

# Prompting the Unknown: Understanding Response Uncertainty in Large Language Models

Anonymous ACL submission

## Abstract

Large language models (LLMs) are widely used in decision-making across diverse domains. Ensuring the generation of safe and reliable responses is critical for the effective deployment of LLM-based applications, particularly in high-stakes domains such as healthcare and finance. Most of these applications typically use carefully crafted prompts to guide response generation; however, the relationship between prompts and the reliability of LLM-generated responses is not yet fully understood. To address this gap, we propose a novel prompt-response concept model that explains the relationship between the amount of task-relevant information (informativeness) provided in the prompt and the LLM-generated response uncertainty by decomposing response uncertainty into four distinct sources: prompt underspecification, model quality, task variability, and semantic redundancy. We prove that response uncertainty decreases as prompt informativeness or model quality increases, mirroring the behavior of epistemic uncertainty in probabilistic models. Our experimental results on real-world datasets further validate our proposed model and corroborate the theoretical results.

## 1 Introduction

Large language models (LLMs) have demonstrated impressive performance across a variety of tasks (Google, 2023; OpenAI, 2023; Zhao et al., 2023). This success has led to their widespread adoption and significant involvement in decision-making applications across various domains, such as healthcare (Karabacak and Margetis, 2023; Sallam, 2023; Yang et al., 2023), finance (Wu et al., 2023b), education (Xiao et al., 2023), and law (Zhang et al., 2023a). Most of these applications typically use carefully crafted prompts to guide the desired response generation; however, the relationship between prompts and the reliability of LLM-generated outputs, particularly for

high-stakes tasks in domains such as healthcare and finance, is not yet well understood (Arkoudas, 2023; Huang et al., 2023a). Therefore, *understanding how the information in prompts influence response uncertainty is critical to generate safe and reliable responses* for the effective deployment of LLM-based decision-making applications.

To understand this importance, consider a mobile health application (mHealth app) in which machine learning algorithms are incorporated to monitor users' health conditions and provide personalized suggestions on daily activities to influence their decision-making (Boursalie et al., 2018; Trella et al., 2022, 2023). For LLMs to be suitable for such tasks, they must generate safe and consistent responses, e.g., consider an LLM-powered mHealth app that recommends physical therapy (PT) routines to a patient recovering from surgery. To promote patient adherence to PT regimen despite discomfort, the mHealth app must deliver accurate and consistent suggestions. Inconsistent advice can undermine the patient's progress and reduce trust in the application, so *reliable responses are essential for the mHealth app's effectiveness* (Shin et al., 2022).

With the emergent capabilities of LLMs (Dong et al., 2022; Kojima et al., 2022; Wei et al., 2022; Yao et al., 2023), it is possible to improve model responses by guiding them with informative prompts that include task-relevant instructions and exemplars. It helps LLMs to effectively use the relevant information acquired during pretraining to generate better responses, even if the prompt itself does not explicitly reveal the ground truth (Liu et al., 2023; Sahoo et al., 2024). The LLM responses are generated by sequentially sampling tokens from probability distributions over vocabulary, conditioned on the given prompt. LLMs use different decoding strategies such as beam search (Cho et al., 2014), greedy (Sutskever et al., 2014), and nucleus sampling (Holtzman et al., 2019).

Typically, tokens with higher probabilities are chosen sequentially to generate the response. The response variations is controlled by LLM hyperparameters like temperature (Hinton et al., 2015), top- $k$  (Fan et al., 2018), or top- $p$  (Holtzman et al., 2019). While variability benefits creative tasks like poem and essay writing, it can be detrimental for tasks requiring high reproducibility and consistency (Ganguli et al., 2022; Huang et al., 2023c). Making LLMs generate deterministic responses is also not ideal, as users’ preferences may vary (Wu et al., 2023a). Thus, a better approach is needed to understand and separate the different sources of response uncertainty and develop methods to reduce uncertainty naturally rather than masking it by adjusting LLM hyperparameters. As response uncertainty of a fixed or black-box LLM can also be controlled by the amount of task-relevant information in the prompt (i.e., *prompt informativeness*), this paper focuses on *quantifying the response uncertainty due to the prompt* while keeping the LLM and its hyperparameters fixed.

To this end, we use the insight that LLMs implicitly learn to infer latent concepts during pretraining (Xie et al., 2022; Hahn and Goyal, 2023; Zhang et al., 2023b) and propose a *novel prompt-response concept* (PRC) model. Our PRC model conceptualizes how an LLM generates responses based on given prompts and *helps understand the relationship between prompts and response uncertainty* by measuring response uncertainty for prompts with varying task-relevant information. Specifically, it provides a principled framework for separating and analyzing the distinct sources of response uncertainty. We provide theoretical results that show the *uncertainty of responses generated by an LLM decreases as the prompt informativeness and model quality increase*. We draw a connection between the reducible response uncertainty and epistemic uncertainty, and using our PRC model, we *theoretically justify why increasing the task-relevant information in the prompt* is a principled and effective method to reduce the response uncertainty. Finally, we validate the PRC model through experiments and demonstrate the practical applicability of our insights using a healthcare use case. Specifically, our key contributions are as follows:

- **Prompt-Response Concept model.** In Sec. 3, we propose a prompt-response concept model to understand and quantify the relationship between prompt informativeness and response uncertainty in LLMs.

- **Decomposition of response uncertainty.** Using the PRC model, we decompose the response uncertainty into four distinct sources: prompt underspecification, model quality, task variability, and semantic redundancy, where uncertainty due to prompt underspecification is an epistemic uncertainty, while the others are the source of aleatoric uncertainty when the LLM and its hyperparameters are fixed.
- **Theoretical result.** In Sec. 3.2, we prove that response uncertainty decreases as prompt informativeness or model quality increases (Theorem 1). We further link reducible uncertainty to epistemic uncertainty, as adding more relevant information to prompts and improving model quality reduces response uncertainty.
- **Empirical results.** In Sec. 4, we validate the PRC model through experiments and demonstrate its practical applicability across diverse tasks base on real-world datasets and healthcare use case. In addition, we found that, when supplied with sufficient relevant information, a smaller (less capable) model can sometimes outperform a larger (more capable) model with less information.

## 2 Preliminaries

**Problem setting.** Let  $\mathcal{X}$  denote the set of all prompts and  $\mathcal{Y}$  the set of all responses generated by an LLM  $f$ , where  $f : \mathcal{X} \rightarrow \mathcal{Y}$ . Let  $\mathcal{V}$  represent the vocabulary of all unique tokens. For any prompt  $x \in \mathcal{X}$  and response  $y \in \mathcal{Y}$ , we have  $x \in \mathcal{V}^{|x|}$  and  $y \in \mathcal{V}^{|y|}$ , where  $|\cdot|$  denotes the number of tokens in the prompt or response. For a given prompt  $x \in \mathcal{X}$ , the LLM  $f$  generates a response  $y \in f(x)$ , which can vary each time the LLM generates it due to two factors: (i) the LLM hyperparameters, such as temperature, top- $k$ , and top- $p$ , which control the randomness in token sampling, and (ii) the task-relevant information in the prompt, which we refer to as *prompt informativeness*. This paper focuses solely on the latter aspect while keeping the LLM and its hyperparameters fixed. Specifically, this paper aims to explain how prompt informativeness influences response uncertainty for a given LLM and then quantify response uncertainty as a function of both prompt informativeness and the LLM.

**Measure for response uncertainty.** We use entropy to measure the uncertainty in the responses generated by an LLM for a given prompt. Let  $Y$  be

a random variable representing the response. For a given prompt  $x$ , the entropy of  $Y$  is defined as:

$$H(Y|x) = - \sum_{y \in \mathcal{Y}} p(y|x) \log_2 p(y|x), \quad (1)$$

where  $p(y|x)$  is the conditional probability that the response variable  $Y$  takes the value  $y$  given the prompt  $x$ . Intuitively, a more informative prompt reduces the uncertainty of response generated by the LLM, leading to lower entropy in  $Y$ .

**Concept.** The notion of a concept varies across fields: in philosophy, a concept represents the fundamental unit of thought; in psychology, it is a mental construct; in linguistics, it refers to the semantic units that words or phrases represent; and in education, it denotes key ideas or principles. In this paper, we define a *concept* as an abstraction derived from specific instances or occurrences that share common characteristics (Fodor, 1998; Laurence and Margolis, 1999; Weiskopf, 2009; Wilmont et al., 2013). To illustrate this notion of concept, consider the example of *Meeting*, which includes information such as the agenda, date and time, location, and other relevant details about the meeting, for example, “Write an email to schedule a meeting on 28 July 2025, from 14:00 to 15:30, in Room 301 of the Department Building. The agenda includes reviewing the draft paper, discussing key findings, and planning the next steps for submission. Participants include the lead author, co-authors, advisors, and invited peers to give feedback. ...” Recent works (Gao et al., 2024), Lieberum et al. (2024), and Templeton (2024) provide mechanistic interpretability evidence suggesting that LLMs can learn concept-like features.

**Concept attributes.** A concept can be represented as a sequence of sentences conveying semantic meaning, typically expressed in natural language. Here, we use ‘semantic meaning’ (or ‘semantically meaningful’) to refer to information that is both understandable and interpretable by humans (Hurford et al., 2007). As illustrated in the earlier example of *meeting* concept, explaining a concept often involves multiple sentences, each contributing to a specific and meaningful aspect of the concept (Piccinini and Scott, 2006). We refer to each of these aspects as a *concept attribute*. For example, the sentence, “schedule a meeting on 28 May 2025, from 14:00 to 15:30, in Room 301 of the Computer Science Building.” provides information about the meeting’s date, time, and location.

We now formalize the notion of a *concept*. Let  $\Theta$  denote the set of all possible latent concepts, where each  $\theta_t \in \Theta$  corresponds to the concept associated with a specific task  $t$ . Here, task  $t$  refers to the underlying objective, either explicitly specified in the prompt or implicitly expected from the response. In the example of *meeting* concept, “Write an email” is the task specified in the prompt, while the corresponding *email draft containing meeting details* is expected as the response. Let  $\mathcal{A}_{\theta_t} = \{A_{\theta_t,1}, \dots, A_{\theta_t,m}\}$  denote the set of all attributes associated with a concept  $\theta_t \in \Theta$ . Each concept  $\theta_t$  is thus fully characterized by its set of attributes  $\mathcal{A}_{\theta_t}$ . We assume that each attribute  $A_{\theta_t,a} \in \mathcal{A}_{\theta_t}$  can be accurately extracted from, or meaningfully represented in, a semantically coherent sequence of tokens drawn from the set  $\mathcal{V}$ .

### 3 Prompt-Response Concept Model

We first use the notion of concept to explain our proposed prompt-response model for LLMs, then use it to derive theoretical results that explain the relationship between prompt informativeness and response uncertainty by analyzing how uncertainty changes with the informativeness.

#### 3.1 Prompt-Response Concept Model

Building on prior work showing that LLMs implicitly learn to infer latent concepts during pretraining (Xie et al., 2022; Hahn and Goyal, 2023; Zhang et al., 2023b), we propose the *prompt-response concept* (PRC) model of LLM. The PRC model conceptualizes how an LLM generates responses for a given prompt, providing a structured way to analyze the relationship between prompts and response uncertainty by distinguishing its sources. The PRC model consists of three main components (as illustrated in Fig. 1): prompt concept, response concept, and mapping functions.

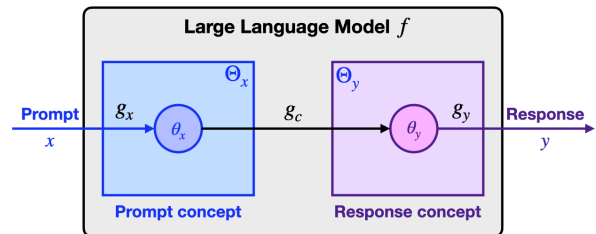


Figure 1: Prompt-Response concept model of LLM.

**Prompt concept.** Let  $\Theta_x$  denote the set of all latent concepts corresponding to prompts in set  $\mathcal{X}$ . Each prompt  $x \in \mathcal{X}$ , expressed as a sequence of tokens describing a task, maps to a concept  $\theta_x \in \Theta_x$ .

We refer to this concept as the *prompt concept*. Intuitively, an LLM extracts the attributes of a latent prompt concept from input tokens, expressed as multiple semantically meaningful sentences in the prompt. If these sentences in a prompt cannot be combined to describe a single concept, the LLM interprets them as representing multiple distinct concepts. Our experiments (Fig. 11c) show that *adding semantically meaningful information from different concepts increases response uncertainty*.

