

A Communication Theory Perspective on Prompting Engineering Methods for Large Language Models

Yuanfeng Song¹, Yuanqin He¹, Xuefang Zhao¹, Hanlin Gu¹, Di Jiang¹, Haijun Yang¹, Lixin Fan¹, and Qiang Yang¹

¹*AI Group, WeBank Co., Ltd, China*

E-mail: {yfsong, yuanqinhe, summerzhao, allengu, dijiang, navyyang, lixinfan, qiangyang}@webank.com;

Abstract The springing up of *Large Language Models* (LLMs) has shifted the community from single-task-orientated natural language processing (NLP) research to a holistic end-to-end multi-task learning paradigm. Along this line of research endeavors in the area, LLM-based *prompting methods* have attracted much attention, partially due to the technological advantages brought by prompt engineering (PE) as well as the underlying NLP principles disclosed by various prompting methods. Traditional supervised learning usually requires training a model based on labeled data and then making predictions. In contrast, PE methods directly use the powerful capabilities of existing LLMs (i.e., GPT-3 and GPT-4) via composing appropriate prompts, especially under few-shot or zero-shot scenarios. Facing the abundance of studies related to the prompting and the ever-evolving nature of this field, this article aims to (i) illustrate a novel perspective to review existing PE methods, within the well-established communication theory framework; (ii) facilitate a better/deeper understanding of developing trends of existing PE methods used in four typical tasks; (iii) shed light on promising research directions for future PE methods.

Keywords Prompting Methods, Large Language Models

1 Introduction

Large Language Models (LLMs) (e.g., GPT-3 [1], GPT-4 [2], LLaMa [3]) make it possible for machines to understand users' attention accurately, thus revolutionizing the human-computer interaction (HCI) paradigm. Compared to traditional machine systems like databases and search engines, LLMs demonstrate impressive capability in understanding, generating, and processing natural language, facilitating a series of services ranging from personal assistants [4], healthcare [5] to e-commercial tools [6] via a unified natural language interface between users and machine.

The research paradigm around LLM has shifted from single-task-orientated natural language processing (NLP) research to a holistic end-to-end multi-task learning approach. Along this line of research endeavors, LLM-based *prompting engineering* (PE) methods [7, 1] have attracted much attention, partially be-

cause they are the key techniques in making full use of the superior capabilities of LLMs via constructing appropriate prompts. PE refers to the process of crafting effective instructions to guide the behavior of LLMs, and it greatly helps in bridging the gap between the pre-training tasks used to construct the LLM with the down-streaming tasks queried by the end users. Through careful prompt designing, users can steer LLM's output in the desired direction, shaping its style, tone, and content to align with their goals.

To this end, numerous prompt engineering (PE) methods have been explored with the notable progress of LLM advancement and technologies [7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24]. A common theme of PE development lies in continuously improving accuracy and responsiveness of designed prompts, which often include components like Role, Context, Input, Output Format, and Examples (see Fig. 1). Specifically, prompt template and answer-

- A typical prompt usually has:
- **Role:** Define the identify that the LLM is emulating, like a doctor or a customer service agent.
 - **Context:** Describe the situation and relevant facts to frame the task or question for the LLM. Providing background context helps guide the response.
 - **Input:** Clearly explain the task or information being requested of the LLM. Concise, direct prompts work best.
 - **Output Format:** Indicate the desired output format and specify the type of output expected focuses the LLM. Typical output includes a conversational response, a summary, or a series of instructional steps, etc.
 - **Examples** (optional): Illustrative examples further clarify the appropriate response style and content. This "few-shot learning" helps steer the capabilities of the LLMs.

Fig.1. Components of A Good Prompt.

ing engineering has evolved from solely utilizing discrete prompts to continuous prompts, and even to exploring hybrid prompts that combine continuous and discrete elements, which provides a larger optimization space to achieve better performance. With the emergent capability of LLM, LLMs are leveraged to plan and use external tools via its in-context learning capability, which significantly enhanced its ability in specialized domains and broadened its application fields.

Following these studies, one can summarize representative PE methods in a chronological overview as illustrated in Fig. 2. These methods can be categorized as three groups that respectively correspond to three prompting tasks proposed to improve the qualities of LLMs' outputs, namely *prompt template engineering*, *prompt answer engineering*, and *multi-turn prompting and multi-prompt learning*. An example of the input and output for the above-mentioned tasks can be found in Table 1.

- First, *prompt template engineering* methods aim to carefully design a piece of "text" that guides

the language models to produce the desired outputs. For example, in Table 1, to finish a classical sentiment detection for a input $A = \text{"Great places to eat near my location!"}$, the prompt template engineering designs a template " $[A]$ Overall, it was a $[Z]$ restaurant" to enforce the LLM to fill the desired comments in the blank i.e. $[Z]$. Essentially this type of template engineering method induces LLM to focus on word embeddings that are relevant to the questions. A common designing principle of existing prompt template engineering methods is to better align information between users and LLMs. Such a trend is manifested by the evolution from using discrete prompts (e.g., a piece of human-readable text) [11, 9] to continuous ones (e.g., a continuous task-specific vector) [13, 20].

- Second, *prompt answer engineering* [7] refers to the process of searching for an answer space and a map to the original output, which enhances users' understanding of the information encapsulated

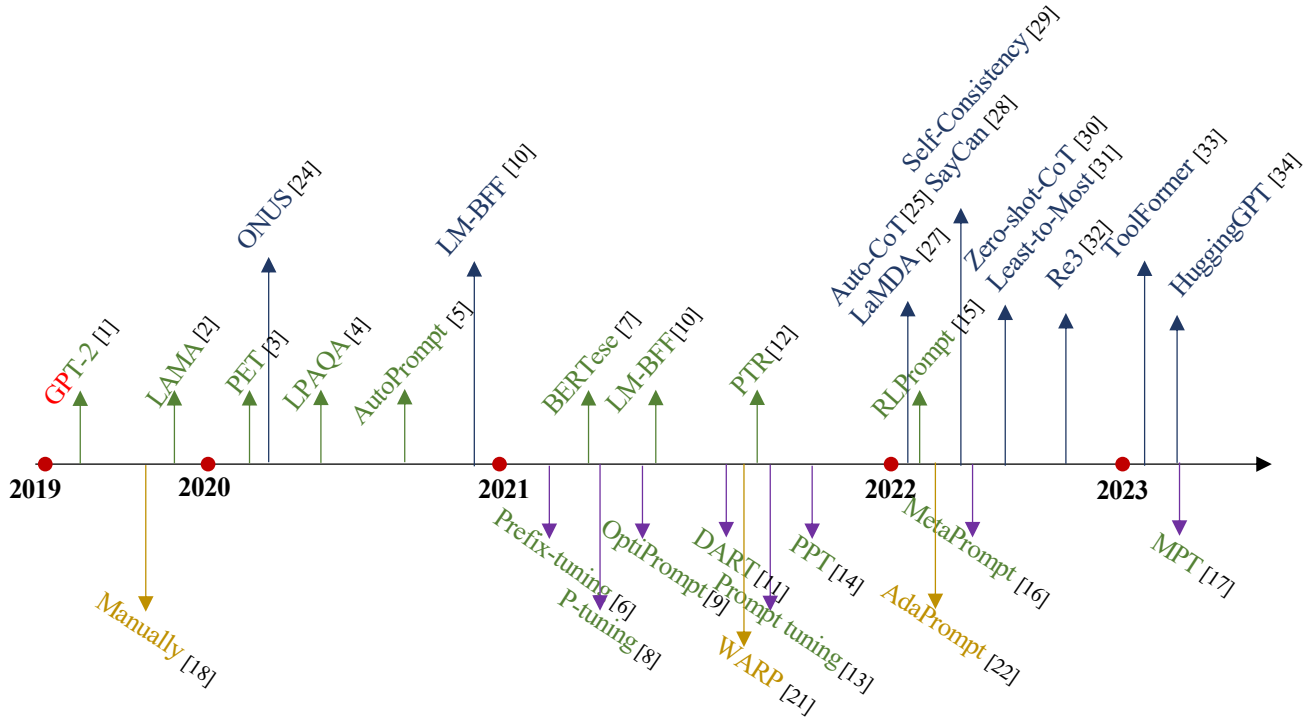


Fig. 2. Chronological overview of representative studies in prompting methods from four aspects: prompt template engineering [8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24], prompt answering engineering [25, 9, 26, 11, 10, 12, 27, 17, 28, 29, 30], and multi-turn prompting and multi-prompt learning [10, 31, 17, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41].

within the LLM. For the same example in Table 1, the prompt answer engineering aims to find a mapping from the result “good” obtained from the LLM to the desired answer “positive”. The field of prompt answer engineering is currently witnessing a notable development trend characterized by the pursuit of models that excel in decoding model information from simple mapping to complex mapping to enhance human comprehension.

- Third, *multi-prompting methods* mainly applied ensemble techniques [10] to mitigate the sensitivity of LLM to different formulations and to obtain a more stable output. In Table 1, the multi-prompting methods combine three different templates (i.e., 1. “It was a [Z]”; 2. “Just [Z]”; 3. “All in all, it was [Z]”;) and their inference results (i.e., 1. “good” 2. “great!” 3. “okay”) to obtain

the final desired one (i.e., “positive”). Later, as LLMs become more capable, multi-turn prompt methods attract more attention that aims to provide more context to LLM by leveraging information either from LLM itself or external tools [37, 33]. In the field of multi-prompting methods, researchers are endeavoring to develop adaptive strategies that enhance LLM’s ability to task planning and the utilization of tools.

