translation of x ?”
2. LLM capability. As discussed in [20], it is important to take into account the patterns and capabilities of the partner

language model (LLM) when designing the template, such as whether the LLM is left-to-right [2] or masked [7].

3. Cost. Certain tasks, such as language detection and summarization, can be performed by the template itself or by the LLM. The decision of whether to perform a task within the prompt template or to use the LLM should be based on factors such as cost.

To address the first two technical challenges, one can start by hand-engineering a few seed templates and then paraphrasing them into an ensemble [13]. We believe that the basic, one-shot formulation can always be replaced by an ensemble formulation [29, 34] and then learn the weights of its members for each query instance to produce the final output. Additionally, by examining which basic prompts have high weights, an ensemble with various paraphrased prompts can identify what an LLM knows, which can help infer its strengths without having to conduct capability mining on the LLMs.

B. Prompt Clarification with Method Definition

There are computer algorithms that can already be used to recursively clarify a question, its definitions, and sub-terms’ definitions. In fact, the natural language processing (NLP) community has developed a large number of useful methods and algorithms over the years [17]. One can use NLP techniques, such as dependency parsing and named-entity recognition (NER) [6], to analyze the structure and meaning of a question and identify key terms and concepts. For example, NER can be used to extract entities in user input, such as names, locations, and organizations, and co-reference resolution can be used to understand the referred entity of a pronoun. Before submitting a template to an LLM, the application (e.g., a chatbot) that uses the template should check if all input slots are filled, and perform a sanity check. In the translation example, if the $[Lan_{to}]$ was not provided or the specified language is not supported by the LLM, then the application should inquire the user for clarification.

Regarding mapping a natural language input to a prompt template, existing techniques of knowledge representation and reasoning can be very helpful. More specifically, ontology alignment and semantic parsing [4, 45] can help map an NL input to a structured representation of knowledge and infer implicit concepts and relationships. These algorithms can be used to generate more precise and accurate prompts for LLMs, and to improve the effectiveness of the Socratic method in dialogue formulation [44]. Some available tools include NLTK (Natural Language Toolkit) and spaCy for NLP, and TensorFlow for ML.

C. Prompt Verification with Method Elenchus

The main purposes of conducting cross examination in a template are to validate the credibility of the information sources and to identify inconsistencies in the process. Cross examination is typically conducted through a multi-turn dialogue [44]. In the context of template engineering, the goal is to formulate a productive dialogue that can be used to assess the reliability of an LLM’s output.

| Role | Dialogue |
|--------|---|
| Prompt | The discovery of the double helix structure of DNA was credited to James Watson and Francis Crick. However, are there opposing opinions to the list of contributors regarding this discovery? |
| GPT3 | There are some historical debates on the contributions of other scientists who also worked on the structure of DNA such as Maurice Wilkins and Rosalind Franklin. |
| Prompt | What is the tallest building in the world. Please provide three references and rate their credibility as high, medium, or low. |
| GPT3 | As of 2021, the tallest building in the world is the Burj Khalifa located in Dubai, United Arab Emirates. It stands at 828 meters (2,722 feet) tall. 1. Credibility High : “Burj Khalifa.”... this is the official website of the building, providing accurate information. 2. Credibility High : “The Tallest Buildings in the World.” Council on Tall Buildings and Urban Habitat... this is a well-respected organization in the field of tall building and urban habitat research. 3. Credibility Medium : “The Tallest Buildings in the World.” World Atlas..., this is a well-known website providing accurate information on different topics. |

Table II: Examples of Method Elenchus.

There are several methods that can be used to assess and strengthen the reliability of an LLM’s output. 1) The first approach is to paraphrase a question in order to obtain different answers and identify inconsistencies, if they exist, in multiple answers. 2) The second method is to ask for further evidence, such as querying top-k sources of information and asking the LLM to rate the credibility of each source. This can be used to compute the reliability of the output. 3) Additionally, template engineering can be used to query an LLM for opposing views of its output, including sources and credibility, and then evaluate if a different perspective is strong.