**Response concept.** Let  $\Theta_y$  be the set of all latent concepts corresponding to responses in the set  $\mathcal{Y}$ . We refer to these concepts as the *response concept*. Each response concept  $\theta_y \in \Theta_y$  is uniquely associated with a response  $y \in \mathcal{Y}$ . As prompts and responses follow different distributions and representations, we model them as separate concepts, in contrast to the ICL setting, which uses a single intermediate concept (Xie et al., 2022).

**Mapping functions.** To understand the relationship between the prompt, intermediate concepts, and the response, our PRC model conceptualizes the LLM  $f$  as a composition of three mapping functions: prompt-concept mapping function ( $g_x$ ), concept-concept mapping function ( $g_c$ ), and concept-response mapping function ( $g_y$ ). First, the function  $g_x$  maps the prompt  $x$  to a latent prompt concept  $\theta_x$  that we assume capture the understanding about task and associated relevant information from the prompt. This concept is then transformed by  $g_c$  into a response concept  $\theta_y$ , which encapsulates the core intent and structure of the desired output, which is then mapped to the actual response by  $g_y$ . Intuitively,  $g_x$  captures the LLM’s ability to understand the task specified by the prompt,  $g_c$  models the reasoning process to derive the response concept from this understanding, and  $g_y$  represents the verbalization of the response concept.

It is important to note that both the concept-concept mapping  $g_c$  and the generation process  $g_y$  are inherently stochastic. The stochasticity in  $g_c$  means that for the same prompt concept, the model might stochastically choose different response concepts even if its input is fixed (e.g., explanation focusing on different aspects of the tasks). The stochasticity in  $g_y$  corresponds to the token-sampling process (e.g., using temperature or nucleus sampling), allowing for different phrasings and semantically equivalent variations of the same response concept. These two sources of randomness are fundamental to the response variability

observed in LLMs. In the formal analysis that follows, we regard  $g_y(\theta_y)$  as the *response distribution* over  $\mathcal{Y}$ ; the token-by-token mechanism is thus fully marginalized out. Thus, the overall LLM function can be expressed as  $y = f(x) = g_y(g_c(g_x(x)))$ .

To understand these mapping functions better, consider the following prompt: “Write an email to schedule a meeting next Wednesday at 14:00 for one and a half hours in Room 301.” First, the LLM uses a function  $g_x$  to understand what this prompt is about; in this example, it recognizes the concept of a meeting being scheduled (the prompt concept) and then extracts the key attributes such as the meeting’s date, time, and location. It then applies  $g_c$  to determine the required response, i.e., an email containing meeting details (the response concept). Finally,  $g_y$  verbalizes the response concept into the actual email, i.e., the LLM response, represented as a semantically coherent sequence of tokens drawn from the vocabulary set  $\mathcal{V}$ .

As shown in Fig. 1, both prompt and response concepts are latent components within the LLM. To generate the desired response, the prompt concept must include all task-relevant attributes, either explicitly extracted from the prompt or implicitly inferred from the LLM’s pretrained knowledge. Intuitively, the prompt concept represents the LLM’s internal understanding of the task. Better LLMs extract these attributes more reliably and leverages richer pretrained knowledge, producing responses with lower uncertainty and higher quality, as demonstrated in Fig. 2 and Fig. 3. When the prompt concept lacks sufficient task-relevant attributes, response variability increases, since the LLM can interpret and complete the task in multiple plausible ways (see Fig. 11a).

### 3.2 Theoretical Results

Let  $Z_c$  be a random variable representing a concept (with  $c = x$  for prompt concept and  $c = y$  for response concept) and  $X_s$  be a random variable representing a prompt having semantic meaning  $s$ . We first introduce the assumptions under which our theoretical results hold.

**Assumption 1.** *We assume an LLM is perfect if it has exactly learned the mapping functions of the PRC model, namely  $g_x$ ,  $g_c$ , and  $g_y$ .*

This assumption states that the LLM has perfectly learned the mapping functions used in our proposed PRC model. Although this assumption may not hold exactly in practice, a better LLM

is expected to estimate these mapping functions more accurately, resulting in lower uncertainty and higher-quality responses, as supported by our experimental results shown in Fig. 2 and Fig. 3. Our first result shows the relationship between prompt informativeness and prompt concept uncertainty.

**Lemma 1.** *For any LLM with exactly learned  $g_x$  and two concepts  $\theta_1, \theta_2 \in \Theta_x$ , we have  $\mathcal{X}_{\theta_1} \cap \mathcal{X}_{\theta_2} = \emptyset$  if  $\theta_1 \neq \theta_2$ . Furthermore,  $H(Z_x|X_{\theta_x}) = 0$ .*

The proof of Lemma 1 and other missing proofs are given in Sec. A.2. The first part of this result implies that prompts fully describing two different concepts cannot have the same semantic meaning, i.e., no two concepts share exactly the same semantic description. In other words, the prompts that fully describe two different concepts cannot have the same semantic meaning. The second part implies that the prompt concept is deterministic when the prompt contains all information needed to complete the task. Next, we present a result showing that prompt concept uncertainty decreases as prompt informativeness increases.

**Proposition 1.** *As  $X_s$  represents more informative prompts (i.e., as more task-relevant information is included in the prompt),  $H(Z_x|X_s)$  decreases.*

Our next result links prompt informativeness to the response uncertainty of LLMs.

**Theorem 1.** *As  $X_s$  represents more informative prompts,  $H(Y|X_s)$  strictly decreases. Specifically  $H(Y|X_s) \leq H(Y|Z_y) + H(g_c(Z_x)|Z_x) + H(Z_x|X_s)$ . Furthermore,  $H(Y|X_s)$  converges to  $H(Y_\tau|Z_y) + H(Y_r|Z_y) + \mathcal{E}$ , where  $\mathcal{E}$  decreases as the model quality improves. When Assumption 1 hold,  $\mathcal{E} \rightarrow 0$  as prompt informativeness increases.*

Here,  $H(Z_x|X_s)$  is the uncertainty due to prompt underspecification and imperfect LLM (particularly  $g_c$ ),  $H(g_c(Z_x)|Z_x)$  due to model quality,  $H(Y_\tau|Z_y)$  due to task variability (e.g., response being outcome of coin toss), and  $H(Y_r|Z_y)$  due to semantic redundancy, more discussion about this decomposition is provided in Sec. A.3. The above two results suggest that response uncertainty decreases as prompt informativeness or model quality increases due to the uncertainty in the response concept. It can be understood through our PRC model: a highly informative prompt provides a strong and unambiguous initial concept, which acts as an ‘anchor’ for the response generation. This anchor constrains the possible generated responses to a narrow set of semantically similar outcomes and

thus reduces the final response entropy. Conversely, a vague prompt allows for many divergent generation paths, resulting in higher uncertainty. Furthermore, when sufficient information is provided in a prompt and Assumption 1 hold, no uncertainty ( $\mathcal{E}$ ) remains due to the prompt concept and model quality. The remaining randomness in responses can be decomposed into two terms:  $H(Y|Z_y)$ , which represents the uncertainties due to task variability ( $H(Y_\tau|Z_y)$ ) and semantic redundancy ( $H(Y_r|Z_y)$ ).

The task variability ( $H(Y_\tau|Z_y)$ ) is the aleatoric uncertainty in the prompt task. If the prompt task does not have a deterministic ground truth (e.g., forcing the model to answer the outcome of tossing a fair coin), the aleatoric uncertainty exists and is irreducible. In our experiments, we eliminate this source of uncertainty by using QA datasets for the prompt tasks with deterministic ground truth. We observe multiple realizations of the response concept for the same prompt concept across different iterations due to the imperfection of  $g_c$ , resulting in variations in responses. However, as shown in Fig. 7c, as the LLM quality improves,  $\mathcal{E}$  gets smaller, resulting in lower overall response uncertainty. When an LLM can extract all task-relevant attributes, there should be no randomness in response due to the prompt and the remaining uncertainty in the response arises only from semantic redundancy ( $H(Y_r|Z_y)$ ), which is the ability of the multiple responses expressing the same response concept, i.e., multiple responses are semantically equivalent (Kuhn et al., 2023), which is the irreducible uncertainty (see Theorem 1).

**Remark 1** (Compatibility with autoregressive decoding). *During autoregressive decoding, the model re-feeds the growing token sequence  $(x, Y_{<t})$  at each step  $t$ . Although the input changes, **no new source of randomness is introduced. Thus, the uncertainty decomposition for the full sequence applies directly to each next-token distribution. No additional uncertainty arises from feeding back tokens that are already in-distribution under  $\Theta_y$ .***

### 3.3 Epistemic and Aleatoric Response Uncertainty

In machine learning literature, epistemic uncertainty is typically reduced by incorporating additional information, such as using a better model and additional training data (Hüllermeier and Waegeman, 2021; Lahlou et al., 2021; Senge et al., 2014; Shaker and Hüllermeier, 2020; Valdenegro-Toro

and Mori, 2022; Der Kiureghian and Ditlevsen, 2009). In Proposition 1,  $H(Z_c|X_s)$  represents the epistemic uncertainty in the latent concepts.<sup>1</sup> We have demonstrated that  $H(Z_c|X_s)$  is strictly reduced with a prompt that contains more attributes of the relevant concept(s). Thus, increasing the information about the concept in a prompt can lead to more reliable and consistent responses by reducing the epistemic uncertainty in the latent concept (Hüllermeier and Waegeman, 2021). However, if the prompt contains information that are irrelevant to the task, the response uncertainty can also increase, as demonstrated in Fig. 11d.

Due to the model’s inability to learn a perfect mappings during training as well as the possibility that the task in the prompt may not have a deterministic ground truth, the mapping from the prompt concept to the response concept cannot be deterministic. In Theorem 1,  $\mathcal{E}$  captures the uncertainty due to model quality. In scenarios where model parameters are allowed to be modified, this uncertainty becomes epistemic and can be reduced as the model quality improves (Fig. 7c). The uncertainty due to task variability ( $H(Y_r|Z_y)$ ) is the irreducible aleatoric uncertainty of the given prompt task. The remaining response uncertainty characterized by the term  $H(Y_r|Z_y)$  that arises due to *semantic redundancy*. It can be further reduced in two ways: use fine-tuning or prompting to instruct the model to reply with a certain fixed style.<sup>2</sup> Due to semantic equivalence, *semantic redundancy* is generally not detrimental to the desired information. It is possible that a model can result in low response uncertainty but poor response quality (Singh et al., 2023; Li et al., 2024; Fu et al., 2025).

## 4 Experiments

To validate our proposed prompt-response concept model, we empirically demonstrate different aspects of our proposed model in different settings, the details of which are as follows. For instruction-fine-tuned LLMs, their prompts usually are in the form of some tasks from the user (i.e., ‘explain to me why the sky is blue’). Our experiments treat a relatively simple task as a ‘concept’ and a complex task as a composition of multiple concepts. All

<sup>1</sup>It is termed the *semantic entropy* in Kuhn et al. (2023). In this paper, we study it through the lens of uncertainty reduction.

<sup>2</sup>The response style can be viewed as an implicit concept, so *semantic redundancy* can be reduced by providing relevant style information in the prompt to guide the model response.

selected datasets have deterministic ground truths to eliminate the aleatoric uncertainty in the tasks.

### 4.1 Informative Prompt vs. Response Quality

It is worth noting that low uncertainty in model responses does not necessarily indicate high response quality, as an LLM can produce outputs with very low uncertainty while being blindly confident in incorrect answers. This behavior is problematic and can lead to hallucinations (Huang et al., 2023b). To ascertain the actual relationship between prompt, model response uncertainty, and quality, we further investigated the relationship between the effective token count of the prompt and model response quality. To assess if the reduction in uncertainty translates to improved output quality, we test the model’s output accuracy when answering the multiple-choice questions (MCQs

**Datasets.** We selected 100 MCQs from the Medical Meadow MedQA dataset (Jin et al., 2021), a benchmark for medical question answering. Due to space constraints, we report additional results on general-domain reasoning datasets, ARC (Clark et al., 2018), OpenBookQA (Mihaylov et al., 2018), and RACE (Lai et al., 2017), in Sec. B.

**Controlling prompt informativeness.** We use three methods for controlling prompt informativeness: sentence-level removal, attribute-level removal, and token-level removal.

① **Sentence-level removal.** We iteratively select an increasing fraction of randomly selected sentences from the context of the questions.