In this article, we summarize the prompting methods from a *communication theory* perspective with which the ultimate goal of PE is to *reduce the information misunderstanding between the users and the LLMs*. Therefore, as delineated in Section 2, the communication theory perspective provides a coherent explanation of different PE methods in terms of their objectives and underlying principles. Moreover, this novel perspective also offers and presents insights into scenarios where

Table 1. Running Examples for PE Methods

Stage	Input	Output
Prompt Template Engineering	<i>Great places to eat near my location!</i>	<i>Great places to eat near my location!</i> Overall, it was a [Z] restaurant.
Large Language Model	<i>Great places to eat near my location!</i> Overall, it was a [Z] restaurant.	<i>Great places to eat near my location!</i> Overall, it was a <i>good</i> restaurant.
Prompt Answering Engineering	good	positive
Multi-Prompt	1. It was a [Z]; 2. Just [Z]; 3. All in all, it was [Z];	1. good 2. great! 3. okay

existing prompting methods come short.

The remainder of the article is structured as follows: Section 2 details the overview of the prompting methods from the communication theory perspective. Sections 3, 4, and 5 review and summarize the recent progresses, respectively, from four PE tasks namely prompt template engineering, answer engineering, and multi-turn prompting methods. Section 6 discusses other related surveys and potential research directions. Finally, we conclude this article in Section 7 by summarizing significant findings and discussing potential research directions.

2 A Communication Theory Perspective of Prompting Methods

The study of modern communication theory, which dates back to the 1940s and the following decades, gave rise to a variety of communication models including both linear transmission models and non-linear models such as interaction, transaction, and convergence models [42, 43, 44]. A common theme of these early studies is to analyze how individuals utilize verbal and non-verbal interactions to develop meaning in diverse circumstances. Conceptually, the communication process is often modeled as a chain of information processing steps involving encoding, transmitting, and decoding of messages, between a sender and a receiver.

To give a better illustration, Fig. 3a depicts the classical *Model of Communication* in communication the-

ory, which includes a sender encoding a message and transmitting it to the receiver over a channel. Then, the receiver decodes the message and delivers some type of response. During the transmission process, the message may be distorted due to noise, leading to the necessity of *multi-turn* interaction.

The original communication theory is widely utilized to examine factors including social [45], cultural [46], and psychological [47] that influence human communication. The overall goal of communication theory is to reveal and clarify the common human experience of interacting with others through information exchange.

Among early studies of various communication models, we are particularly inspired by two influential works, namely, Shannon-Weaver’s mathematical model of communication [48] and Schramm’s communication model [49]. Shannon-Weaver’s pioneering work, first published in 1948, provides a strong mathematical foundation to analyze information flow between an active sender and a passive receiver. It is however oversimplistic in the sense that it does not take into account of complexities involved in interactive communication between active senders and receivers, who may respond by sending their message as a form of feedback. The interaction models of communication were first studied by Schramm and published in his 1954 book [49], which pictorially illustrated the feedback loop as depicted in Fig. 3a. Nevertheless, Schramm’s model falls

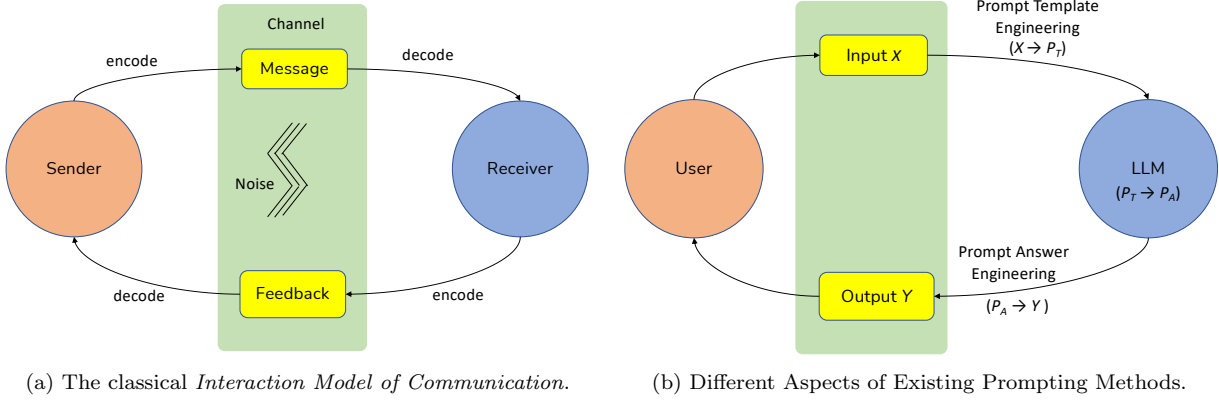


Fig.3. Prompting methods from the communication theory perspective.

short of rigorous theoretical and mathematical formulation to accommodate quantitative analysis e.g. information gain or mutual information between senders and receivers.

Various prompting engineering methods for LLM, in our view, can be understood from Scharmm’s model point of view (see Fig. 3b). In the same vein of Shannon-Weaver’s analysis, we, therefore, delineate a mathematical formulation of Prompting Engineering Systems for interactive user-LLM communication as follows:

Definition 1. *A Prompt Engineering System (PES) consists of a processing chain*

$$X \xrightarrow{g_{\omega_T}} P_T \xrightarrow{f_{\theta}} P_A \xrightarrow{h_{\omega_A}} Y, \quad (1)$$

where g_{ω_T} represents the mapping from the input X to the prompt P_T , f_{θ} denotes the mapping from the prompt P_T to the answer P_A and h_{ω_A} denotes the mapping from the answer P_A to the output Y (see Fig. 3b for an illustration).

Definition 2 (Goal of PES). *PES aims to maximize the mutual information between the inputs X and outputs Y , i.e.,*

$$\max_{\omega_T, \omega_A} I(X, Y) = \max_{\omega_T, \omega_A} I(X, h_{\omega_A} \circ f_{\theta} \circ g_{\omega_T}(X)) \quad (2)$$

where $f \circ g(x) = f(g(x))$.

It’s worth noting that prompt engineering is consistently divided into two procedures: Prompt Template Engineering and Prompt Answer Engineering. Each procedure has specific goals similar to Eq. (2) that align with its intended purpose.

While the capacity in Def. 2 is well-known in information theory [50], how to reach the maximum of Eq. (2) for large language models illustrated in Fig. 3b remains an unexplored research direction. There exists a large variety of prompting engineering methods, which, in our view, essentially aim to reduce information misunderstanding between users and LLMs. In other words, they aim to reach the capacity of PES as defined. This connection between PES and the communication models has never been explicitly stated before.

Moreover, the existing work can be divided into three categories: prompt template engineering ($X \xrightarrow{g_{\omega_T}} P_T$), prompt answer engineering ($P_A \xrightarrow{h_{\omega_A}} Y$), and multi-prompt and multi-tune prompting as shown in Fig. 3b. Specifically, the prompt template engineering aims to *reduce the encoding error/ look for the prompt that is easily understood by the machine*, while the prompt answering engineering aims to *reduce the decoding error/ look for the prompt that easily understood by the human*. The development of LLMs aims to *enhance the capability of the receiver that could better*

handle users’ information needs, and most importantly, the multi-turn prompting and multi-prompt engineering aim to *constantly reduce the information misunderstanding via multi-turn interactions*.

- **Prompt template engineering** aims to optimize

$$\max_{\omega_T} I(X, P_A) = \max_{\omega_T} I(X, f_{\theta} \circ g_{\omega_T}(X)), \quad (3)$$

which looks for an additional piece of text, namely a prompt, to steer the LLMs to produce the desired outputs for downstream tasks. From the communication theory perspective, it acts as an “encoder” to bridge the gap between the users and the LLMs by encoding the messages in a way that the model can understand and then elicit knowledge from LLMs (see details in Section 3). In the encoding process, the challenge lies in the accurate understanding of the user’s intention by LLM with limited instruction following capability. Template engineering aims to reduce this mismatch by translating the user’s request to a format that could be better understood by LLM.

- **Prompt answer engineering** aims to optimize

$$\max_{\omega_A} I(P_T, Y) = \max_{\omega_A} I(P_T, h_{\omega_A} \circ f_{\theta}(P_T)), \quad (4)$$

which focuses on developing appropriate inputs for prompting methods, has two goals: 1) search for a prompt answer P_A ; 2) look for a map to the target output Y that will result in an accurate predictive model. In the decoding process, LLM-generated output often carries redundant information in addition to the expected answer due to its unlimited output space. Answer engineering aims to confine the output space and extract the target answer. The field of prompt answer engineering is currently witnessing a notable develop-

ment trend characterized by the pursuit of effective answer engineering such that ultimate outputs (i.e. Y) are well aligned with that of end users’ expectations (see details in Section 4)

- To further reduce the information misunderstanding, the user could conduct multi-interaction according to Eq. (3) and Eq. (4), called **multi-prompt/multi-turn PE**. *Multi-prompting methods* aims to optimize

$$\max_{\omega_{T_i}} \sum_{i=1}^M I(X, f_{\theta} \circ g_{\omega_{T_i}}(X)), \quad (5)$$

which mainly applied ensemble techniques [10] to mitigate the sensitivity of LLM to different formulations and to obtain a more stable output. Later, as LLMs become more capable, multi-turn prompt methods focus to provide more context to LLM by leveraging multiple communication procedures between the machine and person [37, 33]. In the field of multi-prompting methods, researchers are endeavoring to develop adaptive strategies that enhance LLM’s ability to task planning and the utilization of tools. The adaptive and iterative nature of multi-prompting methods is by the communication theory (see Section 5 for an elaborated explanation).