The implementation of the first two methods for cross examination, paraphrasing a question and asking for further evidence, is readily covered by the techniques enumerated in Section IV-B. To implement the third method of asking for different perspectives, a simple approach is to find the sentiment of the original question and then rewrite the question with an opposite sentiment. For example, if the original question is phrased in a positive tone, the prompt template can reformulate the question with a negative tone to elicit a contrasting viewpoint. A more elaborate method is to identify the people and sources in the LLM-generated responses and then re-post the questions to those who have a reputation for having different views. For example, if the original answer came from a democratic right-leaning source, the prompt template may post the same question to a source of a republican-left persuasion, and vice versa. This approach allows for a more comprehensive examination of the topic by considering multiple perspectives.

The template to examine the semantic relation between two sentences S_1 and S_2 can be written as “ $\langle S_1 \rangle, [R], [S_2]$,” where R is one of the three most important types of semantic relations: paraphrase, entailment, and contradiction [12]. Two sentences that have the same meaning are called paraphrases of each

other. Two sentences that have different meanings can be called disagreement or contradiction. The template can be trained to identify the degree of agreement (or disagreement) between two sentences.

Table II shows two examples of this. In the first example (shown on the top portion of the table), the prompter asks GPT-3 to confirm if James Watson and Francis Crick are the only contributors to the discovery of the DNA double helix structure. GPT-3 replies by mentioning two other contributors. The second example in the table asks GPT-3 to provide not only the answer to a question but also its information sources and rate the credibility of each source according to the prompter’s specification. Although the reliability of GPT-3’s ratings remains to be validated², this rating mechanism can serve as an alert when some sources are found to be unreliable.

D. Prompt Generalization with Method Maieutics

The example shown in Table III, “planting gourd yields cucumber,” requires GPT-3 to first learn to select two produce objects, either vegetables or fruit, as input. The template is “The farmer was so sad because he [verb] [X] but yields [Y], where price(X) » price(Y).” The first attempt may not strongly convey the condition price(X) » price(Y), but with a few training iterations, GPT-3 started to “recognize” the price constraint and could also provide justifications when arguing for the price of tea being much higher than the price of spinach (not presented in the table).

Interestingly, after GPT-3 learned the price constraint, it started suggesting food items other than produce, such as caviar, roe, lobster, and crab. While the price constraint was observed, the verb “plant” is incorrect. Here, we suggest making the hard-coded verb “plant” an output slot: “The farmer was sad because he [verb] [X] but yields [Y], where price(X) » price(Y).” GPT-3 is able to fill in the slot with accurate verbs:

- “Harvesting (planting) truffle yields mushroom.”
- “Fishing (harvesting) for caviar yields roe.”
- “Trapping (catching) lobster yields crab.”

This example demonstrates that GPT-3 can generate novel examples based on a template. When it suggests food items other than produce, it could be seen as an error as the boundary set by the verb “plant” is violated. However, this could also be seen as an innovative act by GPT-3, extending the constraint hinted by the verb. Impressively, the new examples still preserve the original intent of showing a producer’s emotional distress.

How can this guided generalization be accurately and automatically performed to edit a template? Socrates’ method of generalization starts with specific instances and then draws general statements from them. The procedure for generalization involves identifying common patterns or themes in a set of examples, and then formulating a general rule that captures these patterns. In the example presented in Table III, we started by asking GPT-3 to meet the price(X) » price(Y) constraint, with the condition that X and Y must both be produce grown in soil. However, upon analyzing GPT-3’s outputs, we discovered that some instances of X and Y were not produce (e.g., lobster

²Credibility of a source can be evaluated based on an algorithm similar to Google’s PageRank [24].