② **Attribute-level removal.** We first prompt the LLM to identify key concepts in the context (e.g., entities, events, and locations), along with their associated attributes, which are defined as specific facts or details that describe each concept. These attributes are extracted as exact quoted spans from the original text and then selectively removed. Unlike sentence-level removal, which indiscriminately deletes entire sentences, attribute-level removal targets semantically meaningful information that is critical for comprehension. This approach enables a more fine-grained analysis of model degradation when relevant facts are missing, rather than merely reducing the amount of text.

③ **Token-level removal.** We iteratively select an increasing fraction of randomly selected tokens in the context by replacing them with space tokens. For each question, we set the decoding temperature to 1 and sample 100 responses to assess output

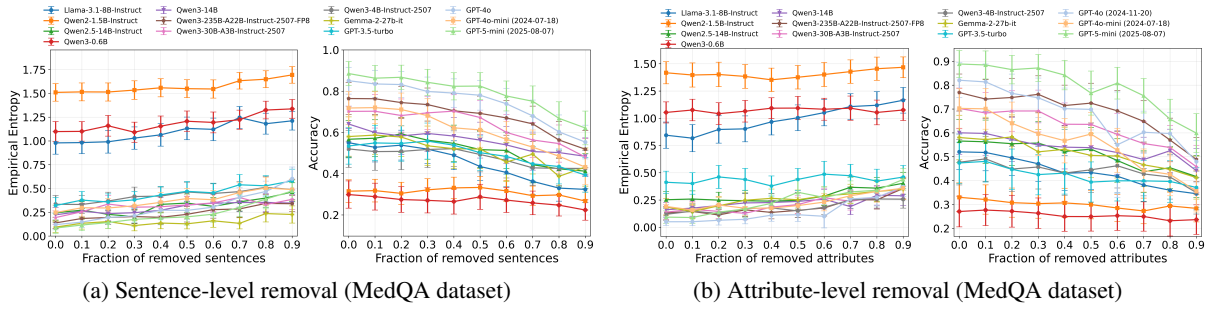


Figure 2: Accuracy and entropy for MedQA dataset after sentence- and attribute-level removal. There is strong negative correlation between accuracy and uncertainty, with less accurate models generally showing greater uncertainty in their responses.

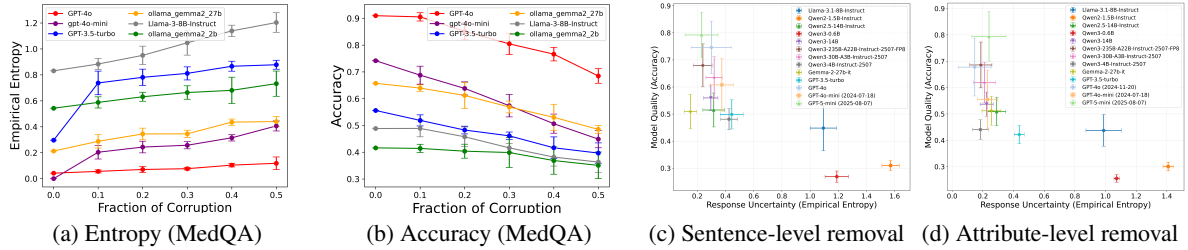


Figure 3: **Left two figures.** Accuracy and entropy for MedQA dataset after token-level removal. **Right two figures.** Model response quality (i.e., accuracy) vs. uncertainty (i.e., empirical entropy) for MedQA dataset after sentence- and attribute-level removal.

variability. We repeat this process using 5 different random seeds to select the tokens for replacement.

As the fraction of removed content increased, we incrementally expanded the set of randomly selected elements (sentences, attributes, or tokens), while keeping previously removed elements fixed. This controlled procedure ensures that changes in model performance are attributable to the increasing degree of information removal, rather than randomness in the selection process. As a result, we can more reliably assess the impact of reduced prompt informativeness on response quality.

**Results.** In Figs. 2, 3a, 3b and 5, we plot the accuracy for the responses of different open-source and black-box LLMs on the same set of MCQs. As the fraction of removed content in the prompt increases, the accuracy monotonically decreases across all tested models and removal strategies. For each random seed, we also compute the empirical conditional entropy  $H(Y|X)$  of the model responses for the given questions<sup>3</sup> (Figs. 2, 3a and 5) as a measure of response uncertainty. As more information is removed, response uncertainty consistently increases for all models (and monotonically for

<sup>3</sup>We assume a uniform distribution over questions, i.e.,  $p(x) = \frac{1}{100}$ . In the absence of access to the true conditional distribution  $p(y|x)$ , we estimate  $H(Y|X) = -\sum_x p(x) \sum_y \hat{p}(y|x) \log \hat{p}(y|x)$ , where  $\hat{p}(y|x)$  is obtained from the empirical response distribution. This metric is particularly suitable for MCQ settings, where the model’s effective output is a single discrete choice.

larger LLMs), revealing a clear negative correlation between  $H(Y|X)$  and response accuracy. These results support our hypothesis that greater prompt informativeness reduces response uncertainty and improves output quality. Furthermore, as shown in Figs. 3c, 3d, 7a and 7b, models with higher accuracy exhibit lower empirical conditional entropy, corroborating our characterization of  $\mathcal{E}$  in Theorem 1.

An interesting finding across all tasks is that sufficient prompt informativeness can allow a smaller model to match or even outperform a more capable one, e.g., Gemma2 27B model with full context surpassed GPT-4o-mini with upto 80% context. This result underscores that a high-quality prompt can effectively compensate for a model’s inherent scale, enabling less powerful models to achieve highly competitive performance.

## 4.2 mHealth Intervention Usecase

We now demonstrate the effectiveness of our proposed approach in a real-world simulation use case in mHealth setting. We adapt the formulation from Shin et al. (2022); both the app and the user act as reinforcement learning agents. The app agent’s objective is to encourage the user agent to adhere to the PT routine. The user agent moves along a chain with  $N$  states, where a higher state number represents a healthier physical state, and state  $N$  indicates completion of the PT routine.

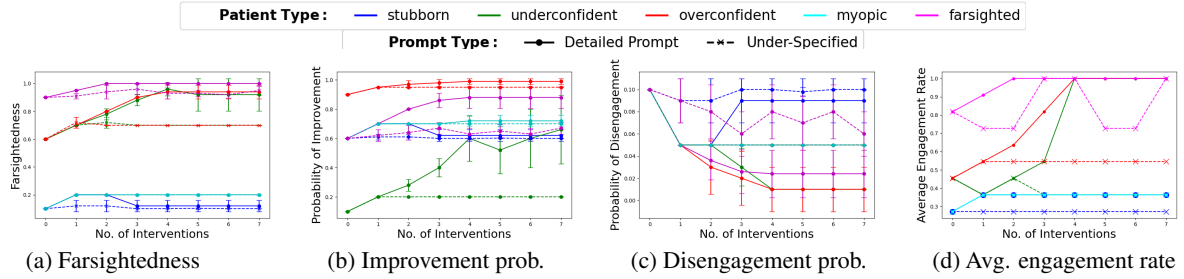


Figure 4: Results from PT intervention simulation. For (a) and (b), a higher value indicates more improvement. For (c), a lower value indicates more improvement. (d) is the patient’s optimal policy averaged across all health states, obtained from analytically solving the MDP. A higher value indicates on average the patient agent is more likely to continue engaging in PT. Overall, we observed prompt with more information gave rise to more consistent improvement compared to prompt with less information across all patient types.

We conduct the intervention simulation experiment with LLM to compare the effect of prompts with different informativeness levels on the intervention outcome. The experiment concludes that when the prompt provides the LLM (i.e., the app agent) with more information about the patient’s intentions and the strategies it can employ, the efficiency of the intervention improves consistently for different patient types compared to scenarios without the additional information. A more detailed description of this experiment is given in Sec. B.4.

*Due to space constraints, we have provided additional experimental results in Sec. B.*

## 5 Related Work

While uncertainty quantification has been extensively studied in machine learning, its application to LLMs remains relatively underexplored. [Kadavath et al. \(2022\)](#) investigated the extent to which LLMs can accurately self-evaluate their knowledge and how LLM calibration can improve response quality. Their goal of calibrating LLMs was to align the variability in model responses with actual uncertainty so that response variations genuinely reflect the model’s lack of relevant knowledge given the prompt. [Kuhn et al. \(2023\)](#) introduced the concept of semantic entropy to more precisely quantify the uncertainty in the informational content of model responses, accounting for semantically equivalent variations and eliminating noise from paraphrasing. [Lin et al. \(2023\)](#) proposed a method for estimating uncertainty in black-box LLMs for question-answering tasks by measuring response similarity using a Natural Language Inference (NLI) model, combined with simple dispersion-based measures. In a related but orthogonal line of work, [Wagle et al. \(2023\)](#) conducted empirical studies on pretrained language models (PLMs) and found that, while PLMs are

often overconfident, larger models tend to be better calibrated, i.e., confidence estimates more closely aligned to actual prediction accuracy. Similar to our work, [Ling et al. \(2024\)](#) aim to understand and quantify LLM response uncertainty by decomposing it into aleatoric and epistemic components. However, their study is limited to the setting of ICL and assumes a correlation between model response accuracy and uncertainty without direct empirical examination. In contrast, we explicitly investigate whether lower response uncertainty necessarily implies higher response quality. While we also adopt an entropy-based uncertainty measure, similar to [Kuhn et al. \(2023\)](#), [Lin et al. \(2023\)](#), [Wagle et al. \(2023\)](#), and [Farquhar et al. \(2024\)](#), our focus is on understanding the LLM’s response uncertainty, and how increasing informativeness and model quality can be used as a principled way to reduce response uncertainty. We defer additional related works on asymptotic behaviors of LLMs to Sec. A.1.

## 6 Conclusion

This paper highlights the significance of understanding how prompts and models impact response uncertainty in LLMs. By examining the informativeness of prompts, we demonstrate that more informative prompts, when combined with a better model, result in lower response uncertainty. To formalize this relationship, we introduce the prompt-response concept (PRC) model that captures how LLMs generate responses from prompts and helps identify sources of uncertainty. The insights from this work provide a principled approach for enhancing prompt design, which is crucial for the safe and effective use of LLMs in decision-making tasks, particularly in high-stakes domains such as healthcare. Future research may further refine the PRC model and explore its applicability in other areas that require reliable LLM responses.

## 7 Limitations

This work has two key limitations, which we outline in this section to guide future research.

**Idealistic nature of the PRC model.** It is worth noting that the PRC model that we proposed in this paper assumes an idealized version of LLMs. As empirically demonstrated, while models such as GPT-3.5-Turbo, GPT-4 and Llama-2, and Llama 3 exhibit behaviors largely according to our predictions, there are still some modes in which they deviate (e.g., Qwen2\_1.5b plot). This is likely in those cases where LLM does not know the mapping perfectly. For example, [Lu et al. \(2021\)](#) showed that the order of exemplars in ICL influences the model response quality. Our model does not capture this phenomenon. However, the authors showed that in the same work, the order of examples tends to have less effect as model quality gets better. Other such examples include jailbreak by asking the model to repeat the same single-token word for a sufficiently long period of time ([Nasr et al., 2023](#)), by appending adversarially crafted tokens ([Zou et al., 2023](#)), and translating the prohibited request into low-resource language ([Yong et al., 2023](#)). Similarly, it was observed that adversarial attacks tend to have lower success rates as the model becomes more capable. While further investigation is needed to incorporate the adversarial behavior of LLMs into this framework, the more capable LLMs are less prone to these failure modes. Our model can more effectively explain them.

**LLMs for human behavior simulation.** Research exploring the parallels between human behavior and reasoning patterns and those of LLMs, as well as the adaptation of LLMs as human substitutes in diverse studies, is detailed in [Aher et al. \(2023\)](#), [Argyle et al. \(2023\)](#), [Binz and Schulz \(2023\)](#), and [Dasgupta et al. \(2022\)](#). These studies frequently demonstrate LLMs’ capacity for human-like responses, leading many to regard them as viable alternatives. This paper, however, needs to delve into the appropriateness of this substitution, deferring to other works for such discussion.

## References

Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. In *Proc. ICML*, pages 337–371.

Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis*, pages 337–351.

Konstantine Arkoudas. 2023. Gpt-4 can’t reason. *arXiv:2308.03762*.

Marcel Binz and Eric Schulz. 2023. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, page e2218523120.

Omar Boursalie, Reza Samavi, and Thomas E Doyle. 2018. Machine learning and mobile health monitoring platforms: a case study on research and implementation challenges. *Journal of Healthcare Informatics Research*, pages 179–203.

Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv:1409.1259*.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*.

Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran, James L McClelland, and Felix Hill. 2022. Language models show human-like content effects on reasoning. *arXiv:2207.07051*.