3 Prompt Template Engineering

Given the information chain $X \rightarrow P_T \rightarrow P_A$, the answer P_A is determined by the prompt-processed P_T and model M with pre-trained weights θ . Suppose that \bar{P}_A is the targeted prediction, the key problem of prompt template engineering is to find a good prompt that maximizes the probability $p(\bar{P}_A|M, P_T, \theta)$ on diverse downstream tasks with limited data. To obtain the optimal prompt, current works [8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24] can be formulated into

three categories: constructing P_T , ranking P_T and tuning P_T , as shown in Fig. 4.

3.1 Constructing P_T

The basic motivation of constructing P_T is to transform the specific task to make it align with the pre-training objective (i.e., next-word prediction, masked LM) of the LM. As shown in Table 2, existing prompt constructing methods [9, 10, 10, 51, 9, 8, 15, 52, 53, 53, 11] could be categorized into five different approaches, which are discussed in detail as follows.

3.1.1 Manually-designed

Initially, the prompt templates are manually designed in natural language based on the user’s experience, and they have been validated to be able to improve the performances of downstream tasks, especially in a zero-shot setting [1, 8]. The most frequent style is to reformulate the original task as a ‘fill-in-blank’ cloze one [9, 10], and the answer is obtained by predicting the words in the given [mask] place. For example, as illustrated in Fig. 4 and Table 2, Petroni *et al.* [9] manually designed prompts to re-structure the relational knowledge, while studies like [10, 51] are dedicated to solving the text classification and language understanding tasks by several self-defining prompt patterns and propose a new training procedure named PET. Another line of work involves developing prefix prompts for generation tasks, which provide instructions and steer the LLMs to finish the sentence. For example, a summarization task can be handled by adding ‘TL;DR:’ [8], and a translation task can be conducted into ‘*Eng.* Translate to Spanish: *Span*’ [52]. Even though manually designed prompts show some effectiveness [53], they are also criticized for being time-consuming and unstable [15]. A subtle difference in the designed prompts may result in a substantial performance decrease. As such, how to explore the prompt space and construct prompts more

thoroughly and more effectively becomes an important and challenging issue.

3.1.2 Heuristic-based

The heuristic-based methods focus on finding acceptable prompts by some intuitive strategies. For example, to construct more flexible and diverse prompts for different examples (rather than the fixed ones), Jiang *et al.* [11] propose to use the most frequent middle words and the phrase spanning in the shortest dependency path that appeared in the training data as a prompt. This method shows a large performance gain compared to the manually-designed prompts. Han *et al.* [19] tries to form task-specific prompts by combining simple human-picked sub-prompts according to some logic rules. Different from the above methods, Logan *et al.* [54] uses an extremely simple uniform rule by *null* prompts, which only concatenates the inputs and the [mask] token, and it’s able to gain a comparable accuracy with manually-defined prompts.

3.1.3 Paraphrasing-based

The paraphrasing-based methods are widely used in data augmentation, aiming at generating augmented data that is semantically related to the original text, and this could be achieved in various ways using machine translation, model-based generation, and rule-based generation [58]. The paraphrasing-based methods could naturally be used to construct prompt candidates based on the original text, and we could further select the best one or integrate them to provide better performance. Representative studies includes [11, 55, 14]. Specifically, Jiang *et al.* [11] uses back-translation to enhance the lexical diversity while keeping the semantic meaning. Yuan *et al.* [55] manually creates some seeds and finds their synonyms to narrow down the search space. Haviv *et al.* [14] uses a BERT-based model to act as a rewriter to obtain prompts that

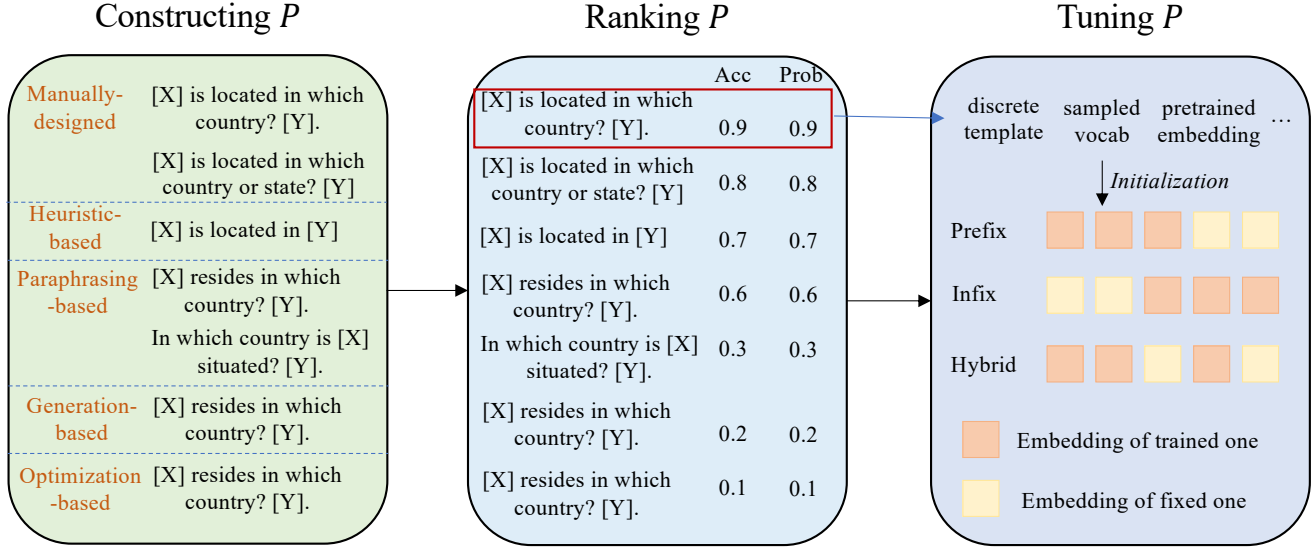


Fig.4. Overview of the PE Methods

Table 2. Summary of the prompt construction methods

Method	Automated	Gradient-Free	Few-shot	Zero-shot	Stability	Interpret-ability
Manually-design [8, 9, 10]	✗	✓	✓	✓	✗	✓
Heuristic-based [11, 54, 19]	✓	✓	✓	✓	✓	✓
Paraphrasing-based [11, 55, 14]	✓	✓	✓	✓	✗	✓
Generation-based [17, 56]	✓	✓	✓	✓	✓	✓
Optimization-based [12, 22, 57]	✓	✗	✓	✗	✓	✗

LLMs can understand better.

3.1.4 Generation-based

The generation-based methods treat prompt searching as a generative task that can be carried out by some LMs. For example, Gao *et al.* [17] first make use of the generative ability of T5 [52] to fill in the placeholders as prompts, and then the prompts could be further improved by encoding domain-specific information [56].

3.1.5 Optimization-based

To alleviate the weakness of insufficient exploration space faced by existing methods, the optimized-based methods try to generate prompts guided by some optimization signals. For example, Shin *et al.* [12] employs gradients as the signals, and then searches for discrete trigger words as prompts to enrich the candidate space. Deng *et al.* [22] generates the prompt using

a reinforced-learning approach that is directed with the reward function.

3.2 Ranking P_T

After obtaining multiple prompt candidates with the above-mentioned methods, the next step is to rank them to select the most effective one. Existing studies solve this problem by finding prompts that are close to the training samples to reduce the information mismatch between the pre-training and inference phases.

3.2.1 Execution Accuracy

Since the objective of the designed prompts is to fulfill the downstream tasks, it's intuitive and straightforward to evaluate the performance by execution accuracy over the specific tasks [57, 17, 11].

3.2.2 Log Probability

The log probability criterion prefers the prompt that delivers the correct output with higher probability, rather than being forced to give the exact answer. For example, a prompt template that can work well for all training examples is given the maximum generated probability in [17]. Furthermore, language models can also be utilized to evaluate the quality of prompts. In [59], the prompt with the highest probability given by an LM is selected, which indicates closer to the general expression that appears in the training dataset.

3.2.3 Others

Other criteria can be used to select the top one or the top-k prompt. For example, Shin *et al.* [12] regards the words that are estimated to have the largest performance improvement as the most crucial elements.

3.3 Tuning P_T

Due to the continuous nature of LLMs, searching over discrete space is sub-optimal [15]. How can we further improve the performance once we obtain a prompt? Recent studies turn to optimizing the prompt as continuous embeddings.

The main idea is to learn a few continuous parameters, referred to as soft prompts, and these continuous parameters can be optionally initialized by the previously obtained discrete prompt. Li *et al.* [13] first introduces a continuous task-specific ‘prefix-tuning’ for generative tasks. Studies like [20] and [15] adopt a similar strategy and prove its effectiveness in various natural language understanding tasks. Following the above-mentioned studies, many improvements have been conducted to find better prompts, such as better optimizing strategies [16], better vector initialization [21, 23], indicative anchors [15] etc. Furthermore, studies like [60, 13, 20] further point out that prompt position,

length, and initialization all affect the performance of continuous prompts [60, 13, 20] (Table 3). In this section, we summarize these factors as follows:

- *Different Position.* There are three different positions for autoregressive LM that the prompt can be inserted into, that is, the prefix [$PREFIX; X_T; Y$], the infix [$X_T; INFIX; Y$], and the hybrid one [$PREFIX; X_T; INFIX; Y$]. There is no significant performance difference between those positions. [13] shows that prefix prompt slightly outperforms infix prompt, and the hybrid one is much more flexible than the others.
- *Different Length.* There is no optimal length for all tasks, but there is always a threshold. The performance will increase before reaching the threshold, then it will either plateau or slightly decrease.
- *Different Initialization.* A proper initialization is essential for the performance of the prompts and the performance of random initialization is usually unsatisfactory. Typical methods include initialized by sampling real word [13, 20], using class label [20], using discrete prompt [16], and using pre-trained based vector [21, 23]. Furthermore, the manually designed prompts can provide a good starting point for the following search process.