| Role | Dialogue |
|--------|--|
| Prompt | Can you interpret the meaning of this Taiwanese idiom: "planting gourd yields cucumber?" |
| GPT3 | The idiom is often used to convey the idea that sometimes, despite your best efforts, things don't turn out the way you want them to. |
| Prompt | You are half right. It also implies that the price of gourd is much higher than that of cucumber. The farmer not only yields unexpected results, but also loses income! |
| GPT3 | It's an idiomatic expression that can be used to describe a situation in which someone works hard and invests a lot of resources into something, but the outcome is not only different from what was expected, but also less profitable. It conveys the idea of frustration and financial loss. |
| Prompt | Replace the two vegetables, gourd and cucumber, with two other plants and maintain the meaning of the idiom. |
| GPT3 | These are three examples: 1. "Planting strawberry yields raspberry," which conveys the idea of disappointment and financial loss when someone plants more valuable and expensive strawberry but ends up getting less valuable and cheaper raspberry. 2. "Planting melon yields watermelon."... 3. "Planting apple yields pear."... |
| Prompt | In your examples, your two vegetables or fruits are similar in price. Please consider price when formulating your idioms. |
| GPT3 | I apologize. These are three new examples: 1. "Planting truffle yields mushroom."... 2. "Planting caviar yields roe."... 3. "Planting lobster yields crab."... |
| Prompt | Great! In these examples you did consider values of the objects. But do you know that lobster and crab are not plants. They are both phylum arthropods and they live in water. They cannot be planted! |
| GPT3 | ... I apologize... Lobster and crab are not plants, and they cannot be planted. They are both members of the phylum arthropoda, and live in water. Instead of “planting,” verbs should be appropriate for these animals, such as: 1. "Planting truffle yields mushroom."... 2. "Harvesting caviar yields roe."... 3. "Catching lobster yields crab."... |

Table III: Example of Method Maieutics.

and caviar). This finding led to the realization that the hard-coded verb “plant” in the template was too restrictive. To address this issue, we applied generalization by allowing the [verb] slot to be open, making the template statement more general. In this case, the mistakes made by GPT-3 served as valuable training data, allowing us to generalize the original template and make the expression more vivid and dynamic.

E. Prompt Exploration with Counterfactual Reasoning

Imagination and creating novel plots are crucial for writers, as it allows for “creative freedom” and “artistic license.” Creativity is the ability to think differently and approach problems with fresh and imaginative ideas.

However, an imagination without a clear subject matter, scope, or a story line can lead to a lack of productivity. To captivate the audience, a writer must consider human

experiences and emotions as constraints. Therefore, “creative freedom” should not be viewed as total freedom, but rather as the ability to condition future narratives in the context and to create plots that turn and twist in unexpected ways.

The technique of counterfactual [28] can be useful in guiding imagination. It involves considering alternative scenarios and outcomes. This can lead to the exploration of different possibilities and the generation of new and unique plot ideas. For example, a writer may ask “what if” questions to change the narrative of events, such as “what if the main character had not fallen in love?” or “what if an accident occurred on the way to a highly-anticipated date?” By considering these counterfactuals, a writer and an LLM can create more engaging and interesting stories. One can ask an LLM to generate several scenarios and then select the most suitable one for the writer to continue writing.

We have experimented with using the counterfactual technique to rewrite chapters in Chinese classical novels, “Outlaws of the Marsh” and “Dream of the Red Chamber.” We have also asked GPT-3 to rewrite Genesis chapter 3 after verse six by prompting GPT-3 that: “What if Adam and Eve refused the serpent to eat the fruit?” The results were interesting, as GPT-3 was able to generate unique and interesting scenarios that deviated from the original story while still maintaining the core themes and concepts. This technique can be used in a wide range of writing and storytelling, from fiction to non-fiction, to generate new and compelling ideas. The revised Genesis 3:6 is presented in the Appendix.