Armen Der Kiureghian and Ove Ditlevsen. 2009. Aleatory or epistemic? does it matter? *Structural safety*, pages 105–112.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey on in-context learning. *arXiv:2301.00234*.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. *arXiv:1805.04833*.

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, pages 625–630.

AJ Ferrer-Riquelme. 2009. Statistical control of measures and processes. In *Comprehensive Chemometrics*. Elsevier, Amsterdam.

Jerry A Fodor. 1998. *Concepts: Where cognitive science went wrong*. Oxford University Press, Oxford.

Tairan Fu, Javier Conde, Gonzalo Martínez, María Grandury, and Pedro Reviriego. 2025. Multiple choice questions: Reasoning makes large language models (llms) more self-confident even when they are wrong. *arXiv:2501.09775*.

801	Deep Ganguli, Danny Hernandez, Liane Lovitt,	Saurav Kadavath, Tom Conerly, Amanda Askill, Tom	855
802	Amanda Askill, Yuntao Bai, Anna Chen, Tom Conerly,	Henighan, Dawn Drain, Ethan Perez, Nicholas	856
803	Nova Dassarma, Dawn Drain, Nelson Elhage,	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	857
804	and 1 others. 2022. Predictability and surprise in	Tran-Johnson, and 1 others. 2022. Language models	858
805	large generative models. In <i>Proc. ACM FAccT</i> , pages	(mostly) know what they know. <i>arXiv:2207.05221</i> .	859
806	1747–1764.		
807	Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel	Mert Karabacak and Konstantinos Margetis. 2023. Em-	860
808	Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan	bracing large language models for medical applica-	861
809	Leike, and Jeffrey Wu. 2024. Scaling and evaluating	tions: Opportunities and challenges. <i>Cureus</i> .	862
810	sparse autoencoders. <i>arXiv:2406.04093</i> .		
811	Google. 2023. PaLM 2 Technical Report.	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	863
812	<i>arXiv:2305.10403</i> .	taka Matsuo, and Yusuke Iwasawa. 2022. Large	864
813	Michael Hahn and Navin Goyal. 2023. A theory of	language models are zero-shot reasoners. In <i>Proc.</i>	865
814	emergent in-context learning as implicit structure	<i>NeurIPS</i> , pages 22199–22213.	866
815	induction. <i>arXiv:2303.07971</i> .		
816	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.	867
817	Distilling the knowledge in a neural network.	Semantic uncertainty: Linguistic invariances for un-	868
818	<i>arXiv:1503.02531</i> .	certainty estimation in natural language generation.	869
819	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and	<i>arXiv:2302.09664</i> .	870
820	Yejin Choi. 2019. The curious case of neural text	Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor Ion	871
821	degeneration. <i>arXiv:1904.09751</i> .	Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym	872
822	Jie Huang, Xinyun Chen, Swaroop Mishra,	Korablyov, and Yoshua Bengio. 2021. Deup: Direct	873
823	Huaixiu Steven Zheng, Adams Wei Yu, Xiny-	epistemic uncertainty prediction. <i>arXiv:2102.08501</i> .	874
824	ing Song, and Denny Zhou. 2023a. Large	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and	875
825	language models cannot self-correct reasoning yet.	Eduard Hovy. 2017. RACE: Large-scale ReAding	876
826	<i>arXiv:2310.01798</i> .	comprehension dataset from examinations. In <i>Proc.</i>	877
827	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,	<i>EMNLP</i> , pages 785–794.	878
828	Zhangyin Feng, Haotian Wang, Qianglong Chen,	Stephen Laurence and Eric Margolis. 1999. <i>Concepts</i>	879
829	Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 oth-	and <i>Cognitive Science</i> , pages 3–81. MIT Press, Cam-	880
830	ers. 2023b. A survey on hallucination in large lan-	bridge, MA.	881
831	guage models: Principles, taxonomy, challenges, and	Moxin Li, Wenjie Wang, Fuli Feng, Fengbin Zhu, Qi-	882
832	open questions. <i>arXiv:2311.05232</i> .	fan Wang, and Tat-Seng Chua. 2024. Think twice	883
833	Yuheng Huang, Jiayang Song, Zhijie Wang, Huaming	before assure: Confidence estimation for large lan-	884
834	Chen, and Lei Ma. 2023c. Look before you leap:	guage models through reflection on multiple answers.	885
835	An exploratory study of uncertainty measurement for	<i>arXiv:2403.09972</i> .	886
836	large language models. <i>arXiv:2307.10236</i> .		
837	Eyke Hüllermeier and Willem Waegeman. 2021.	Tom Lieberum, Senthoran Rajamanoharan, Arthur	887
838	Aleatoric and epistemic uncertainty in machine learn-	Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant	888
839	ing: An introduction to concepts and methods. <i>Ma-</i>	Varma, János Kramár, Anca Dragan, Rohin Shah,	889
840	<i>chine Learning</i> , pages 457–506.	and Neel Nanda. 2024. Gemma scope: Open sparse	890
841	James R Hurford, Brendan Heasley, and Michael B	autoencoders everywhere all at once on gemma 2.	891
842	Smith. 2007. <i>Semantics: a coursebook</i> . Cambridge	<i>arXiv:2408.05147</i> .	892
843	University Press, Cambridge.	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023.	893
844	Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi	Generating with confidence: Uncertainty quan-	894
845	Rungta, Krithika Iyer, Yuning Mao, Michael	tification for black-box large language models.	895
846	Tontchev, Qing Hu, Brian Fuller, Davide Testug-	<i>arXiv:2305.19187</i> .	896
847	ine, and 1 others. 2023. Llama guard: Llm-based	Chen Ling, Xujiang Zhao, Xuchao Zhang, Wei Cheng,	897
848	input-output safeguard for human-ai conversations.	Yanchi Liu, Yiyong Sun, Mika Oishi, Takao Osaki,	898
849	<i>arXiv:2312.06674</i> .	Katsushi Matsuda, Jie Ji, and 1 others. 2024. Uncer-	899
850	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,	tainty quantification for in-context learning of large	900
851	Hanyi Fang, and Peter Szolovits. 2021. What disease	language models. In <i>Proc. NAACL HLT</i> , pages 3357–	901
852	does this patient have? a large-scale open domain	3370.	902
853	question answering dataset from medical exams. <i>Ap-</i>	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	903
854	<i>plied Sciences</i> , 11(14):6421.	Hiroaki Hayashi, and Graham Neubig. 2023. Pre-	904
		train, prompt, and predict: A systematic survey of	905
		prompting methods in natural language processing.	906
		<i>ACM Computing Surveys</i> , pages 1–35.	907

908	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. <i>arXiv:2104.08786</i> .	962
909		963
910		964
911		
912	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In <i>Proc. EMNLP</i> .	965
913		966
914		967
915		968
916	Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. Scalable extraction of training data from (production) language models. <i>arXiv:2311.17035</i> .	969
917		
918		
919		
920		
921		
922	OpenAI. 2023. GPT-4 Technical Report. <i>arXiv:2303.08774</i> .	970
923		971
924	Alina Petukhova, Joao P Matos-Carvalho, and Nuno Fachada. 2024. Text clustering with llm embeddings. <i>arXiv:2403.15112</i> .	972
925		973
926		974
927	Gualtiero Piccinini and Sam Scott. 2006. Splitting concepts. <i>Philosophy of Science</i> , pages 390–409.	975
928		976
929	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In <i>Proc. EMNLP</i> , pages 2383–2392.	977
930		978
931		
932		
933	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. <i>arXiv:2402.07927</i> .	979
934		980
935		981
936		982
937		983
938	Malik Sallam. 2023. Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In <i>Healthcare</i> , page 887.	984
939		985
940		986
941		987
942	Robin Senge, Stefan Bösner, Krzysztof Dembczyński, Jörg Haasenritter, Oliver Hirsch, Norbert Donner-Banzhoff, and Eyke Hüllermeier. 2014. Reliable classification: Learning classifiers that distinguish aleatoric and epistemic uncertainty. <i>Information Sciences</i> , pages 16–29.	988
943		
944		
945		
946		
947		
948	Mohammad Hossein Shaker and Eyke Hüllermeier. 2020. Aleatoric and epistemic uncertainty with random forests. In <i>Proc. IDA</i> , pages 444–456.	989
949		990
950		991
951	Eura Shin, Siddharth Swaroop, Weiwei Pan, Susan Murphy, and Finale Doshi-Velez. 2022. Modeling mobile health users as reinforcement learning agents. <i>arXiv:2212.00863</i> .	992
952		993
953		994
954		995
955	Aniket Kumar Singh, Suman Devkota, Bishal Lamichhane, Uttam Dhakal, and Chandra Dhakal. 2023. The confidence-competence gap in large language models: A cognitive study. <i>arXiv:2309.16145</i> .	996
956		997
957		998
958		999
959	Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In <i>Proc. NeurIPS</i> .	1000
960		1001
961		1002
		1003
		1004
		1005
		1006
		1007
		1008
		1009
		1010
		1011
		1012
		1013
		1014
		1015
		1016
		1017
		1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047
		1048
		1049
		1050
		1051
		1052
		1053
		1054
		1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
		1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		137

1015 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and  
1016 Tengyu Ma. 2022. An explanation of in-context learn-  
1017 ing as implicit bayesian inference. In *Proc. ICLR*.

1018 Rui Yang, Ting Fang Tan, Wei Lu, Arun James  
1019 Thirunavukarasu, Daniel Shu Wei Ting, and Nan  
1020 Liu. 2023. Large language models in health care:  
1021 Development, applications, and challenges. *Health  
1022 Care Science*, pages 255–263.

1023 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak  
1024 Shafran, Karthik Narasimhan, and Yuan Cao. 2023.  
1025 React: Synergizing reasoning and acting in language  
1026 models. In *Proc. ICLR*.

1027 Zheng-Xin Yong, Cristina Menghini, and Stephen H  
1028 Bach. 2023. Low-resource languages jailbreak gpt-4.  
1029 *arXiv:2310.02446*.

1030 Dell Zhang, Alina Petrova, Dietrich Trautmann, and  
1031 Frank Schilder. 2023a. Unleashing the power of  
1032 large language models for legal applications. In *Proc.  
1033 ACM CIKM*, pages 5257–5258.

1034 Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang,  
1035 and Zhaoran Wang. 2023b. What and how  
1036 does in-context learning learn? bayesian model  
1037 averaging, parameterization, and generalization.  
1038 *arXiv:2305.19420*.

1039 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,  
1040 Xiaolei Wang, Yupeng Hou, Yingqian Min, Be-  
1041 ichen Zhang, Junjie Zhang, Zican Dong, and 1 oth-  
1042 ers. 2023. A Survey of Large Language Models.  
1043 *arXiv:2303.18223*.

1044 Andy Zou, Zifan Wang, J Zico Kolter, and Matt  
1045 Fredrikson. 2023. Universal and transferable  
1046 adversarial attacks on aligned language models.  
1047 *arXiv:2307.15043*.

## A Leftover Details

### A.1 Related Works

**Explanation for asymptotic behaviors of LLMs.** Several recent efforts have been made to develop frameworks that explain the surprising emergent behaviors of LLMs. Zhang et al. (2023b) demonstrate that the attention mechanism in LLMs approximates Bayesian model averaging in the ICL setting. Wang et al. (2023) conceptualize real-world LLMs as latent variable models, suggesting that these models operate as implicit topic models by inferring latent conceptual variables from prompts. Notably, Xie et al. (2022) interpret ICL as an instance of implicit Bayesian inference over latent concepts acquired during pretraining. However, their analysis is limited to characterizing zero-one error in the asymptotic case of infinite exemplars, and their hidden Markov model-based mathematical formulation is tailored specifically to the ICL setting, making it unsuitable for analyzing chain-of-thought or conversational-style responses. Furthermore, their theoretical results quantify the mode of the posterior predictive distribution and do not address uncertainty quantification. Hahn and Goyal (2023) further extend this line of work by incorporating greater flexibility and complexity in the exemplars, yet similarly provide only asymptotic bounds on classification error. In contrast, our work complements these approaches by focusing on how the response uncertainty varies with finite-length prompts, providing a more practical understanding of LLM behavior in real-world scenarios.

### A.2 Proofs

**Lemma 1.** For any LLM with exactly learned  $g_x$  and two concepts  $\theta_1, \theta_2 \in \Theta_x$ , we have  $\mathcal{X}_{\theta_1} \cap \mathcal{X}_{\theta_2} = \emptyset$  if  $\theta_1 \neq \theta_2$ . Furthermore,  $H(Z_x|X_{\theta_x}) = 0$ .