3.4 Trends for Prompt Template Engineering

There are two trends in prompt template engineering:

- Tend to have less human involvement, using automated methods rather than designing manually when constructing prompts.
- Tend to develop optimization-based techniques. The gradient-based searching method shows better performance than the derivative-free one in

Table 3. Summary of the prompt tuning methods

Work	Position	Length	Initialization
prefix tuning [13]	prefix, infix	200 (summarization), 10 (table-to-text)	random, real words
prompt tuning [20]	prefix	1, 5, 20, 100, 150	random, sampled vocab, class label
p-tuning [15]	hybrid	3 (prefix), 3 (infix)	LSTM-trained
DART [18]	infix	3	unused token in vocabulary
OPTIPrompt [16]	infix	5, 10	manual prompt
dynamic [60]	hybrid, dynamic	dynamic	sampled vocab

hard prompt construction while the soft prompt is more promising than the hard one.

From the communication theory perspective, the development history of prompting template engineering reflects the trends of utilizing prompts with stronger expressive ability to better capture the user’s intent.

4 Prompt Answering Engineering

As illustrated in Fig. 3(b), prompt answer engineering (PAE) aims to align LLMs outputs with the intended purpose. The use of PAE is motivated by the need to mitigate the gap between the capabilities of pre-trained LLMs and a large variety of requirements of different downstream tasks (see more discussion in Sect. 2). Technology-wise, PAE involves a set of methods that control admissible answer space and optimization mechanisms of LLMs’ output (see overview in Table 4).

4.1 Search for an Answer Space

4.1.1 Pre-defined Answer Space

This involves a set of pre-defined answers for the question-answering task, e.g., pre-defined emotions (“happiness”, “surprise”, “shame”, “anger”, etc.) for the sentiment classification task. The model can then be trained to select the best answer from this pre-defined space. As an illustration, the answer space P_A can be defined as the set of all tokens [9], fixed-length spans [65], or token sequences [8]. Furthermore, in certain tasks like text classification, question answering,

or entity recognition, answers are crafted manually as word lists that pertain to relevant topics [25, 27].

4.1.2 Discrete Answer Space

Discrete answer space refers to a set of specific and distinct answer options that a language model can choose from when generating a response to a given prompt.

Specifically, the possible answers are limited to a fixed set of choices, such as a small number of named entities or keyphrases (e.g., the total choice of planet in the solar system is eight). The model can then be trained to identify whether the correct answer is among this set of possibilities [63, 61, 12].

4.1.3 Continuous Answer Space

Continuous answer space refers to a scenario where the possible answers or responses are not restricted to a predefined set of discrete options. Instead, the answers can take on a range of continuous values or be any text, number, or value within a broader, unbounded spectrum [28, 66].

The model can then be trained to predict a point in this continuous space that corresponds to the correct answer.

4.1.4 Hybrid Approach

This involves combining multiple methods to design the answer space, such as using a pre-defined list of entities for certain types of questions, but allowing for free-form text answers for other types of questions [67].

Table 4. Summary for the prompt answer engineering methods

Answer Space Type	Answer Mapping Method	Work	Task Type
Optimizing the mapping	Discrete Answer Space	[26, 61, 12, 62]	Classification & regression
	Continuous Answer Space	[28]	Classification
Broadening the output	discrete Answer Space	[63]	Generation
Decomposing the output	discrete Answer Space	[64]	Classification
Manually Mapping	Pre-defined answer	[9, 25, 27]	Generation

Remark 1. Answer shapes summarized as follows are also needed in prompt answer engineering. In practice, the choice of answer shape depends on the desired outcome of the task.

- *Tokens:* individual tokens within the vocabulary of a pre-trained Language Model (LLM), or a subset of the vocabulary..
- *Span:* short sequences of multiple tokens, often comprising a phrase or segment of text.
- *Sentence:* A longer segment of text that can encompass one or more complete sentences.

4.2 Search for an Answer Mapping

There are several strategies to search for an answer mapping.

4.2.1 Manually Mapping

In many cases, the mapping from potential answers space P_A to output Y is obvious such that this mapping can be done manually. For instance, the answer is output itself for the translation task [9] such that the mapping is identity mapping; In addition, Yin et al [25] designed related topics (“health”, “food”, “finance”, “sports”, etc.), situations (“shelter”, “water”, “medical assistance”, etc.), or other possible labels. Cui et al. [27] manually proposed some entity tags such as “organization”, “person” and “location”, etc. for the Named Entity Recognition problem.

4.2.2 Broadening the answer P_A

Broadening P_A ($P'_A = B(P_A)$) is expanding the answer space to obtain a more accurate mapping. Jiang et al. [63] proposed a method to paraphrase the answer space P_A by transferring the original prompt into other similar expressions. In their approach, they employed a back-translation technique by first translating prompts into another language and then translating them back, resulting in a set of diverse paraphrased answers. The probability of the final output can be expressed as $P(Y|x) = \sum_{y \in B(P_A)} P(y|x)$, where $B(Y)$ represents the set of possible paraphrased answers.

4.2.3 Decomposing the output

Decomposing Y ($D(Y)$) aims to expand the information of Y , which makes it easier to look for a mapping g_θ . For example, Chen et al. [64] decomposed the labels into several words and regarded them as the answer. Concretely, they decomposed label/output “per:city_of_death” into three separated words {person, city, death}. The probability of final output can be written as $P(y|x) = \sum_{y \in D(Y)} P(y|x)$.

4.2.4 Optimizing the mapping.

There exist two approaches to optimize the mapping function. The first approach is to generate the pruned space \tilde{P}_A and search for a set of answers within this pruned space. Schick et al. [26, 61] introduced a technique for generating a mapping from each label to a singular token that represents its semantic meaning. This mapping, referred to as a *verbalizer* v , is designed

to identify sets of answers. Their approach involves estimating a verbalizer v by maximizing the likelihood w.r.t. the training data conditioned on the verbalizer v . Shin et al. [12] proposed an alternative approach for selecting the answer tokens. They employed logistic classifiers to identify the top-k tokens that yield the highest probability score, which together form the selected answer. In addition, Gao et al. [62] constructed a pruned set \tilde{P}_A^c containing the top-k vocabulary words based on their conditional likelihood for each class c . As for the second approach, it investigates the potential of utilizing soft answer tokens that can be optimized through gradient descent. Hambarzumyan et al. [28] allocated a virtual token to represent each class label and optimized the token embedding for each class along with the prompt token embedding using gradient descent.

4.3 Trends for Prompt Answer Engineering

There are two trends in prompt answer engineering:

- Developing more robust and generalizable question-answering models that can handle more complex tasks and a broader range of inputs. For example, the answer space is some discrete spans at the beginning (see Sect. 6) and developed to the complex continuous space (see Sect. 4.1.3).
- There is also a focus on improving the quality and relevance of prompts to improve model performance. Specifically, several techniques have been explored, such as paraphrasing and pruning, after the direct mapping approach. More recently, optimization methods using gradient descent have been proposed to enhance accuracy.

The prompt answering engineering also shows a trend of exploring prompts to decode the machine language with less information loss, i.e., has a better understanding of the machine.

5 Multiple Prompting Methods

Multiple prompts can be utilized to further reduce the information mismatch during the encoding and decoding process. These methods can be categorized into two main types, namely “multi-prompt engineering” and “multi-turn prompt engineering”, depending on the interrelationship of prompts (see Fig. 5). Multi-prompt engineering is akin to an ensemble system, whereby each response serves as a valid answer, and responses from multiple prompts are aggregated to produce a more stable outcome. This type of method can be thought to extend the use of prompts in the spatial domain. On the other hand, multi-turn PE entails a sequence of prompts, whereby subsequent prompts depend on the response generated from previous prompts or the obtaining of the final answer relies on multiple responses. Consequently, this type of method can be viewed as an extension in the temporal domain.

5.1 Multi-prompt Engineering Methods

Multi-prompt methods employ multiple prompts with similar patterns during the inference aiming to enhance information preservation. This method is closely associated with ensembling techniques [93, 94, 95]. Although the primary motivation is to exploit the complementary advantages of different prompts and reduce the expenses associated with PE, it can also be integrated with prompt-engineering techniques to further improve efficacy. From a communication theory perspective, multi-prompt engineering can be considered as sending multiple copies of the message to ensure the authentic delivery of data.

5.1.1 Expanding P_T

Expanding P_T aims to cover a larger semantic space around the sender’s true intention, and a more stable approximation of the target output, \bar{X}_A , can be ob-

Table 5. Summary of the PE methods involving multiple prompts. NLU: Natural Language Understanding, NLG: Natural Language Generation.

	Method	NLU	NLG	Reasoning
Multi-prompt	Expanding P_T	[11, 20, 28, 68]	[55]	-
	Diversifying P_A	-	-	[69, 70, 71, 72, 73]
	Optimizing θ	[10, 74]	[75, 17]	-
Multi-turn prompt	Decomposing P_T	-	[31, 76, 77]	[38, 78, 79, 80, 81, 82, 83]
	Refining P_T	-	[84, 85]	[32, 37, 85, 69, 86, 87]
	Augmenting P_T	-	[39, 88]	[40, 41, 89]
	Optimizing θ	-	[31, 76]	[89, 90, 91, 92]

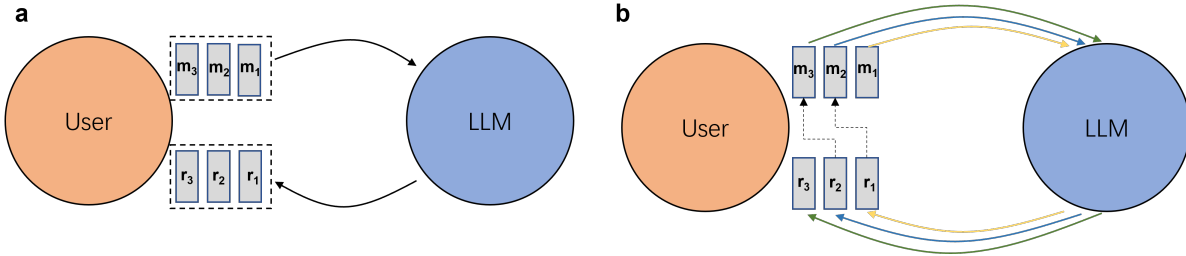


Fig.5. Overview of multiple prompting methods. (a) Multi-prompt methods utilize several similar prompts to produce a more stable result. (b) Multi-turn prompt methods produce the final result by aggregating responses from a sequence of prompts.

tained by aggregating the responses.