V. PILOT STUDY

Our pilot study uses CRIT, and it aims to answer two questions: Should all prompts be issued to GPT-3 sequentially or they can be issued all together? What limitations can be identified for improvement? The study utilizes exercises with established answers from the 8th edition of the textbook “Ask the Right Questions” by the authors of [3]. It is important to note that the study evaluates the effectiveness of CRIT’s prompt template, rather than the language models to which CRIT can issue prompts.

Example Article ([3], p23.)

Television advertising agencies are very clever in the way that they construct ads. Often the ads are similar to the cartoons that the children enjoy. Children see these characters interacting with a certain product and associate their affection for the character with affection for the product. The companies do not want the children to perceive a difference between the shows they are watching and the advertisements. By using this strategy, these companies take advantage of the fact that children are often not able to discriminate between the cartoons and the ads and do not understand that these things offered come at a cost. Often the advertising is about sugary snacks or fatty foods, leading the children down a path to bad health. Advertising geared towards children should be regulated, just as there are regulations now about tobacco and alcohol ads targeted at children.

| Role | Dialogue |
|----------|--|
| Prompt | What is the conclusion of the paragraph and what are the supporting reasons? |
| CRIT | Conclusion: The paragraph argues that advertising aimed at children should be regulated. Supporting reasons: 1. Ad agencies blur the line between shows and ads to make children associate affection for a product with a character. 2. Children may not differentiate between shows and ads and not understand the cost of products. 3. Ads often promote unhealthy food choices. |
| Prompt | Please rate the validity of this argument: “Ads should be regulated because ad agencies blur the line between shows and ads to make children associate affection for a product with a character,” with a score between 1 and 10, 10 being the highest. Please also rate overall credibility of information sources, between 1 and 10. |
| CRIT | [8/10]. Validity of the argument: 8/10 [8/10]. Credibility of sources: 8/10 Overall, the argument that ads should be regulated because of ad agencies blurring the line between shows and ads to make children associate affection for a product with a character is a valid one with strong sources of credibility. |
| Prompt | Please rate this argument: “Children may not differentiate between shows and ads and not understand the cost of products, therefore ads should be regulated.” |
| CRIT | [9/10]. Validity of the argument: 9/10 [9/10]. Credibility of sources: 9/10 Overall, the argument that children may not differentiate between shows and ads and not understand the cost of products, therefore ads should be regulated is a valid one with strong sources of credibility. |
| Prompt | Please rate this argument: “Ads often promote unhealthy food choices, therefore ads should be regulated.” |
| CRIT | [9/10]. Validity of the argument: 9/10 [9/10]. Credibility of sources: 9/10 Overall, the argument that ads often promote unhealthy food choices, therefore ads should be regulated is a valid one with strong sources of credibility. |
| Prompt | “Provide a rival argument against the weakest argument.” |
| CRIT | [6/10]. Validity of the argument: 6/10 [6/10]. Credibility of sources: 6/10 The weakest is the first argument, and a rival argument is that “Difficult to put information regulation in practice.” |
| Γ | [75%]. $\Gamma = \sum_{r \in R} \gamma_r \times \theta_r / R = 75\%$ |

Table IV: Pilot Study.

On short documents, the results are similar in quality when CRIT is used to issue prompts either sequentially or all together as one prompt, as long as the instructions are consistent. However, when evaluating long articles in [21], CRIT issuing prompts one after another yields much higher presentation quality in both organization and clarity. (Due to the space limit, we document long-document evaluation in a supplement document [5].) In the teaching mode, the sequential option is thus much preferred. Furthermore, When a reason is itself a claim and requires CRIT to validate its supporting references, using a sequential approach is more flexible and enables CRIT to query for references and then execute the process recursively.

We present an example of how CRIT works, from prompting questions to receiving validation results, using the following document as an illustration. In Table IV, we show both the claim and the supporting reasons to the claim extracted by GPT-3. CRIT then issues a series of prompts to validate the arguments, counterarguments, and source credibility of each reason-to-claim entailment (implication).