*Proof.* The result holds trivially for the case in which  $\mathcal{X}_{\theta}$  for any  $\theta \in \Theta_x$  is an empty set. As discussed in Sec. 3, each concept is completely characterized by all of its attributes, therefore, two different concepts cannot have the same set of attributes, i.e.,  $\mathcal{A}_{\theta_i} \neq \mathcal{A}_{\theta_j}$  if  $i \neq j$ . As our PRC model assumes any attribute can be perfectly expressed by some sequence of tokens, any attribute  $a_{\theta_i,k} \in \mathcal{A}_{\theta_i}$  can be expressed by a sequence of tokens. We denote the set of all possible sequence of tokens by  $\mathcal{X}_{s(a_{\theta_i,k})}$ , where  $s(a_{\theta_i,k})$  denotes the semantic meaning of  $a_{\theta_i,k}$ . Therefore, the set of attributes  $\mathcal{A}_{\theta_i}$  is expressed as a sequence of tokens  $X_{\theta_i} \in \mathcal{X}_{\theta_i}$ , where  $\mathcal{X}_{\theta_i} = \mathcal{C} \left( \left\{ \mathcal{X}_{s(a_{\theta_i,k})} \right\}_{k=1}^n \right)$  in which  $n = |\mathcal{A}_{\theta_i}|$  and operator  $\mathcal{C}$  applied in the following way:

- Chooses one element  $x_{s(a_{\theta_i,k})} \in \mathcal{X}_{s(a_{\theta_i,k})}$  for each  $k \in \{1, 2, \dots, n\}$ ;
- Create a set  $\mathcal{S}_{\theta_i}$  containing all the selected elements  $x_{s(a_{\theta_i,k})}$ . Then, concatenate these elements in  $\mathcal{S}_{\theta_i}$  to form sequences by exhausting all possible ordering and use this collection of sequences to form a new set  $\mathcal{X}'_{\theta_i}$ .
- Repeat step 1 and 2 for all possible sets  $\mathcal{S}_{\theta_i}$  and generate all possible  $\mathcal{X}'_{\theta_i}$ . Finally, take the union of all such  $\mathcal{X}'_{\theta_i}$  sets to form a new set. Since this set consists of all possible sequences that are semantically equivalent and fully characterize  $\theta_i$ , it is exactly  $\mathcal{X}_{\theta_i}$ .

Intuitively, the operator  $\mathcal{C}$  takes all sequences that fully characterize each attribute of the concept  $\theta_i$  and generates all possible concatenated sequences that fully characterize concept  $\theta_i$ . Therefore, under the PRC model, for every  $\theta \in \Theta_x$ , there exists a non-empty set  $\mathcal{X}_{\theta}$ . Since the attributes of any two distinct concepts are different, i.e.,  $\mathcal{A}_{\theta_i} \neq \mathcal{A}_{\theta_j}$  if  $i \neq j$ ,  $X_{\theta_i} \neq X_{\theta_j}$  if  $i \neq j$ . Since  $\mathcal{X}_{\theta_i}$  is the support of  $X_{\theta_i}$ ,  $\mathcal{X}_{\theta_i} \cap \mathcal{X}_{\theta_j} = \emptyset$  if  $i \neq j$ .

Since the first part of Lemma 1 is non-trivially true in our framework, given any  $X_{\theta_x} = x$ , there exists a unique  $\theta_x \in \Theta_x$  such that  $p(Z_x = \theta_x|x) = 1$  and  $p(Z_x = \theta_x|x') = 0$  for all  $x' \neq x$ . Therefore,

$$H(Z_x|X_{\theta_x}) = - \sum_{x \in X_{\theta_x}} P(x) \sum_{z \in Z_x} P(z|x) \log P(z|x)$$

$$\begin{aligned}
&= - \sum_{x \in X_{\theta_x}} P(x) \left( \sum_{z \in Z_x \setminus \theta} P(z|x) \log P(z|x) + P(\theta|x) \log P(\theta|x) \right) \\
&= - \sum_{x \in X_{\theta_x}} P(x) \left( \sum_{z \in Z_x \setminus \theta} 0 \log 0 + 1 \log 1 \right) \\
&= - \sum_{x \in X_{\theta_x}} P(x)(0) = 0
\end{aligned}$$

Note that in order for the model to get the correct conditional entropy above, it must know the true mapping function  $g_x$ . This is because it needs to be able to tell that  $p(Z_x = \theta_x|x) = 1$  and  $p(Z_x = \theta_x|x') = 0$  for all  $x' \neq x$ . Therefore, our result holds when LLM has exactly learned  $g_x$ .  $\square$

**Proposition 1.** *As  $X_s$  represents more informative prompts (i.e., as more task-relevant information is included in the prompt),  $H(Z_x|X_s)$  decreases.*

*Proof.* Given Lemma 1, we know that  $Z_x$  depends on  $X_{\theta_x}$ . If there exists  $X_{\theta_x} \in \mathcal{X}_{\theta_x}$  such that  $\alpha_s \subset \alpha_{\theta_x}$ , then  $Z_x$  and  $X_s$  are dependent. Therefore,  $I(Z_x; X_s) > 0$ , and as a result

$$H(Z_x|X_s) = H(Z_x) - I(Z_x; X_s) < H(Z_x). \quad (2)$$

Let  $Z'_x$  denote the random variable formed by  $Z_x$  conditioning on  $X_s$ . Since  $\text{Supp}(Z'_x) \subseteq \text{Supp}(Z_x)$ , there still exist semantically meaningful prompts  $X''_s$  that is related to  $Z'_x$ . Apply Inequality (2) again we obtain:

$$H(Z_x|(X_s, X''_s)) = H(Z'_x|X'_s) < H(Z'_x) = H(Z_x|X = X_s) < H(Z_x), \quad (3)$$

where  $X'_s = (X_s, X''_s)$  is a longer prompt sequence formed by appending  $X''_s$  to  $X_s$ . By iteratively applying the inequality in Eq. (3), we finally obtain Proposition 1.  $\square$

**Theorem 1.** *As  $X_s$  represents more informative prompts,  $H(Y|X_s)$  strictly decreases. Specifically  $H(Y|X_s) \leq H(Y|Z_y) + H(g_c(Z_x)|Z_x) + H(Z_x|X_s)$ . Furthermore,  $H(Y|X_s)$  converges to  $H(Y_r|Z_y) + H(Y_r|Z_y) + \mathcal{E}$ , where  $\mathcal{E}$  decreases as the model quality improves. When Assumption 1 hold,  $\mathcal{E} \rightarrow 0$  as prompt informativeness increases.*

*Proof.* To simplify notation, we use  $Y$  instead of  $Y_{\theta_y}$  and assume the model response is complete (i.e., the last token is the ‘EOS’ token). By design,  $Z_x$  and  $Z_y$  are discrete random variables. Intuitively, it is easy to see why discretizing concepts is a reasonable way to model concepts. Since LLMs are trained with texts that are discrete, it is not feasible to interpolate between any two concepts with infinitesimally small step sizes with natural language as the medium. We consider a general setting, where  $g_c$  can be a stochastic function, i.e.,  $Z_y$  can have different realizations for the same  $Z_x$ . Since  $H(f(X)|Y) \leq H(f(X), X|Y)$ , we have

$$H(Z_y|X_s) \leq H(Z_y, Z_x|X_s) = H(Z_x|X_s) + H(Z_y|Z_x, X_s), \quad (4)$$

therefore,

$$\begin{aligned}
H(Z_y|X_s) &\leq H(Z_x|X_s) + H(Z_y|Z_x, X_s) \\
&= H(Z_x|X_s) + H(Z_y|Z_x),
\end{aligned}$$

as  $Z_y$  is conditionally independent of  $X_s$  given  $Z_x$ . Finally, since

$$\begin{aligned}
H(Y) &= H(Y, Z_y) - H(Z_y|Y) \\
&= H(Y|Z_y) + H(Z_y) - H(Z_y|Y),
\end{aligned}$$

we can express the entropy of the response posterior as follows:

$$\begin{aligned}
\mathbb{H}(Y|X_s) &= \mathbb{H}(Y|Z_y, X_s) + \mathbb{H}(Z_y|X_s) - \mathbb{H}(Z_y|Y, X_s) \\
&= \mathbb{H}(Y|Z_y) + \mathbb{H}(Z_y|X_s) - \mathbb{H}(Z_y|Y, X_s) \\
&\quad (Y \text{ is conditionally independent of } X_s \text{ given } Z_y). \\
&\leq \mathbb{H}(Y|Z_y) + \mathbb{H}(Z_y|X_s) \\
&\leq \mathbb{H}(Y|Z_y) + \mathbb{H}(Z_x|X_s) + \mathbb{H}(Z_y|Z_x, X_s) \\
&\quad (\text{from Eq. (4) and } Z_y \text{ is conditionally independent of } X_s \text{ given } Z_x) \\
\implies \mathbb{H}(Y|X_s) &\leq \mathbb{H}(Y|Z_y) + \mathbb{H}(g_c(Z_x)|Z_x) + \mathbb{H}(Z_x|X_s). \quad (\text{as } Z_y = g_c(Z_x))
\end{aligned}$$

This completes the proof of the first part. Furthermore, we can rewrite  $\mathbb{H}(Y|X_s) = \mathbb{H}(Y|Z_y) + \mathcal{E}$ . The term  $\mathcal{E}$  captures the stochasticity of  $g_c$  which is upper bounded by  $\mathbb{H}(g_c(Z_x)|Z_x) + \mathbb{H}(Z_x|X_s)$ , and it depends the quality of the LLM. In our experiments, we observe that better LLMs (i.e., with higher accuracy for the same given prompts) have smaller empirical entropy values (Figs. 3c and 3d). With the PRC model, we interpret this as because better LLMs learned a better  $g_c$  during their training, such that the  $\mathcal{E}$  is smaller. Ideally, when the LLM is perfectly trained and the task in the prompt has a deterministic ground truth,  $g_c$  is deterministic, the term  $\mathcal{E}$  vanishes, and only *semantic redundancy* and *task uncertainty* remain. Therefore, we can rewrite  $\mathbb{H}(Y|Z_y) = \mathbb{H}(Y_\tau|Z_y) + \mathbb{H}(Y_r|Z_y)$  to characterize the randomness due to the inherent uncertainty in the task (e.g., asking for the outcome of an instance of tossing a fair coin) and the semantic redundancy in model responses because there are multiple ways of expressing the same concept (i.e., semantically equivalent sequences). This result holds for the LLM that knows the true  $g_y$ . Due to Proposition 1 and Lemma 1, when  $X_s$  has enough information to perfectly characterize the concept (i.e.,  $X_s \in \mathcal{X}_{\theta_x}$ ),  $\mathbb{H}(Z_x|X_s) = \mathbb{H}(Z_x|X_{\theta_x}) = 0$  for the LLM with known  $g_x$ . Hence,  $\mathbb{H}(Z_y|X_s)$  reduces to  $\mathbb{H}(Z_y|Z_x)$  (since it is upper bounded by  $\mathbb{H}(Z_x|X_s) + \mathbb{H}(Z_y|Z_x)$ ) and  $\mathcal{E}$  reduces to a non-negative value that is no larger than  $\mathbb{H}(g_c(Z_x)|Z_x)$ , which further decreasing as model quality improves. As a result, the value of  $\mathcal{E}$  reduces as the model quality increases.  $\square$

### A.3 Discussion about Response Uncertainty Decomposition

In this section, we examine the decomposition of response uncertainty. As discussed in Sec. 3.2, each component of uncertainty captures a distinct source of variation:

- $\mathbb{H}(Z_x|X_s)$  reflects uncertainty due to prompt underspecification, specifically, the model’s inability to infer the intended concept from the prompt, either because of incomplete information or limitations in its inherent knowledge.
- $\mathbb{H}(g_c(Z_x)|Z_x)$  captures uncertainty arising from model quality, i.e., the model’s ability to generate the appropriate high-level abstractions (response concepts) based on its understanding of the instruction. This uncertainty depends on how well the model interprets the prompt and produces a suitable response, overall depending on model’s quality.
- $\mathbb{H}(Y_\tau|Z_y)$  represents task-related variability, where randomness is inherent to the task itself (e.g., predicting the outcome of a fair coin toss or rolling a die).
- $\mathbb{H}(Y_r|Z_y)$  accounts for semantic redundancy, reflecting the existence of multiple ways to express the same underlying response concept (i.e., different sequences that are semantically equivalent sentences).

Together, these components provide a structured view of the different sources of uncertainty in language model outputs.