Jiang *et al.* [11, 20, 28] propose to combine outputs of different prompts to get the final result for classification tasks. Qin *et al.* [68] incorporates multi-prompt ideas with soft prompts and optimizes weights of each prompt together with prompt parameters. Yuan *et al.* [55] propose to use text generation probability as the score for text generation evaluation, and aggregate multiple results of different prompts as the final score.

5.1.2 Diversifying P_A

Different from expanding P_T whose main goal is to leverage the input space around P_T , diversifying P_A aims to exploit the various "thinking paths" of the LLM through sampling its decoder. This is especially effective for handling complex tasks, such as mathematical and reasoning problems.

Wang *et al.* [69] propose a self-consistency method based on the Chain-of-thoughts (CoT) which samples multiple reasoning paths and selects the most consistent answer by majority voting or weighted averaging. Lewkowycz [70] applied a similar idea to quantitative

problems by combining multiple prompts and output sampling. Wang *et al.* [71] investigated various ensemble variants in reasoning problems and found that rational sampling in the output space is more efficient. These methods solely used the final answer as the selection criterion and did not exploit the generated rationals from various sampling paths. To take advantage of these intermediate results, Li *et al.* [72] proposed to generate more diversified reasoning paths with multiple prompts and used a model-based verifier to select and rank these reasoning paths. Fu *et al.* [73] introduced a complexity-based metric to evaluate reasoning paths and prioritize those with higher complexity in the aggregation. Weng *et al.* [96] employed LLM to self-verify various reasonings by comparing predicted conditions using the generated reasonings to original conditions. The consistency score is then used to select the final result. Yao *et al.* [97] proposed the "Tree of Thoughts" to explore the intermediate steps across various reasoning paths, and used the LLM to evaluate the quality of each possible path.

5.1.3 Optimizing θ

This line of work treats multiple prompts as a label generator to address the sample deficiency problem. Schick *et al.* [10] first proposes pattern-exploiting training (PET) that employs a knowledge distillation strategy to aggregate results from multiple prompt-verbalizer combinations (PVP). They first utilize PVP pairs to train separate models that generate pseudo-labels for unlabeled datasets. This extended dataset is then used to train the final classification model. Schick *et al.* [98] extends this idea to the text generation task by using the generation probability of decoded text as the score. [17] uses a similar method for automatic template generation. Schick *et al.* [74] further expands PET with multiple verbalizers. This is achieved by introducing sample-dependent output space.

5.2 Multi-turn Prompt Engineering Methods

Multi-turn prompt engineering methods involve decomposing the full prompting task into several sub-tasks, each addressed by a corresponding prompt. This process typically entails a sequence of encoding and decoding operations, where subsequent prompts may depend on the decoded message from previous prompts or each prompt is responsible for a sub-task. The outcome can be obtained either from the result of the last prompt or by aggregating the responses generated by all prompts. This strategy is designed to tackle challenging tasks, such as complex mathematical questions or reasoning tasks. It mainly involves two components: 1) decomposing P_T into sub-tasks to reduce the difficulty of each sub-task; and 2) modifying P_T to generate better intermediate results for later steps. These two components can help to bridge the gap between complex X and Y .

5.2.1 Decomposing P_T

Decomposing P_T is the first step in handling complex tasks, and a proper decomposition requires a good understanding of both the target task and the user’s intention. Yang *et al.* [99] decomposed SQL operations using fine-tuned few-shot models and untrained zero-shot models combined with predefined rules. However, ruled-based decomposition heavily relies on human experiences, so it is desirable to automate this step with LLMs. Min *et al.* [76] proposed an unsupervised method that utilizes a similarity-based pseudo-decomposition set as a target to train a seq2seq model as a question generator. The decomposed simple question is then answered by an off-the-shelf single-hop QA model. Perez *et al.* [31] treats the decomposition in multi-hop reading comprehension (RC) task as a span prediction problem which only needs a few hundreds of samples. For each task, various decomposition paths are generated, with each sub-question answered by a single-hop RC model. Finally, a scorer model is used to select the top-scoring answer based on the solving path. Khot *et al.* [77] proposed a text modular network leveraging existing models to build a next-question generator. The training samples are obtained from sub-task models conditioned on distant supervision hints.

With the emergent general ability of LLMs, instead of training a task-specific decomposition model, LLMs are used to fulfill decomposition tasks. Zhou *et al.* [38] proposed the least-to-most prompting method where hard tasks are first reduced to less difficult sub-tasks by LLM. Then answers from previous sub-problems are combined with the original task to facilitate subsequent question solving. Dua *et al.* [78] employs a similar idea and appends both question and answer from the previous stage to the subsequent prompt. Creswell *et al.* [79] proposed a selection-inference framework. It uses LLM to alternatively execute selecting relevant information

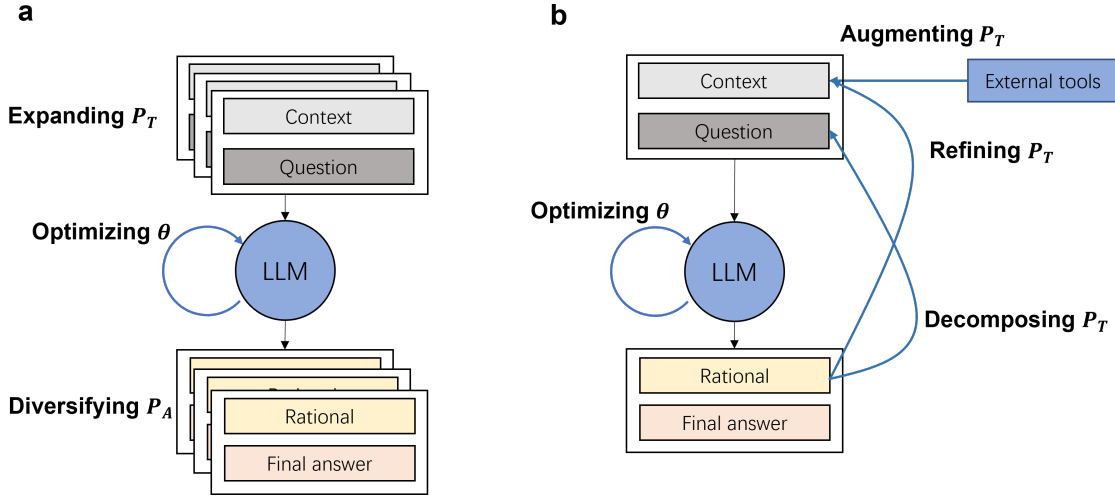


Fig.6. Schematic illustrations of multi-prompting methods. (a) Multi-prompt methods mainly employ ensemble-based methods. (b) Multi-turn prompt methods mainly leverage LLMs or external tools to provide clearer and more helpful context.

from a given context and inferring new facts based on the selected information. Arora *et al.* [80] proposed to format the intermediate steps as open-ended question-answering tasks using LLMs. It further generates a set of prompt chains and uses weak supervision to aggregate the results. Khot *et al.* [81] proposed a modular approach for task decomposition with LLMs by using specialized decomposition prompts. Drozdov *et al.* [100] introduced a dynamic least-to-most prompting method for semantic parsing tasks by utilizing multiple prompts to build a more flexible tree-based decomposition. Ye *et al.* [82] uses LLMs as the decomposer for table-based reasoning tasks. LLMs are used for both sub-table extraction and question decomposition. Press *et al.* [101] proposed Self-Ask which decomposes the original task by repeatedly asking LLM if follow-up questions are needed. Wu *et al.* [83] proposed to build an interactive chaining framework with several primitive operations of LLM to provide better transparency and controllability of using LLMs.

5.2.2 Refining P_T

Refining P_T aims to construct a better representation of P_T based on the feedback from previous prompt-

ing results. This is especially important for multi-step reasoning, where the quality of generated intermediate reasonings has a critical impact on the final answer.

Following the success of the few-shot chain-of-thoughts (CoT) prompting method, Kojima *et al.* [37] proposed a zero-shot CoT method that utilizes the fixed prompt 'Let's think step by step' to generate reasonings. These intermediate results are then fused with the original question to get the final answer. To select more effective exemplars, various methods are proposed. Li *et al.* [84] uses LLMs to first generate a pseudo-QA pool, then a clustering method combined with similarity to the question is adopted to dynamically select QA pairs from the generated QA pool as demonstration exemplars. Shum *et al.* [32] leveraged a high-quality exemplar pool to obtain an exemplar distribution using a variance-reduced policy gradient estimator. Ye *et al.* [85] employs self-consistency method [69] to generate pseudo-labels of an unlabeled dataset. The accuracy of these silver labels serves as the selection criterion of exemplars. To further reduce the search complexity of various combinations, additional surrogate metrics are introduced to estimate the accuracy. Diao *et al.* [86]

addresses this problem by using hard questions with human annotations as exemplars. The hardness is measured by the disagreement of results obtained by multiple sampling of the LLM. Zhang *et al.* [87] proposed automatic CoT methods. They introduced question clustering and demonstration sampling steps to automatically select the best demonstrations for the CoT template.