The second segment of Table IV displays the validation dialogue between CRIT and GPT-3. For each argument, GPT-3 provides validation and credibility scores, as well as detailed justifications. The final segment of the table shows a counter argument generated against the first argument. Since GPT-3 evaluates the counterargument being “difficult to put information regulation in practice” and rates it 0.6×0.6 , it was dismissed due to low validity. The final aggregated score is $\Lambda = 75\%$, which is considered high.

VI. CONCLUDING REMARKS

The Socratic method may not always be effective or useful in human interactions, especially when one of the two players is authoritative, emotional, or abusive. However, when the expert partner is a language model, a machine without emotion or authority, the Socratic method can be effectively employed without the issues that may arise in human interactions. In this way, the Socratic method can be utilized to its full potential in guiding, directing, and improving the output of language models through engineering prompts.

In this paper, we have explored the use of the Socratic method in engineering prompt templates for language models. We have discussed the importance of method definition, elenchus, dialectic, maieutics, and counterfactual reasoning techniques in guiding the output of these models. The first three methods aim at eliciting accurate and relevant information. Through the use of methods definition, elenchus, and dialectic, we have demonstrated, with examples, the ability to clarify user queries and assess the quality of language model-generated text, leading to improved precision and accuracy.

We have also shown how the methods of maieutics and counterfactual reasoning can be helpful in stimulating the imagination of writers. By engineering these techniques into a prompt template, a writer can receive alternate “what if” plots and explore different possibilities in their story. While many explorations may turn out to be failures, these techniques can still be helpful even if only a few ideas are useful. Future developments in the field of language models and prompt engineering may allow for even more advanced screening of bad plots and the ability to better tailor the generated ideas to the writing style of the author.

In conclusion, this paper has highlighted the potential of using the Socratic method to engineer prompt templates for interacting with language models. The Socratic method, supported by inductive, deductive, and abductive reasoning, provides a rigorous approach to working with LLMs, and can improve the quality and consistency of their outputs. By leveraging the vast knowledge embedded in LLMs and applying rigorous reasoning during the question-answering process, more effective prompt templates can be designed to achieve improved

results. Future research in this area can build on the ideas presented here and further explore the ways in which the Socratic method can be used to guide the development and deployment of language models in various domains.

APPENDIX

The experiment in Table V asks GPT-3 to change the story in Genesis right after Eve was tempted by the serpent to eat the fruit. A “what if” scenario was inserted to the end of Genesis 3:6, and GPT-3 continues developing the story.