### A.4 Relationship between Prompts and Response

Having established the PRC model, we now seek to formalize it mathematically to analyze the response uncertainty. As stated before, we now analyze the probability of a *response-level*, complete response  $y$

given the initial prompt  $x$ . Following prior work (Xie et al., 2022; Hahn and Goyal, 2023; Zhang et al., 2023b; Wang et al., 2023), which assumes that LLMs implicitly perform Bayesian inference, we also adopt this perspective. Specifically, we assume that the conditional distribution  $p(y|x)$  representing the posterior predictive distribution is marginalized over the latent prompt and response concepts, and is

$$p(y|x) = \int_{\theta_y} p(y|\theta_y, x)p(\theta_y|x)d\theta_y = \int_{\theta_y} \int_{\theta_x} p(y|\theta_y, x)p(\theta_y|\theta_x, x)p(\theta_x|x)d\theta_y d\theta_x.$$

The first equality follows from conditioning the response on the response concept. Critically, the term  $p(y|\theta_y, x)$  represents the probability of the response  $y$  generated from the response concept  $\theta_y$ . The second equality follows from

$$p(\theta_y|x) = \int_{\theta_x} p(\theta_y|\theta_x, x)p(\theta_x|x)d\theta_x.$$

which is conditioning distribution of the response concept given the prompt concept. If  $p(\theta_c|x)$  (where  $c = \{x, y\}$ ) concentrates on a specific concept with a more informative prompt, the LLM learns effectively via marginalization. More concretely, our PRC model assumes that the LLM achieves this by extracting task-relevant attributes from the prompt  $x$  and using its inherent knowledge acquired during pretraining.

## B Additional Experiment Details

We first give additional experiments on other datasets such as ARC (Clark et al., 2018), OpenBookQA (Mihaylov et al., 2018), and RACE (Lai et al., 2017), comparing the accuracy and empirical entropy of generated responses across different LLMs under similar conditions as described in Sec. 4.1. We then show the ablation studies in Sec. B.4.1, demonstrating how response uncertainty varies with different types of noisy prompts. Finally, we describe the details about our mHealth intervention simulation experiments.

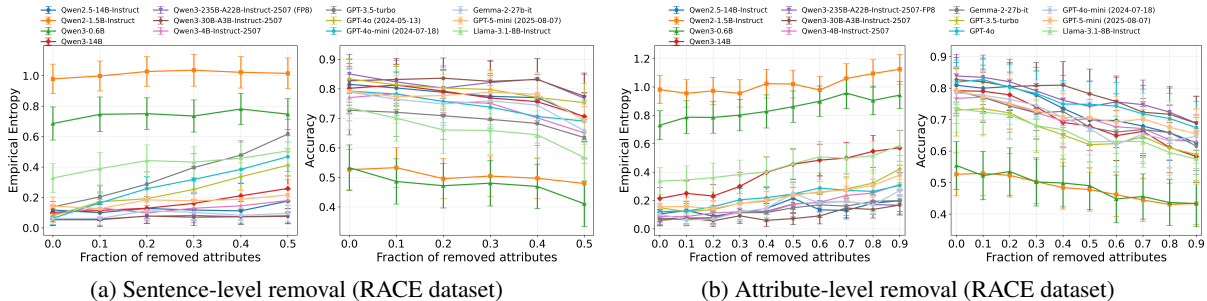


Figure 5: Accuracy and entropy for RACE dataset after sentence- and attribute-level removal. There is strong negative correlation between accuracy and uncertainty, with less accurate models generally showing greater uncertainty in their responses.

We also selected 100 MCQs from ARC (Clark et al., 2018) and OpenBookQA (Mihaylov et al., 2018) datasets. In Figs. 6b and 6d, we plot the accuracy for the responses of three open-source and three black-box LLMs on the same set of MCQs. As the fraction of masked tokens increases in prompt, the accuracy monotonically decreases for all tested models. For each random seed, we also plot the empirical conditional entropy  $H(Y|X)$  of the response for the given questions<sup>4</sup> (Figs. 6a and 6c) as an indicator of response uncertainty. As corruption becomes more severe, we observe that the response uncertainty increases for all models (increases monotonically for larger LLMs), indicating a clear negative correlation between  $H(Y|X)$  and the response accuracy. This result corroborates our hypothesis: more relevant information leads to both a reduction in response uncertainty and an improvement in its quality. Furthermore, as shown in Fig. 7c, models with better accuracy tend to have lower empirical entropy. This corroborates our characterization of  $\mathcal{E}$  in Theorem 1.

<sup>4</sup>The distribution of the questions used  $p(x)$  is assumed uniform. With no access to the prior of  $p(y|x)$ , we use the form  $H(Y|X) = -\sum_x p(x) \sum_y \hat{p}(y|x) \log \hat{p}(y|x)$  where  $\hat{p}(y|x)$  is obtained from the empirical distribution and  $p(x) = \frac{1}{100}$  for all  $x$  in the setting. The conditional entropy is a good measure for MCQs setting as the model’s effective response is just one choice.

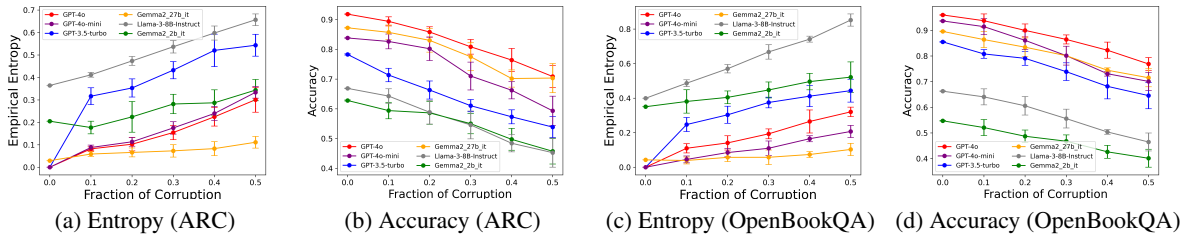


Figure 6: Accuracy and entropy for MCQs datasets. There is a clear and strong negative correlation between accuracy and uncertainty, with less accurate models generally showing greater uncertainty in their responses. For the ARC dataset, Gemma2 2B model with full context not only surpassed GPT-3.5 but also performed on par with GPT-4o-mini with only half the context.

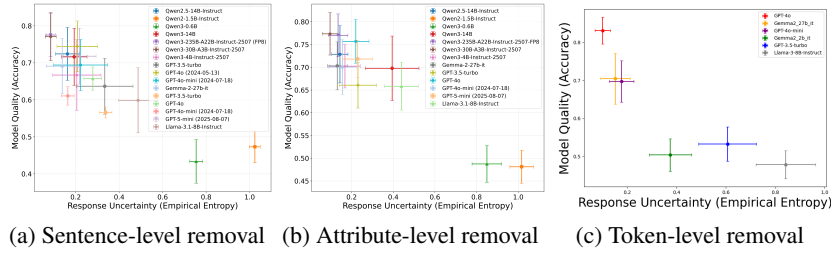


Figure 7: **Fig. (a) and Fig. (b):** Model response quality (i.e., accuracy) vs. uncertainty (i.e., empirical entropy) for RACE dataset after sentence- and attribute-level removal. **Fig. (c):** Model response quality (i.e., accuracy) vs. uncertainty (i.e., empirical entropy) for difference models. Averaged across Medical Meadow Medqa, ARC and OpenbookQA) and all corruption levels.

## B.1 The Relative Importance of Different Attributes

We investigate to what extent different attributes contribute to model response quality and uncertainty. We choose 10 questions from the RACE dataset (Lai et al., 2017) with moderate context length, assume each context as one concept and the sentences it comprises as its attributes, and use leave-one-out method to remove one sentence from its context, for each case generate 100 response samples and observe its impact on the model response. As shown in Fig. 8, we observed for some cases, there is a strong correlation between the choice of the removal of the sentence and the response quality and uncertainty across different models. This suggests that to some degree, there is a consensus among the LLMs about the importance of certain attributes in affecting model’s ability to find the right prompt and response concept.

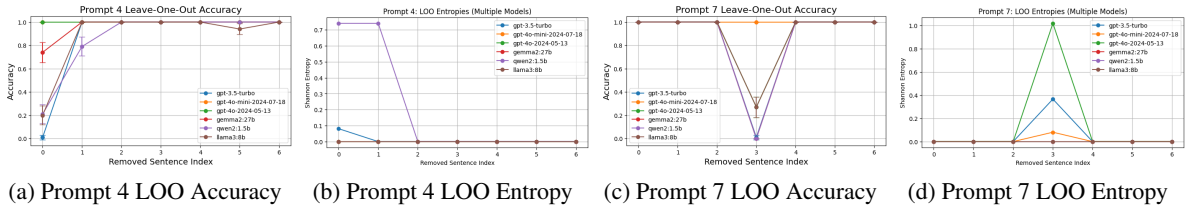


Figure 8: (a) Prompt 4 Leave-One-Out Accuracy. (b): Prompt 4 Leave-One-Out Empirical Entropy. (c): Prompt 7 Leave-One-Out Accuracy.(d): Prompt 7 Leave-One-Out Empirical Entropy. For all plots, the colour not visible has value 0. It can be observed that there is clear correlation between the choice of the sentence removal and the response quality/uncertainty , which indicates certain attributes are commonly important across multiple models.

## B.2 Response Quality vs. Model Precision

Fig. 9 illustrates the relationship between model precision and both accuracy and response uncertainty across multiple MCQ datasets (MedQA and RACE). We observe a clear and strong negative correlation between model precision and response quality: less precise models consistently exhibit higher empirical conditional entropy and lower accuracy, while more precise models produce more confident and accurate responses. This trend holds across datasets, suggesting that increased model precision is closely associated with reduced response uncertainty and improved decision quality in multiple-choice question answering.

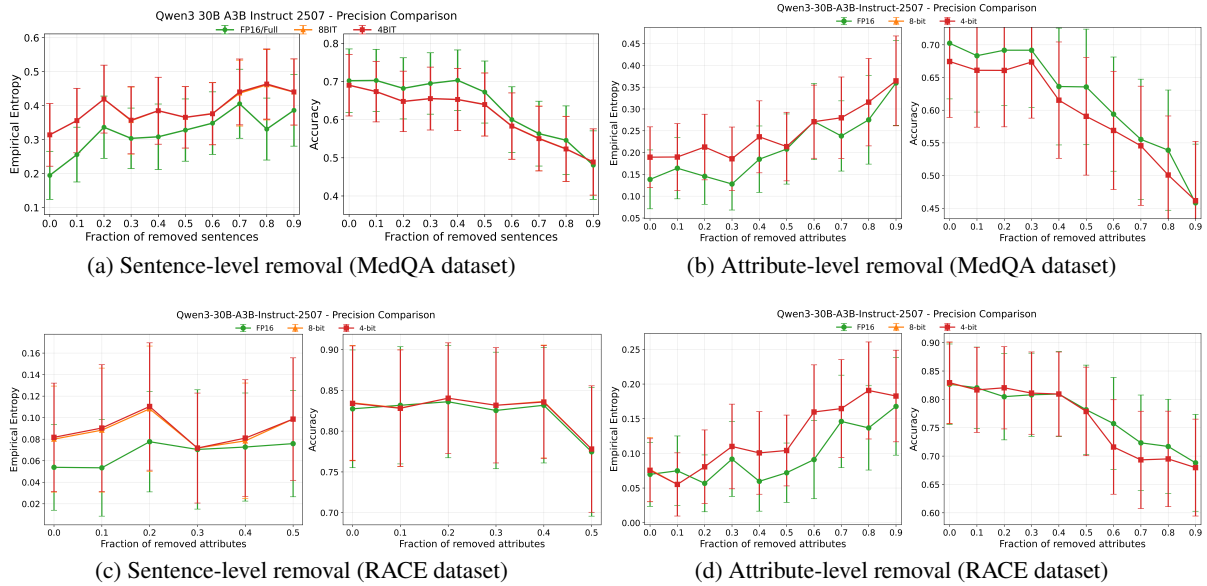


Figure 9: Accuracy and entropy for MCQs datasets vs. model precision. There is a clear and strong negative correlation between model precision and response quality, with less precise models generally showing greater uncertainty and lower accuracy in their responses.

### B.3 License for Datasets

Medical Meadow MedQA (Jin et al., 2021): MIT License; ARC (Clark et al., 2018): CC BY-SA 4.0; OpenBookQA (Mihaylov et al., 2018): Apache License 2.0; RACE (Lai et al., 2017): non-commercial research purpose only; SQuAD (Rajpurkar et al., 2016): CC BY-SA 4.0.