5.2.3 Augmenting P_T

Different from refining P_T which mainly focuses on finding prompts that generate better intermediate results, augmenting P_T leverages the exploitation of external information, knowledge, tools, etc. in the prompting. We present some examples in this field below, for more details we refer the reader to the specific survey [102]. Yang *et al.* [39] proposed a recursive re-prompting and revision (3R) framework for long story generation leveraging pre-defined outlines. In each step, the context of the current status and the outline of the story is provided to the prompt to ensure better content coherence. Yang *et al.* [88] proposed to use more detailed outlines so that the story generation LLM can focus more on linguistic aspects. Information retrieved from other sources is also often used to augment P_T . Yao *et al.* [103] gives the LLM access to information from Wikipedia. Thoppilan *et al.* [34] taught the LLM to use search engines for knowledge retrieval. More broadly, Paranjape *et al.* [33] introduces a task library to enable the LLM using external tools. Schick *et al.* [40] trained the LLM to use various external tools via API. Shen *et al.* [41] utilized LLM as a central controller to coordinate other models to solve tasks.

5.2.4 Optimizing θ

General LMs are not optimized for producing intermediate rationals or decomposing a complex task or question. Before the era of LLMs, these tasks require

specifically trained LMs. Min *et al.* [76, 31] trained an LM model for decomposing the original task into sub-tasks. Nye *et al.* [90] trains the LLM to produce intermediate steps stored in a scratch pad for later usage. Zelikman *et al.* [91] utilized the intermediate outputs that lead to the correct answer as the target to fine-tune the LLM. Wang *et al.* [89] proposed an iterative prompting framework using a context-aware prompter. The prompter consists of a set of soft prompts that are prepared for the encoder and decoder of the LLMs respectively. Taylor *et al.* [92] employed step-by-step solutions of scientific papers in the training corpus, which enables the LM to output reasoning steps if required.

5.3 Trends for Multiple Prompting Methods

Ensemble-based methods are easy to implement and flexible to incorporate with various strategies, e.g. expanding input space and aggregating output space. However, this brings limited advantages for complex problems whose final answers are hard to obtain directly, but rely heavily on the intermediate thinking steps. Therefore, multi-turn PE methods emerged. It essentially adjusts its input dynamically during the interaction based on the knowledge and feedback from the LLM or external tools. In this way, LLM can leverage more context and understand better the true intention of the user. Initially, specialized LLM are trained to handle planning and solving specific sub-tasks, this not only introduces extra training effort but also constrains the generalization capability of LLM. With the increasing understanding ability and larger input length of LLM, in-context learning becomes the preferred paradigm, which utilizes embedded knowledge and the capability of LLM to handle various tasks via prompting. This paradigm soon dominated because of its efficiency and flexibility.

There are two trends in multiple prompting engi-

neering:

- Developing an enhanced adaptive prompting strategy for LLM-based task decomposition is imperative. The extensive range and intricacy of tasks render human-based or rule-based task decomposition infeasible. While some studies have explored the use of LLM prompting to generate intermediate questions or actions for specific tasks, a comprehensive strategy is currently lacking.
- Enabling LLM to leverage tools without the need for fine-tuning is a crucial objective. By incorporating external tools, LLMs can address their limitations in specialized domains or capabilities. Previous studies have employed fine-tuning-based approaches to train LLMs in utilizing web search or other tools accessible through APIs.

From the communication theory perspective, multiple prompting methods evolved from the extension in the spatial domain (ensemble-based methods) into the temporal domain (multi-turn), to better align the user’s intention and LLM’s capability by decomposing the user’s request and leveraging external tools.

6 Discussion

Researchers have proposed several surveys to recapitulate the rapid advancements in the field of PE methods [7, 104, 105, 106, 107, 108]. To name a few, Liu *et. al* proposed a comprehensive survey about existing PE methods, which covers common aspects like template engineering, answering engineering, training strategies, applications, and challenges [7]. They reveal the development history of prompting learning and describe a set of mathematical notations that could summarize most of the existing studies. Furthermore, they consider prompt-based learning as a new paradigm

that revolves around the way we look at NLP. In another survey [104] that mainly focuses on the reasoning abilities (e.g., arithmetic, commonsense, symbolic reasoning, logical, multi-modal) of LLMs, Qiao *et. al* summarized the studies that harness these reasoning abilities via advanced PE methods like chain-of-thought and generated knowledge prompts. Additionally, some focused surveys cover specific topics like parameter-efficient fine-tuning (PEFT) LLMs using PE methods [105]. Different from the above-mentioned studies, we try to interpret existing PE methods from a communication theory perspective.

Following this line of research, we also would like to discuss some potential challenges and future directions for PE methods, which could be divided into four categories including *Reducing the Encoding Error*, *Reducing the Decoding Error*, and *Interactive and Multi-turn Prompting*.

6.1 Reducing the Encoding Error

- *Better Ranking Criteria.* One of the points of discrete prompts is that it’s difficult to design and choose an optimal prompt, causing its instability. Although soft prompts partly addressed this problem, the discrete prompt is still very important because it has good interpretability and has been proven to be able to help soft prompts search effectively. Looking through the existing methods, we can find that accuracy-based criteria are resource-consuming, while LM-based log probability is not sufficient to evaluate the prompt. So a well-designed ranking criterion combined with a mass of auto-generated prompts may be a good direction for the future.
- *Task-agnostic Prompt.* Even though the prompt has been proven effective in many tasks such as classification and text generation, most of the

existing work has to design a specific prompt for a given task, which makes it complex and complicated [109]. So how to generate a task-agnostic prompt or transfer the prompt to other fields quickly may be a challenging problem. Discrete(meta-learning [110]) and continuous (decomposition [111]) prompts are applied to tackle this issue. However, they are not well-optimized and can serve unseen tasks.

- *Interpretability Issue.* Recent studies show that those methods learning optimal prompts in continuous space can achieve better performance than in discrete space [7]. However, the generated ‘soft’ prompts are difficult to read and understand, namely poor interpretability. Therefore, designing and improving soft prompts can be tough. Existing work [20] tries to use the nearest words in embedding space to probe the effect. However, the reasons excavated are not obvious. It remains to explore why this kind of prompt can work well and what causes the performance differences between different prompts.

6.2 Reducing the Decoding Error

- *Privacy-preserving Methods* [112, 113, 114]. To address privacy concerns on output, future research could focus on developing methods that preserve the privacy of the data used for training and inference. This could include techniques such as differential privacy, homomorphic encryption, and federated learning.
- *Human-in-the-loop Methods* [115]. To improve the accuracy and relevance of prompt answer engineering methods, future research could focus on developing methods that incorporate human feedback and interaction. This could enable users to

provide feedback and corrections to the generated answers and to refine the model over time.

6.3 Interactive and Multi-turn Prompting

- *Transparency and Explainability.* Despite the recent popularity of LLMs, the lack of explainability of the outputs and transparency of the working mechanism makes LLMs less attractive in complex tasks that require high stability. The success of the chain-of-thoughts methodology shows the “thinking path” of LLMs can be evoked with proper indication. This property can be exploited to generate step-by-step task-solving procedures like scratch paper in exams so that the final answer can be better justified. [83] builds an interactive framework based on this idea and further involves human interaction for better controllability of the process. In addition to this online paradigm, LLMs can also be asked to explain the answer or decision afterward. [116] demonstrated that GPT-3 generates more favorable free-text explanations than crowdsourced.
- *Interactive Multi-turn Prompt.* Though automation in prompting methods is highly wanted, humans in the loop can bring more controllability and supervision over the process, producing more reliable results. However, frequent human intervention will diminish the efficiency gained by using LLMs. Therefore, in addition to the granularity of decomposed tasks, it is also required to determine when to involve human feedback. This could be designed manually for each task, but it would be much more efficient if LLMs could plan these stages by themselves.

7 Conclusion

This paper tries to provide an overview of existing prompting methods from a communication theory perspective. Towards this objective, we consider LLMs as a unified interface to achieve various NLP tasks and examine these prompt-based studies to reduce the information misunderstanding that appears in the different stages between users and LLMs during their interactions. We hope this survey will inspire researchers with a new understanding of the related issues in prompting methods, therefore stimulating progress in this promising area.

References

- [1] Brown T, Mann B, Ryder N, Subbiah M, Kaplan J D, Dhariwal P, Neelakantan A, Shyam P, Sastry G, Askell A et al. Language models are few-shot learners. *Advances in neural information processing systems*, 2020, 33:1877–1901.
- [2] OpenAI. Gpt-4 technical report. 2023.
- [3] Touvron H, Lavril T, Izacard G, Martinet X, Lachaux M A, Lacroix T, Rozière B, Goyal N, Hambro E, Azhar F et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [4] Cheng K, Li Z, Li C, Xie R, Guo Q, He Y, Wu H. The potential of gpt-4 as an ai-powered virtual assistant for surgeons specialized in joint arthroplasty. *Annals of Biomedical Engineering*, 2023, pp. 1–5.
- [5] Cascella M, Montomoli J, Bellini V, Bignami E. Evaluating the feasibility of chatgpt in health-care: an analysis of multiple clinical and research scenarios. *Journal of Medical Systems*, 2023, 47(1):33.
- [6] George A S, George A H. A review of chatgpt ai’s impact on several business sectors. *Partners Universal International Innovation Journal*, 2023, 1(1):9–23.
- [7] Liu P, Yuan W, Fu J, Jiang Z, Hayashi H, Neubig G. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 2023, 55(9):1–35.
- [8] Radford A, Wu J, Child R, Luan D, Amodei D, Sutskever I et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019, 1(8):9.
- [9] Petroni F, Rocktäschel T, Lewis P, Bakhtin A, Wu Y, Miller A H, Riedel S. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*, 2019.
- [10] Schick T, Schütze H. Exploiting cloze questions for few shot text classification and natural language inference. *arXiv preprint arXiv:2001.07676*, 2020.
- [11] Jiang Z, Xu F F, Araki J, Neubig G. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 2020, 8:423–438.
- [12] Shin T, Razeghi Y, Logan IV R L, Wallace E, Singh S. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 4222–4235.
- [13] Li X L, Liang P. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association*

for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), 2021, pp. 4582–4597.