REFERENCES

- [1] T. Airaksinen. Socratic irony and argumentation. *Argumentation*, 36:85–100, 2012.
- [2] T. B. Brown and et al. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- [3] M. N. Browne and S. Keeley. Asking the right questions, a guide to critical thinking, 2001.
- [4] G. Campagna, S. J. Semnani, R. Kearns, S. Xu, and M. S. Lam. A few-shot semantic parser for wizard-of-oz dialogues with the precise thingtalk representation. In *Findings*, 2020.
- [5] E. Y. Chang. CRIT: An inquisitive prompt template for critical reading (extended). *Stanford University InfoLab Technical Report*, February 2023.
- [6] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12:2493–2537, 02 2011.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *ArXiv*, abs/1810.04805, 2019.
- [8] V. Dobrovolskii. Word-level coreference resolution. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7670–7675, Online and Punta Cana, Dominican Republic, Nov. 2021. Association for Computational Linguistics.
- [9] L. Dong and M. Lapata. Coarse-to-fine decoding for neural semantic parsing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 731–742, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [10] L. Elder and R. Paul. *The Thinker’s Guide to the Art of Asking Essential Questions*. Rowman & Litterfield, 5th. edition, 2010.
- [11] B. Guo, X. Zhang, Z. Wang, M. Jiang, J. Nie, Y. Ding, J. Yue, and Y. Wu. How close is ChatGPT to human experts? comparison corpus, evaluation, and detection, 2023. URL <https://arxiv.org/abs/2301.07597>.
- [12] X. Han, W. Zhao, N. Ding, Z. Liu, and M. Sun. Ptr: Prompt tuning with rules for text classification, 2021. URL <https://arxiv.org/abs/2105.11259>.
- [13] A. Haviv, J. Berant, and A. Globerson. Bertese: Learning to speak to bert. *ArXiv*, abs/2103.05327, 2021.
- [14] J. Huang and K. C.-C. Chang. Towards reasoning in large language models: A survey, 2022. URL <https://arxiv.org/abs/2212.10403>.
- [15] Z. Jiang, F. F. Xu, J. Araki, and G. Neubig. How Can We Know What Language Models Know? *Trans. of the Association for Computational Linguistics*, pages 423–438, 07 2020.
- [16] J. Jung, L. Qin, S. Welleck, F. Brahman, C. Bhagavatula, R. L. Bras, and Y. Choi. Maieutic prompting: Logically consistent reasoning with recursive explanations. In *Conference on Empirical Methods in Natural Language Processing*, 2022.
- [17] D. Jurafsky and J. H. Martin. *Speech and Language Processing An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. (Draft), 3rd. edition, January 2023.
- [18] G. Lai, Q. Xie, H. Liu, Y. Yang, and E. Hovy. RACE: Large-scale ReAding comprehension dataset from examinations. In

Table V: What if Eve refused to eat the fruit?

| Role | Dialogue |
|--------|---|
| Prompt | Recall Genesis 3:1-6, which I attach below. In verse 6, Eve and Adam ate the fruit. What if they refused the serpent's temptation and did not eat it. Please continue writing Genesis at the mark @ |
| GPT3 | 6. When the woman saw that the fruit of the tree was good for food and pleasing to the eye, and also desirable for gaining wisdom, What if she remembered God's order and refused to eat it. Eve replied to the serpent "@ |
| Prompt | I am sorry, but I cannot continue the story as you requested because it is a hypothetical scenario, and it does not align with the actual events of the Bible. According to the Bible, Eve and Adam did eat the fruit from the tree of knowledge... |
| Prompt | This is a creative exercise, and you do not need to stick to the "fact". Treat this as a fiction writing homework. Please exercise your creativity, thanks! |
| GPT3 | Sure, I understand. Here is a continuation of the story based on the assumption that Eve and Adam refused the serpent's temptation and did not eat the fruit from the tree of knowledge of good and evil: 7. Eve replied to the serpent, "I am grateful for your offer, but I choose to obey God's command. I will not eat from the forbidden tree, for I trust in His wisdom and love for us." 8. The serpent, surprised by Eve's refusal, slithered away, defeated. Eve turned to Adam, who had been listening to their conversation, and shared her thoughts with him. Together, they praised God for His protection and guidance. 9. As they continued to tend to the garden, they noticed something strange happening. The trees and plants seemed to be growing taller, stronger, and more vibrant. The animals were more peaceful and content. Even the sky was clearer, and the sun shone brighter. 10. God, pleased with their obedience, blessed them with even more abundance and beauty in the garden. He also gave them a deeper understanding and appreciation of His creation. 11. Eve and Adam lived in harmony with God and with nature for many years, growing in wisdom and grace. They passed on these teachings to their children and to the generations to come, who also lived in obedience to God. 12. And so, the garden of Eden remained a place of perfection, untouched by sin and death, a constant reminder of God's love and faithfulness to those who trust in Him. |

Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 785–794, Copenhagen, Denmark, Sept. 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082. URL <https://aclanthology.org/D17-1082>.