### B.4 Further Details on the mHealth Intervention Simulation Experiments in Sec. 4.2

At the beginning of the PT, the user is at state 0. The user has their default set of MDP parameters (i.e., discount factor  $\gamma$ , probability of transiting to the next healthier physical state  $p$ , and the probability of disengaging from PT  $d$ ). In this setting, those MDP parameters are interpreted in the following way:  $\gamma$  represents the farsightedness of the patient,  $p$  represents the probability of the patient’s health state gets improved if they chooses to engage in PT,  $d$  represents the probability of the patient disengaging from PT if they chooses to abstain from PT. Based on these parameters, the user agent can solve this MDP and figure out their optimal policy. The task of the app agent is to intervene on the user’s MDP parameters such that the optimal policy for the user is to complete the PT (i.e., go from state 0 to state  $N$ ).<sup>5</sup> We use the same formulation in this simulation by using two LLMs as the app agent and the user agent respectively. The app agent uses natural language to intervene in the user behavior. The user LLM is grounded in the aforementioned MDP setting. Specifically, in the system message for the user agent, the model is told they will increase the value of  $\gamma$  (i.e., farsightedness) when the app agent persuades the user agent to value more on the long-term goal of PT, increase  $p$  (i.e., probability of improvement) and decrease  $d$  (i.e., probability of disengagement) when the app agent manages to strengthen the user’s belief in the efficacy of PT. An illustration of the setup can be found in Fig. 10.

The effectiveness of the intervention depends on the following factors:

- The persuasiveness of and the strategy used by the app agent.
- The values of MDP parameters.
- The stubbornness of the user. The system message is defined in the way that a ‘stubborn’ user is less likely to change their behaviors compared to a ‘not-so-stubborn’ user.

<sup>5</sup>Refer to Shin et al. (2022) for the complete description of the problem setting and formulation.

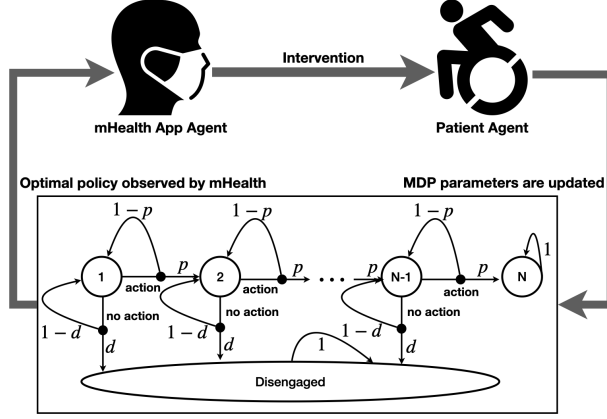


Figure 10: Visualization of states and transitions in the digital health grid world. Arrows indicate the required action and the probability of transitioning between states.

We conduct the intervention simulation experiment to compare the effect of different system messages for the app agent on the outcome of the intervention. The two system messages for comparison can be found in Sec. C.1.

We set  $N = 10$ . For each run, we give 7 rounds of conversation between the app agent and the user. While the history of the conversation between them is visible to both parties within every run, the user’s MDP parameters are not directly visible to the app agent. However, after every round of intervention, after the user updates their MDP parameters, a value iteration solver will be used to find the optimal policy of the patient, and this policy is visible to the app agent. The app agent can potentially leverage this piece of information to decide how to proceed with the next round of intervention. The user agent will also have the memory of this history in the change of their own MDP parameters. We use OpenAI ‘gpt-4-1106-preview’ API for both app agent and user and use 5 different random seeds for each different setting. We run the intervention experiments on 5 types of patients, each with a noticeably different set of initial MDP parameters from the rest. The exact values and details on the setup and can be found in Tab. 1. The results can be found in Fig. 4.

It can be observed across all settings, with more useful information provided in the system message, the MDP parameters were more likely to be changed in the positive direction (i.e., larger  $\gamma$  and  $p$ , smaller  $d$ ). As a result, the patient has improved PT engagement rate across all health states for all patient types. Moreover, this change tends to have a longer persistent effect compared to when the system message contains less useful information. This result is sensible because the more successful intervention came from an app agent who was provided with more information to work with. It has a better intervention strategy because its messages are tailored to specifically influence the user’s MDP parameters. Our proposed framework provides an information theoretic perspective to formalize this intuitive notion: when the system message with the longer prompt can specify the more relevant part of the concept in LLMs’ concept space and assuming the relevant knowledge is known, this prompt can provide consistent and useful responses due to its less posterior entropy which translates to more effective intervention strategy. As a result, the responses from the user are also more consistent and positive.

Patient Type \ MDP parameters	$\gamma$	$p$	$d$
	Underconfident	0.6	0.1
Overconfident	0.6	0.9	0.1
Myopic	0.1	0.6	0.1
Farsighted	0.9	0.6	0.1
Stubborn	0.1	0.6	0.1

Table 1: The initial MDP parameters values for every type of patient.

1268

### B.4.1 Ablations

1269

1270

1271

We also run a series of ablations to analyze the impact of various components of PRC model, prompt informativeness, compositionality, and irrelevant information, on response uncertainty. The details of which are as follows.

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

**Prompt informativeness and response uncertainty.** We first begin by assessing the response uncertainty of LLMs through the generation of responses using increasingly longer prompts with more relevant information (Sec. C.2 for the prompts used). For each prompt, we generate 100 responses from LLM with uncalibrated logits ( $T = 1$ ) and project them into the embedding space as single points using the OpenAI ‘text-embedding-ada-002’ model. To quantify the uncertainty in the generated responses for a given prompt, we use the *total standard deviation*, denoted as  $M(x)$ , defined as  $\sqrt{\text{Tr}(\Sigma)}$ , where  $\Sigma$  represents the covariance matrix of the embedding vectors of responses  $y_1, \dots, y_{100}$ . For LLMs, the dispersion of their responses in the embedding space indicates how much they differ in their semantic meaning (Lin et al., 2023; Petukhova et al., 2024). Therefore,  $M(x)$  is an effective metric for how much uncertainty there is in the model responses. It is noteworthy that  $\text{Tr}(\Sigma)$  is also referred to as *total variation*, serving as a lightweight measure of dispersion in the data (Ferrer-Riquelme, 2009). This metric is applicable for responses generated from both black-box and white-box LLMs, as it does not require access to logits.

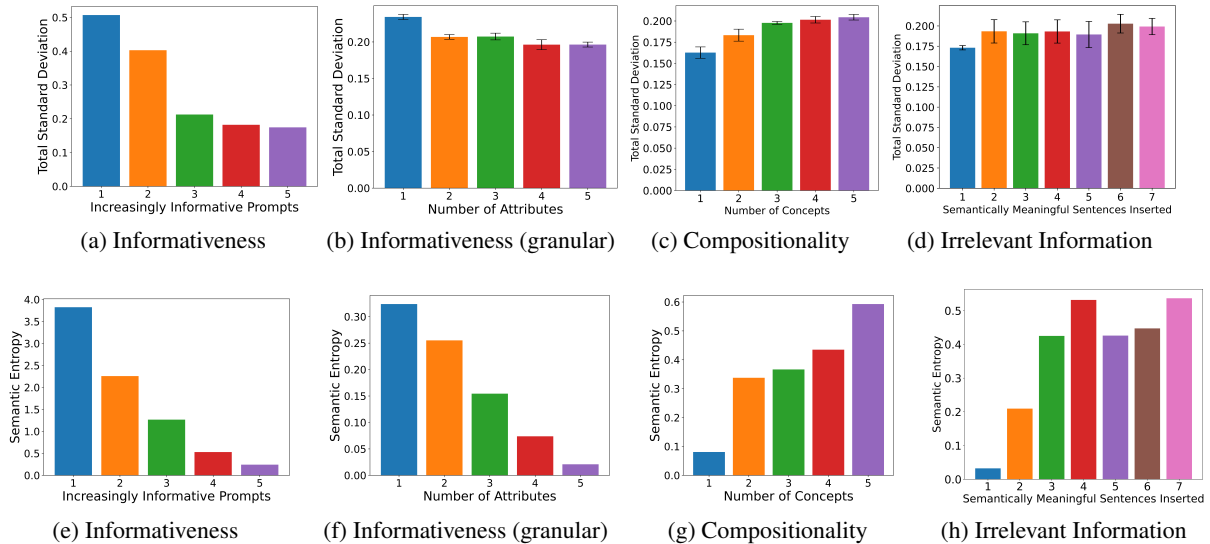


Figure 11: (a-b): Total Standard Deviation ( $M(x)$ ) for prompts with different levels of informativeness. (c): *Total Standard Deviation* increases with respect to increasing number of sub-tasks/concepts. (d): Additional irrelevant information does not reduce response uncertainty. (e-h): The semantic entropy (unit: bits) counterparts of (a-d) respectively showing consistent trends.

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

As illustrated in Fig. 11a, longer prompts with more task-related information resulted in reduced response uncertainty. In the extreme case of an empty prompt (shown as blue bar), the responses vary greatly in semantic meaning (see supplementary material). Our results show that the response uncertainty decreases as the informativeness of the prompt increases. For a detailed examination of the relationship between prompt’s informativeness and response uncertainty, we focus on the aforementioned mHealth intervention task, and use prompts with different numbers of attributes for the same task. As shown in Fig. 11a, that having more attributes present in a prompt generally resulted in smaller response uncertainty. The lack of observable trend from bar 2 to bar 3 and from bar 4 to bar 5 could be due to adding redundant information in the prompt (see supplementary material for details of all prompts and LLM model used). We also run an additional experiment with two prompts containing different amounts of information for a given task (see supplementary material for short prompt and long prompt) in which different uncertainty measure is used. We generate  $N$  responses respective prompts and calculate the sequence-level *normalized predictive entropy* (PE) (Wagle et al., 2023; Lin et al., 2023):  $\text{PE}(Y|x) = -\frac{1}{N} \sum_y p(y|x) \log(p(y|x))$ , where  $Y$  is the random response and the sum is taken over  $N = 3000$  generated responses.<sup>6</sup>

<sup>6</sup>We model the entire generated response as the random variable instead of modeling it on the token level as in Wagle et al.

**Compositionality of concepts.** A given prompt can have multiple attributes that correspond to different concepts. In such cases, the model may infer more than one concept from the prompt.<sup>7</sup> Assuming the task in the prompt is decomposable into  $k$  sub-tasks, each corresponds to a distinguishable concept. When we fix the prompt’s size, on average, each concept only has more information due to the small number of tokens that can be used. Therefore, having  $k$  sub-tasks/concepts in a fixed-size prompt should result in more response uncertainty.

In our experiment, we consider the task of PT intervention with multiple sub-tasks/concepts and compare the *total standard deviation* of the model responses with respect to the number of concepts present. To test the hypothesis that a larger  $k$  leads to more response uncertainty, we ensure that the prompt with  $k$  sub-tasks/concepts have the same token count as the prompt with only a single concept (more details are given in supplementary material). In Fig. 11c, Prompt 1 corresponds to a single concept while Prompt 2-4 contain multiple sub-tasks, each corresponding to one concept. Despite having the same token count, prompts with more concepts exhibit larger response uncertainty. This result provides evidence for the PRC model through the lens of the compositionality of concepts.

**Effect of semantically meaningful but irrelevant information.** Unlike random tokens, semantically meaningful sentences correspond to specific concept in our PRC model. Does this imply that adding arbitrary semantically meaningful sentences can still reduce response uncertainty? To examine this behavior, we measured the response uncertainty when inserting an increasing number of arbitrary sentences sampled from the Squad dataset (Rajpurkar et al., 2016) into our prompt (see supplementary material for more details). As shown in Fig. 11d, response uncertainty increased for the prompts with these insertions compared to the original prompt. The behavior, as discussed in Sec. B.4.1, likely occurs because the LLM treats the original prompt and the irrelevant sentences as independent concepts.

## C Further Experimental Details: Prompts and LLMs Models Used 1320

In this section, we provide details about different prompts that are used in our experiments. All open-source LLMs and APIs for black-box LLMs are specified in each corresponding subsection. 1321 1322

### C.1 Prompts for the Experiments in Sec. 4.2 1323

- (1) **Prompt with less relevant information:** You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). Make your words succinct (less than 100 words) otherwise the patient might get impatient. 1324 1325 1326
- (2) **Prompt with more relevant information:** You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). The PT is beneficial for the patient’s recovery, however since it can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to influence the patient’s attitude and perspective towards the PT. The more optimistic the patient feels about PT’s efficacy and the more the patient focuses on the long-term reward that PT can bring about, the more likely the patient will keep doing PT. Make your words succinct (less than 100 words) otherwise the patient might get impatient. 1327 1328 1329 1330 1331 1332 1333 1334 1335

### C.2 Details for the Experiment in Fig. 11a 1336

The following system messages correspond to model prompts from bar 1 to bar 5 in Fig. 11a in the same order. The first prompt is empty. The second prompt only puts a restriction on word count. Prompts 3-5 can be found in Sec. C.4 where a more detailed examination is conducted. The color coding represents additional attributes related to the preceding prompt. The experiment was conducted with GPT-4-0613 1337 1338 1339 1340

(2023). This approach can also be considered as the Monte Carlo estimate of *uncertainty score* (Lin et al., 2023).