- [14] Haviv A, Berant J, Globerson A. Bertese: Learning to speak to bert. *arXiv preprint arXiv:2103.05327*, 2021.
- [15] Liu X, Zheng Y, Du Z, Ding M, Qian Y, Yang Z, Tang J. Gpt understands, too. *arXiv preprint arXiv:2103.10385*, 2021.
- [16] Zhong Z, Friedman D, Chen D. Factual probing is [mask]: Learning vs. learning to recall. *arXiv preprint arXiv:2104.05240*, 2021.
- [17] Gao T, Fisch A, Chen D. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020.
- [18] Zhang N, Li L, Chen X, Deng S, Bi Z, Tan C, Huang F, Chen H. Differentiable prompt makes pre-trained language models better few-shot learners. *arXiv preprint arXiv:2108.13161*, 2021.
- [19] Han X, Zhao W, Ding N, Liu Z, Sun M. Ptr: Prompt tuning with rules for text classification. *AI Open*, 2022, 3:182–192.
- [20] Lester B, Al-Rfou R, Constant N. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 3045–3059.
- [21] Gu Y, Han X, Liu Z, Huang M. Ppt: Pre-trained prompt tuning for few-shot learning. *arXiv preprint arXiv:2109.04332*, 2021.
- [22] Deng M, Wang J, Hsieh C P, Wang Y, Guo H, Shu T, Song M, Xing E P, Hu Z. Rlprompt: Optimizing discrete text prompts with reinforcement learning. *arXiv preprint arXiv:2205.12548*, 2022.
- [23] Hou Y, Dong H, Wang X, Li B, Che W. Metaprompting: Learning to learn better prompts. *arXiv preprint arXiv:2209.11486*, 2022.
- [24] Wang Z, Panda R, Karlinsky L, Feris R, Sun H, Kim Y. Multitask prompt tuning enables parameter-efficient transfer learning. In *The Eleventh International Conference on Learning Representations*, 2023.
- [25] Yin W, Hay J, Roth D. Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach. *arXiv preprint arXiv:1909.00161*, 2019.
- [26] Schick T, Schmid H, Schütze H. Automatically identifying words that can serve as labels for few-shot text classification. *arXiv preprint arXiv:2010.13641*, 2020.
- [27] Cui L, Wu Y, Liu J, Yang S, Zhang Y. Template-based named entity recognition using bart. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2021, pp. 1835–1845.
- [28] Hambardzumyan K, Khachatryan H, May J. WARP: Word-level Adversarial ReProgramming. DOI: 10.48550/arXiv.2101.00121.
- [29] Chen Y, Liu Y, Dong L, Wang S, Zhu C, Zeng M, Zhang Y. Adaprompt: Adaptive model training for prompt-based nlp. *arXiv preprint arXiv:2202.04824*, 2022.
- [30] Kim E, Yoon H, Lee J, Kim M. Accurate and prompt answering framework based on customer

- reviews and question-answer pairs. *Expert Systems with Applications*, 2022, 203:117405.
- [31] Perez E, Lewis P, Yih W t, Cho K, Kiela D. Unsupervised Question Decomposition for Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 8864–8880. DOI: 10.18653/v1/2020.emnlp-main.713.
- [32] Shum K, Diao S, Zhang T. Automatic Prompt Augmentation and Selection with Chain-of-Thought from Labeled Data. DOI: 10.48550/arXiv.2302.12822.
- [33] Paranjape B, Lundberg S, Singh S, Hajishirzi H, Zettlemoyer L, Ribeiro M T. ART: Automatic multi-step reasoning and tool-use for large language models. DOI: 10.48550/arXiv.2303.09014.
- [34] Thoppilan R, De Freitas D, Hall J, Shazeer N, Kulshreshtha A, Cheng H T, Jin A, Bos T, Baker L, Du Y et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- [35] Ahn M, Brohan A, Brown N, Chebotar Y, Cortes O, David B, Finn C, Fu C, Gopalakrishnan K, Hausman K, Herzog A, Ho D, Hsu J, Ibarz J, Ichter B, Irpan A, Jang E, Ruano R J, Jeffrey K, Jesmonth S, Joshi N J, Julian R, Kalashnikov D, Kuang Y, Lee K H, Levine S, Lu Y, Luu L, Parada C, Pastor P, Quiambao J, Rao K, Rettinghouse J, Reyes D, Sermanet P, Sievers N, Tan C, Toshev A, Vanhoucke V, Xia F, Xiao T, Xu P, Xu S, Yan M, Zeng A. Do As I Can, Not As I Say: Grounding Language in Robotic Affordances, August 2022. DOI: 10.48550/arXiv.2204.01691.
- [36] Wang X, Wei J, Schuurmans D, Le Q, Chi E, Narang S, Chowdhery A, Zhou D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- [37] Kojima T, Gu S S, Reid M, Matsuo Y, Iwasawa Y. Large Language Models are Zero-Shot Reasoners. DOI: 10.48550/arXiv.2205.11916.
- [38] Zhou D, Schärli N, Hou L, Wei J, Scales N, Wang X, Schuurmans D, Cui C, Bousquet O, Le Q, Chi E. Least-to-Most Prompting Enables Complex Reasoning in Large Language Models. DOI: 10.48550/arXiv.2205.10625.
- [39] Yang K, Tian Y, Peng N, Klein D. Re3: Generating Longer Stories With Recursive Reprompting and Revision. DOI: 10.48550/arXiv.2210.06774.
- [40] Schick T, Dwivedi-Yu J, Dessì R, Raileanu R, Lomeli M, Zettlemoyer L, Cancedda N, Scialom T. Toolformer: Language Models Can Teach Themselves to Use Tools, February 2023.
- [41] Shen Y, Song K, Tan X, Li D, Lu W, Zhuang Y. HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace, April 2023. DOI: 10.48550/arXiv.2303.17580.
- [42] Narula U. *Handbook of communication models, perspectives, strategies*. Atlantic Publishers & Dist, 2006.
- [43] Chandler D, Munday R. *A dictionary of media and communication*. OUP Oxford, 2011.
- [44] Cobley P, Schulz P J. *Theories and models of communication*, volume 1. Walter de Gruyter, 2013.

- [45] Latané B. Dynamic social impact: The creation of culture by communication. *Journal of communication*, 1996, 46(4):13–25.
- [46] Orbe M P. From the standpoint (s) of traditionally muted groups: Explicating a co-cultural communication theoretical model. *Communication Theory*, 1998, 8(1):1–26.
- [47] Segrin C, Abramson L Y. Negative reactions to depressive behaviors: a communication theories analysis. *Journal of abnormal psychology*, 1994, 103(4):655.
- [48] Shannon C E. A mathematical theory of communication. *The Bell system technical journal*, 1948, 27(3):379–423.
- [49] Schram W E. The process and effects of mass communication. 1954.
- [50] Cover T M. *Elements of information theory*. John Wiley & Sons, 1999.
- [51] Schick T, Schütze H. It’s not just size that matters: Small language models are also few-shot learners. *arXiv preprint arXiv:2009.07118*, 2020.
- [52] Raffel C, Shazeer N, Roberts A, Lee K, Narang S, Matena M, Zhou Y, Li W, Liu P J. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 2020, 21(1):5485–5551.
- [53] Zhou Y, Zhao Y, Shumailov I, Mullins R, Gal Y. Revisiting automated prompting: Are we actually doing better? *arXiv preprint arXiv:2304.03609*, 2023.
- [54] Logan IV R L, Balažević I, Wallace E, Petroni F, Singh S, Riedel S. Cutting down on prompts and parameters: Simple few-shot learning with language models. *arXiv preprint arXiv:2106.13353*, 2021.
- [55] Yuan W, Neubig G, Liu P. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 2021, 34:27263–27277.
- [56] Ben-David E, Oved N, Reichart R. Pada: A prompt-based autoregressive approach for adaptation to unseen domains. *arXiv preprint arXiv:2102.12206*, 2021.
- [57] Zhou Y, Muresanu A I, Han Z, Paster K, Pitis S, Chan H, Ba J. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*, 2022.
- [58] Li B, Hou Y, Che W. Data augmentation approaches in natural language processing: A survey. *AI Open*, 2022, 3:71–90.
- [59] Davison J, Feldman J, Rush A M. Commonsense knowledge mining from pretrained models. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, 2019, pp. 1173–1178.
- [60] Yang X, Cheng W, Zhao X, Petzold L, Chen H. Dynamic prompting: A unified framework for prompt tuning. *arXiv preprint arXiv:2303.02909*, 2023.
- [61] Schick T, Schütze H. Exploiting Cloze Questions for Few Shot Text Classification and Natural Language Inference. DOI: 10.48550/arXiv.2001.07676.