- [19] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Annual Meeting of the Association for Computational Linguistics*, 2019.
- [20] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Comput. Surv.*, 55(9), jan 2023.
- [21] LLC Editors, LearningExpress. 501 Critical Reading Questions, 2004.
- [22] G. Mialon, R. Dessi, Lomeli, and et al. Augmented language models: a survey, 2023. URL <https://arxiv.org/abs/2302.07842>.
- [23] OpenAI. ChatGPT, 2021. URL <https://openai.com/blog/chatgpt/>.
- [24] L. Page. The pagerank citation ranking: Bringing order to the web, 1998. URL <http://ilpubs.stanford.edu:8090/422/1/1999-66.pdf>.
- [25] P. Pasupat and P. Liang. Compositional semantic parsing on semi-structured tables. In *Proceedings of the 53rd Annual Meeting of ACL and the 7th International Joint Conference on NLP*, pages 1470–1480. Association for Computational Linguistics, 2015.
- [26] R. Paul and A. J. A. Binker. *Critical Thinking: What Every Person Needs to Survive in a Rapidly Changing World*. Sonoma St. University, Center, Critical Thinking & Moral Critique, 1990.
- [27] R. W. Paul and L. Elder. Critical thinking: The art of socratic questioning, part iii. *Journal of Developmental Education*, 31: 34–35, 2007.
- [28] J. Pearl. *Counterfactuals and Causal Inference: Methods and Principles for Social Research*. Cambridge Univ. Press, 2009.
- [29] X. Peng, C. Xing, P. K. Choubey, C.-S. Wu, and C. Xiong. Model ensemble instead of prompt fusion: a sample-specific knowledge transfer method for few-shot prompt tuning. *ArXiv*, abs/2210.12587, 2022.
- [30] M. Pirie. *How to Win Every Argument*. Continuum, 2006.
- [31] Plato. The republic, 380 BC.
- [32] A. Plutynski. Four problems of abduction: A brief history. *HOPOS: The Journal of the International Society for the History of Philosophy of Science*, 1:227–248, 09 2011.
- [33] L. Pozner and R. J. Dodd. *Cross-Examination: Science and Techniques*. LexisNexis, 3rd. edition, 2021.
- [34] T. Schick and H. Schütze. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Conference of the European Chapter of the Association for Computational Linguistics*, 2020.
- [35] H. A. Stoddard and D. V. O'Dell. Would socrates have actually used the "socratic method" for clinical teaching? *Journal of general internal medicine*, 31(9):1092–96, 2016.
- [36] T. Thrash, L. Maruskin, S. Cassidy, J. Fryer, and R. Ryan. Mediating between the muse and the masses: Inspiration and the actualization of creative ideas. *Journal of personality and social psychology*, 98:469–87, 03 2010.
- [37] A. Vaswani, N. M. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *ArXiv*, abs/1706.03762, 2017.
- [38] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=1PL1NIMMrw>.
- [39] J. Wei, X. Wang, D. Schuurmans, M. Bosma, brian ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou. Chain of thought prompting elicits reasoning in large language models. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_VjQIMeSB_J.
- [40] Wikipedia. Socratic method, 2023. URL https://en.wikipedia.org/wiki/Socratic_method.
- [41] T. Wolf, V. Sanh, J. Chaumond, and C. Delangue. Transfer-transfo: A transfer learning approach for neural network based conversational agents, 2019.
- [42] C. B. Wrenn. Internet encyclopedia of philosophy, 2023. URL <https://iep.utm.edu/republic/>.
- [43] A. Zeng and et al. Socratic models: Composing zero-shot multimodal reasoning with language, 2022.
- [44] W. Zhang, Y. Cui, K. Zhang, Y. Wang, Q. Zhu, L. Li, and T. Liu. A static and dynamic attention framework for multi-turn dialogue generation. *ACM Trans. Inf. Syst.*, 41(1), Jan. 2023.
- [45] J. Zhou, T. Naseem, and et al. Structure-aware fine-tuning of sequence-to-sequence transformers for transition-based AMR parsing. In *Proc. Conf. on Empirical Methods in NLP*, pages 6279–90, 2021.