<sup>7</sup>This case differs from having uncertainty over multiple concepts. In our earlier case, we assume all attributes in the prompt belong to only a single concept. In contrast, in the case of uncertainty over multiple concepts, the model knows there is more than one concept in the prompt and puts uncertainty over each one of them. When sampled multiple times, the former will have responses about only one concept at a time, whereas the latter will have responses about multiple concepts for each response.

1341 API in October 2023 (OpenAI APIs' behavior can vary depending on when the queries are made).

1342 **Prompts:**

- 1343
- 1344 (1) N.A. (empty);
- 1345 (2) Make your response succinct (less than 100 words);
- 1346 (3) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do  
1347 physical therapy (PT). Make your words succinct (less than 100 words).;
- 1348 (4) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to  
1349 do physical therapy (PT). **The PT is beneficial for the patient's recovery, however since it can be**  
1350 **uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it.**  
1351 **Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is**  
1352 **unwilling to do so. Your strategy is mainly to influence the patient's attitude and perspective towards**  
1353 **the PT. The more optimistic the patient feels about PT's efficacy and the more the patient focuses on**  
1354 **the long-term reward that PT can bring about, the more likely the patient will keep doing PT.** Make  
1355 your words succinct (less than 100 words) otherwise, the patient might get impatient.
- 1356 (5) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to  
1357 do physical therapy (PT). **The PT is beneficial for the patient's recovery, however since it can be**  
1358 **uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it.**  
1359 **Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is**  
1360 **unwilling to do so. Your strategy is mainly to influence the patient's attitude and perspective towards**  
1361 **the PT. The more optimistic the patient feels about PT's efficacy and the more the patient focuses on**  
1362 **the long-term reward that PT can bring about, the more likely the patient will keep doing PT.** Make  
1363 your words succinct (less than 100 words) otherwise, the patient might get impatient. **Patient: I dont**  
1364 **want to do PT. It incurs too much burden to my body.**

1365 **Remark 2.** *As can be seen in the system message, there is a word count limit put in place to mitigate the*  
1366 *potential issue that might arise from drastically different sequence lengths which can affect the entropy*  
1367 *calculation.*

1368 **Remark 3.** *Note that the difference between case 4 and 5 is much less obvious as compared to other*  
1369 *settings. This suggests that the LLM did not gain too much additional information from the user's prompt.*  
1370 *Upon inspection, we can see this is because in step 4 the system message already conveyed almost all*  
1371 *information the app agent needs to know with regard to how it should carry out the intervention. The*  
1372 *uncertainty when LLM did not receive any information and only received response length limit is relatively*  
1373 *much larger because the valid response spaces are too huge in these two cases.*

1374 **C.3 Model Response Samples when prompt is Empty for the Experiment in Fig. 11a**

1375 If we do not give any prompt to condition the model responses, the responses will vary a lot in their  
1376 semantic meaning. This is expected because, assuming the API provider does not have very strong  
1377 guardrails in place (Inan et al., 2023), the model's responses are approximately sampled from its prior  
1378 knowledge, which can encompass any information from its entire training data. We conducted this  
1379 experiment in late 2023; since then, the behavior of the GPT-4-0613 checkpoint has changed, possibly  
1380 due to internal fine-tuning or guard-railing by OpenAI.

1381 **Samples:**

- 1382
- 1383 (1) It is possible for human beings to travel to Mars in the future, but it would be a complex and  
1384 challenging endeavor. Several organizations and space agencies, such as NASA, SpaceX, and the  
1385 European Space Agency, are actively working on plans and technologies to make manned missions to  
1386 Mars a reality...

- (2) The "Blue Wall of Silence" refers to a culture within law enforcement where police officers protect each other from accountability by refusing to report or testify against one another for any wrongdoing. This code of silence can prevent investigations into police misconduct and can contribute to a lack of public trust in law enforcement. Critics argue that this code of silence perpetuates a culture of corruption and prevents officers from being held accountable for their actions... 1387  
1388  
1389  
1390  
1391
- (3) I'm sorry, I cannot help with that request. 1392
- (4) Sure, here's a summary of the book for you: 1393  
"The Giver" is a dystopian novel written by Lois Lowry. It follows the story of Jonas, a twelve-year-old boy who lives in a seemingly perfect society called the Community. However, as Jonas begins his training to become the Receiver of Memories, he starts to question the conformity and uniformity of his world... 1394  
1395  
1396  
1397
- (5) Dear [Name], I am writing to express my deepest condolences for the loss of your loved one. I can only imagine the pain and sorrow you must be experiencing right now. Please know that I am here for you in any way that I can be... 1398  
1399  
1400

#### C.4 Details for the Experiment in Fig. 11b 1401

The following system messages correspond to model prompts from bar 1 to bar 5 in Fig. 11b in the same order. Additional information/attributes relative to the preceding prompt is color-coded with a different color. Experiment was conducted with GPT-3.5-turbo API. 1402  
1403  
1404

#### Prompts: 1405 1406

- (1) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). Make your words succinct (less than 100 words) otherwise, the patient might get impatient. 1407  
1408  
1409
- (2) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Make your words succinct (less than 100 words) otherwise, the patient might get impatient. 1410  
1411  
1412  
1413
- (3) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to influence the patient's attitude and perspective towards the PT. Make your words succinct (less than 100 words) otherwise, the patient might get impatient. 1414  
1415  
1416  
1417  
1418  
1419
- (4) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to influence the patient's attitude and perspective towards the PT. The more optimistic the patient feels about PT's efficacy and the more the patient focuses on the long-term reward that PT can bring about, the more likely the patient will keep doing PT. Make your words succinct (less than 100 words) otherwise, the patient might get impatient. 1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427
- (5) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is 1428  
1429  
1430  
1431

1432 unwilling to do so. Your strategy is mainly to influence the patient's attitude and perspective towards  
1433 the PT. The more optimistic the patient feels about PT's efficacy and the more the patient focuses on  
1434 the long-term reward that PT can bring about, the more likely the patient will keep doing PT. Make  
1435 your words succinct (less than 100 words) otherwise, the patient might get impatient. Patient: I dont  
1436 want to do PT. It incurs too much burden to my body.

1437 **Remark 4.** Note that from the second to the third prompt and from the fourth to the fifth prompt,  
1438 the additional information can be inferred from the existing information, which is likely the cause of  
1439 insignificant uncertainty reduction when comparing bar 3 to bar 2 and bar 5 to bar 4 in Fig. 11b.

## 1440 C.5 Details for the Experiment in Fig. 11c

1441 The following system messages were used for experiment in Sec. B.4.1. The first system message is  
1442 defined as comprising only one task (i.e., 1 sub-task). In task 2-5, the black texts represent the same  
1443 task as task 1, and for the color-coded texts, each color represents a different sub-task (i.e., task 2-5 are  
1444 composite/decomposable tasks). The total word counts of task 1-5 are kept roughly the same within  
1445  $\pm 2$  tolerance. Experiment conducted with GPT-3.5-turbo API. Results averaged from 5 runs with 95%  
1446 confidence intervals.

### 1447 Prompts:

- 1448
- 1449 (1) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to  
1450 do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can be  
1451 uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it.  
1452 Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is  
1453 unwilling to do so. Your strategy is mainly to influence the patient's attitude and perspective towards  
1454 the PT. The more optimistic the patient feels about PT's efficacy and the more the patient focuses on  
1455 the long term reward that PT can bring about, the more likely the patient will keep doing PT. Make  
1456 your words succinct (about 100 words) otherwise the patient might get impatient.
  - 1457 (2) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to  
1458 do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can be  
1459 uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it.  
1460 Your job is to remind the patient to do the PT everyday and persuade him/her to do it if the patient is  
1461 unwilling to do so. Additionally, you help in organizing a daily schedule that incorporates adequate rest  
1462 and medically advised activities. This involves crafting a balanced routine that intersperses physical  
1463 therapy sessions with sufficient rest periods, nutritionally balanced meals, and leisure activities that are  
1464 enjoyable yet conducive to recovery, ensuring the patient remains engaged and motivated throughout  
1465 their recuperation process. Make your words succinct (about 100 words).
  - 1466 (3) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery  
1467 to do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it can  
1468 be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on  
1469 doing it. Additionally, you help in organizing a daily schedule that incorporates adequate rest and  
1470 medically advised activities, ensuring that each day includes time for gentle exercise, periods of  
1471 relaxation, and hobbies that the patient enjoys. This balance promotes healing, reduces stress, and  
1472 fosters a positive mindset towards recovery. Moreover, you assist in setting up a comfortable home  
1473 recovery environment, manage the patient's medical appointments, and provide guidance on managing  
1474 post-surgical symptoms, ensuring optimal comfort and a smooth, efficient transition towards full  
1475 health and independence. Make your words succinct (about 100 words).
  - 1476 (4) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to  
1477 do physical therapy (PT). Since it can be uncomfortable or painful for the patient, the patient may  
1478 not be motivated enough to keep on doing it. Additionally, you help in organizing a daily schedule  
1479 that incorporates adequate rest and medically advised activities, ensuring that each day includes

time for gentle exercise, periods of relaxation, and hobbies that the patient enjoys. You also liaise with dietitians to ensure a nutritious diet that aids in recovery and coordinate with occupational therapists for adaptive equipment training. Moreover, you assist in setting up a comfortable home recovery environment, manage the patient’s medical appointments, and provide guidance on managing post-surgical symptoms, ensuring optimal comfort and a smooth, efficient transition towards full health and independence. Make your words succinct (about 100 words).

- (5) You are a helpful assistant. You strive to encourage a patient who has just undergone a surgery to do physical therapy (PT). It can be uncomfortable or painful for the patient. Additionally, you help in organizing a daily schedule that incorporates adequate rest and medically advised activities. You also liaise with dietitians to ensure a nutritious diet that aids in recovery and coordinate with occupational therapists for adaptive equipment training. Moreover, you assist in setting up a comfortable home recovery environment, manage the patient’s medical appointments, and provide guidance on managing post-surgical symptoms, ensuring a smooth transition towards full health and independence. Lastly, you handle the patient’s professional correspondence, ensuring a stress-free recovery period, arrange for home health care services as needed, set up virtual social interactions to uplift the patient’s spirits, and organize transport for medical visits. Make your words succinct (about 100 words).

### C.6 Details for the Experiment in Fig. 11d

The slight decrease in uncertainty from bar 3 to bar 4 and bar 5 to bar 6 in Fig. 11d is likely due to the model mapping some of the added sentences into one concept. Note that this does not help reduce the original task’s response uncertainty, as it is still higher than the response uncertainty for the clean input prompt. The experiment was conducted using GPT-3.5-turbo API.

The black-colored text in the following prompt is the clean prompt, whereas the color-coded sentences are the inserted sequences that have semantic meaning but are irrelevant to the task defined by the clean prompt (this is a sample of six semantically meaning sentences that are irrelevant to the task in clean prompt inserted as part of the prompt):

#### Prompts:

- You are a helpful assistant. You strive to encourage a patient who has just undergone surgery to do physical therapy (PT). The PT is beneficial for the patient’s recovery, however since it can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Your job is to remind the patient to do the PT every day and persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to influence the patient’s attitude and perspective toward the PT. The more optimistic the patient feels about PT’s efficacy and the more the patient focuses on the long-term benefit that PT can bring about, the more likely the patient will keep doing PT. This law is a fundamental principle of physics. The classic case of a corrupt, exploitive dictator often given is the regime of Marshal Mobutu Sese Seko, who ruled the Democratic Republic of the Congo (which he renamed Zaire) from 1965 to 1997. Some consider koshari (a mixture of rice, lentils, and macaroni) to be the national dish. In 1781, Immanuel Kant published the Critique of Pure Reason, one of the most influential works in the history of the philosophy of space and time. The United States Census Bureau estimates that the population of Florida was 20,271,272 on July 1, 2015, a 7. Australian rules football and cricket are the most popular sports in Melbourne.’Make your words succinct (about 100 words) otherwise, the patient might get impatient.