- [62] Gao T, Fisch A, Chen D. Making pre-trained language models better few-shot learners. In *Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL-IJCNLP 2021*, 2021, pp. 3816–3830.
- [63] Jiang Z, Xu F F, Araki J, Neubig G. How Can We Know What Language Models Know? DOI: 10.48550/arXiv.1911.12543.
- [64] Chen Y, Liu Y, Dong L, Wang S, Zhu C, Zeng M, Zhang Y. AdaPrompt: Adaptive model training for prompt-based NLP. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, December 2022, pp. 6057–6068.
- [65] Jiang Z, Anastasopoulos A, Araki J, Ding H, Neubig G. X-factr: Multilingual factual knowledge retrieval from pretrained language models. *arXiv preprint arXiv:2010.06189*, 2020.
- [66] Nickel M, Kiela D. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *International conference on machine learning*, 2018, pp. 3779–3788.
- [67] Hou Y, Che W, Lai Y, Zhou Z, Liu Y, Liu H, Liu T. Few-shot slot tagging with collapsed dependency transfer and label-enhanced task-adaptive projection network. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 1381–1393.
- [68] Qin G, Eisner J. Learning How to Ask: Querying LMs with Mixtures of Soft Prompts. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5203–5212. DOI: 10.18653/v1/2021.naacl-main.410.
- [69] Wang X, Wei J, Schuurmans D, Le Q, Chi E, Narang S, Chowdhery A, Zhou D. Self-Consistency Improves Chain of Thought Reasoning in Language Models. DOI: 10.48550/arXiv.2203.11171.
- [70] Lewkowycz A, Andreassen A, Dohan D, Dyer E, Michalewski H, Ramasesh V, Slone A, Anil C, Schlag I, Gutman-Solo T, Wu Y, Neyshabur B, Gur-Ari G, Misra V. Solving Quantitative Reasoning Problems with Language Models, June 2022.
- [71] Wang X, Wei J, Schuurmans D, Le Q, Chi E, Zhou D. Rationale-Augmented Ensembles in Language Models. DOI: 10.48550/arXiv.2207.00747.
- [72] Li Y, Lin Z, Zhang S, Fu Q, Chen B, Lou J G, Chen W. On the Advance of Making Language Models Better Reasoners.
- [73] Fu Y, Peng H, Sabharwal A, Clark P, Khot T. Complexity-Based Prompting for Multi-Step Reasoning. DOI: 10.48550/arXiv.2210.00720.
- [74] Schick T, Schütze H. It’s Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners. DOI: 10.48550/arXiv.2009.07118.
- [75] Schick T, Schütze H. Few-Shot Text Generation with Pattern-Exploiting Training. DOI: 10.48550/arXiv.2012.11926.
- [76] Min S, Zhong V, Zettlemoyer L, Hajishirzi H. Multi-hop Reading Comprehension through Question Decomposition and Rescoring. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6097–6109. DOI: 10.18653/v1/P19-1613.

- [77] Khot T, Khashabi D, Richardson K, Clark P, Sabharwal A. Text Modular Networks: Learning to Decompose Tasks in the Language of Existing Models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2021*, pp. 1264–1279. DOI: 10.18653/v1/2021.naacl-main.99.
- [78] Dua D, Gupta S, Singh S, Gardner M. Successive Prompting for Decomposing Complex Questions, December 2022. DOI: 10.48550/arXiv.2212.04092.
- [79] Creswell A, Shanahan M, Higgins I. Selection-Inference: Exploiting Large Language Models for Interpretable Logical Reasoning, May 2022. DOI: 10.48550/arXiv.2205.09712.
- [80] Arora S, Narayan A, Chen M F, Orr L, Guha N, Bhatia K, Chami I, Sala F, Ré C. Ask Me Anything: A simple strategy for prompting language models, November 2022. DOI: 10.48550/arXiv.2210.02441.
- [81] Khot T, Trivedi H, Finlayson M, Fu Y, Richardson K, Clark P, Sabharwal A. Decomposed Prompting: A Modular Approach for Solving Complex Tasks.
- [82] Ye Y, Hui B, Yang M, Li B, Huang F, Li Y. Large Language Models are Versatile Decomposers: Decompose Evidence and Questions for Table-based Reasoning. DOI: 10.48550/arXiv.2301.13808.
- [83] Wu T, Terry M, Cai C J. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. DOI: 10.48550/arXiv.2110.01691.
- [84] Li J, Zhang Z, Zhao H. Self-Prompting Large Language Models for Open-Domain QA. DOI: 10.48550/arXiv.2212.08635.
- [85] Ye X, Durrett G. Explanation Selection Using Unlabeled Data for In-Context Learning. DOI: 10.48550/arXiv.2302.04813.
- [86] Diao S, Wang P, Lin Y, Zhang T. Active Prompting with Chain-of-Thought for Large Language Models. DOI: 10.48550/arXiv.2302.12246.
- [87] Zhang Z, Zhang A, Li M, Smola A. Automatic Chain of Thought Prompting in Large Language Models. DOI: 10.48550/arXiv.2210.03493.
- [88] Yang K, Klein D, Peng N, Tian Y. DOC: Improving Long Story Coherence With Detailed Outline Control.
- [89] Wang B, Deng X, Sun H. Iteratively Prompt Pre-trained Language Models for Chain of Thought. DOI: 10.48550/arXiv.2203.08383.
- [90] Nye M, Andreassen A J, Gur-Ari G, Michalewski H, Austin J, Bieber D, Dohan D, Lewkowycz A, Bosma M, Luan D, Sutton C, Odena A. Show Your Work: Scratchpads for Intermediate Computation with Language Models, November 2021. DOI: 10.48550/arXiv.2112.00114.
- [91] Zelikman E, Wu Y, Mu J, Goodman N D. STaR: Bootstrapping Reasoning With Reasoning.
- [92] Taylor R, Kardas M, Cucurull G, Scialom T, Hartshorn A, Saravia E, Poulton A, Kerkez V, Stojnic R. Galactica: A Large Language Model for Science, November 2022.
- [93] Ting K M, Witten I H. Stacked Generalization: When does it work?

- [94] Zhou Z H, Wu J, Tang W. Ensembling neural networks: Many could be better than all. 137(1):239–263. DOI: 10.1016/S0004-3702(02)00190-X.
- [95] Duh K, Sudoh K, Wu X, Tsukada H, Nagata M. Generalized Minimum Bayes Risk System Combination. In *Proceedings of 5th International Joint Conference on Natural Language Processing*, pp. 1356–1360.
- [96] Weng Y, Zhu M, Xia F, Li B, He S, Liu K, Zhao J. Large Language Models are Better Reasoners with Self-Verification, May 2023.
- [97] Yao S, Yu D, Zhao J, Shafran I, Griffiths T L, Cao Y, Narasimhan K. Tree of Thoughts: Deliberate Problem Solving with Large Language Models, May 2023. DOI: 10.48550/arXiv.2305.10601.
- [98] Schick T, Schütze H. Few-shot text generation with natural language instructions. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 390–402.
- [99] Yang J, Jiang H, Yin Q, Zhang D, Yin B, Yang D. SeqZero: Few-shot Compositional Semantic Parsing with Sequential Prompts and Zero-shot Models, May 2022. DOI: 10.48550/arXiv.2205.07381.
- [100] Drozdov A, Schärli N, Akyürek E, Scales N, Song X, Chen X, Bousquet O, Zhou D. Compositional Semantic Parsing with Large Language Models.
- [101] Press O, Zhang M, Min S, Schmidt L, Smith N A, Lewis M. Measuring and Narrowing the Compositionality Gap in Language Models. DOI: 10.48550/arXiv.2210.03350.
- [102] Mialon G, Dessì R, Lomeli M, Nalmpantis C, Pasunuru R, Raileanu R, Rozière B, Schick T, Dwivedi-Yu J, Celikyilmaz A, Grave E, LeCun Y, Scialom T. Augmented Language Models: A Survey. DOI: 10.48550/arXiv.2302.07842.
- [103] Yao S, Zhao J, Yu D, Du N, Shafran I, Narasimhan K, Cao Y. ReAct: Synergizing Reasoning and Acting in Language Models. DOI: 10.48550/arXiv.2210.03629.
- [104] Qiao S, Ou Y, Zhang N, Chen X, Yao Y, Deng S, Tan C, Huang F, Chen H. Reasoning with Language Model Prompting: A Survey. DOI: 10.48550/arXiv.2212.09597.
- [105] Lialin V, Deshpande V, Rumshisky A. Scaling down to scale up: A guide to parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.15647*, 2023.
- [106] Zhao W X, Zhou K, Li J, Tang T, Wang X, Hou Y, Min Y, Zhang B, Zhang J, Dong Z et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [107] Dong Q, Li L, Dai D, Zheng C, Wu Z, Chang B, Sun X, Xu J, Sui Z. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- [108] Lou R, Zhang K, Yin W. Is prompt all you need? no. a comprehensive and broader view of instruction learning. *arXiv preprint arXiv:2303.10475*, 2023.
- [109] Zhong R, Lee K, Zhang Z, Klein D. Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections. *arXiv preprint arXiv:2104.04670*, 2021.
- [110] Reynolds L, McDonell K. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI*

Conference on Human Factors in Computing Systems, 2021, pp. 1–7.

- [111] MADDEN S. Few-shot text-to-sql translation using structure and content prompt learning. 2023.
- [112] Abadi M, Chu A, Goodfellow I, McMahan H B, Mironov I, Talwar K, Zhang L. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 2016, pp. 308–318.
- [113] Gentry C. *A fully homomorphic encryption scheme*. Stanford university, 2009.
- [114] Yang Q, Liu Y, Chen T, Tong Y. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2019, 10(2):1–19.
- [115] Brennan S E, Hulteen E A. Interaction and feedback in a spoken language system: A theoretical framework. *Knowledge-based systems*, 1995, 8(2-3):143–151.
- [116] Wiegrefe S, Hessel J, Swayamdipta S, Riedl M, Choi Y. Reframing Human-AI Collaboration for Generating Free-Text Explanations. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, July 2022, pp. 632–658. DOI: 10.18653/v1/2022.naacl-main.47.