UNDERSTANDING THE RELATIONSHIP BETWEEN PROMPTS AND RESPONSE UNCERTAINTY IN LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) are widely used in decision-making, but their reliability, especially in critical tasks like healthcare, is not well-established. Therefore, understanding how LLMs reason and make decisions is crucial for their safe deployment. This paper investigates how the uncertainty of responses generated by LLMs relates to the information provided in the input prompt. Leveraging the insight that LLMs learn to infer latent concepts during pretraining, we propose a prompt-response concept model that explains how LLMs generate responses and helps understand the relationship between prompts and response uncertainty. We show that the uncertainty decreases as the prompt's informativeness increases, similar to epistemic uncertainty. Our detailed experimental results on real-world datasets validate our proposed model.

1 INTRODUCTION

025 026

000

001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

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 various decision-making applications, such as healthcare (Karabacak & Margetis, 2023; Sallam, 2023; Yang et al., 2023), education (Xiao et al., 2023), finance (Wu et al., 2023b), and law (Zhang et al., 2023a). However, despite their rapid adoption, the reliability of LLMs in handling high-stakes tasks has yet to be demonstrated (Arkoudas, 2023; Huang et al., 2023a). The reliability is particularly critical in domains such as healthcare, where model responses can have immediate and significant impacts on human behavior and hence their well-being (Ji et al., 2023). Thus, understanding LLMs' reasoning and decision-making processes and how they influence response uncertainty is critical for their safe and reliable deployment.

To understand this importance, consider the mobile health (mHealth) application in which machine 037 learning algorithms are integrated to monitor users' health conditions and provide advice on daily activities (Boursalie et al., 2018; Trella et al., 2022; 2023). Providing suggestions that can influence users' health is a form of intervention in human decision-making. For LLMs to be suitable for such 040 use cases, they should be accurate and provide consistent intervention strategies, e.g., consider an 041 LLM-powered mHealth app that suggests physical therapy (PT) routines to a patient recovering 042 from surgery. mHealth ensures the patient adheres to their PT regimen during rehabilitation despite 043 the discomfort it may cause. mHealth must provide not only good but consistent suggestions to 044 encourage PT adherence. Any inconsistent behaviors from the app could undermine any progress 045 made. Conversely, providing accurate and consistent responses helps make the system more reliable and trustworthy (Shin et al., 2022). 046

Traditionally, when model architecture is fixed, model improvement relies on better hyperparameters (Bischl et al., 2023), such as using a more suitable optimizer (Hassan et al., 2023), or train it with more/better data (Simon et al., 2023). Due to the emergent capability of LLMs (Dong et al., 2022), such as in-context learning (ICL) (Kojima et al., 2022) and chain-of-thought (CoT) prompting (Wei et al., 2022), it is possible to improve the model responses by guiding it with informative prompts having relevant instructions and exemplars. Through these techniques, the LLMs can more effectively use the relevant information acquired from the training data to generate better responses, even if the prompt itself does not directly reveal the ground truth (Liu et al., 2023; Sahoo et al., 2024).

054 The response generated by LLMs is a series of tokens sampled from probability vectors of tokens 055 using various heuristics (Radford et al., 2018; 2019; Brown et al., 2020), such as beam search, nucleus 056 sampling, and greedy decoding. Typically, tokens with higher probabilities are chosen sequentially 057 to produce the final response. The response variations are controlled by LLM hyperparameters 058 such as temperature (T), top-k, or top-p. While response variations benefit creative tasks like poem and essay writing, they can be detrimental for tasks requiring high reproducibility and consistency (Ganguli et al., 2022; Huang et al., 2023c). However, making LLMs generate deterministic responses 060 is not ideal, as users preferences may vary (Wu et al., 2023a). Hence, a better approach is needed to 061 understand the sources of response uncertainty and develop methods to reduce it naturally rather than 062 masking it by adjusting LLM hyperparameters. 063

For a fixed or black-box LLM, response uncertainty can be controlled by two ways: adjusting LLM's hyperparameters (e.g., temperature) to control the randomness in generated response and providing more task-relevant information in input prompt (prompt informativeness). This paper focuses on the response uncertainty due to the input prompt while keeping the LLM hyperparameters fixed. Hence, we address the following question: *How is the amount of relevant information about a task in the input prompt related to the uncertainty in the response generated by an LLM?* We use the insight that LLMs implicitly learn to infer latent concepts during pretraining (Xie et al., 2022; Hahn & Goyal, 2023; Zhang et al., 2023b) and propose a novel *prompt-response concept (PRC) model* in Section 2.

Our PRC model conceptualizes how an LLM generates responses based on given prompts and 072 helps understand the relationship between prompts and response uncertainty by measuring response 073 uncertainty for prompts with varying information about the task. We provide theoretical results that 074 show the uncertainty of responses generated by an LLM decreases as the prompt informativeness 075 and model quality increase. We draw connection between the reducible response uncertainty and 076 epistemic uncertainty, and using our PRC model, we explain why adding more relevant information 077 to the prompt to a better trained LLM is a principled and effective method to reduce the response uncertainty. Finally, we corroborate the validity of the PRC model via experiments and provide a 079 simulation for a healthcare use case to demonstrate its efficacy. Specifically, our contributions can be 080 summarized as follows:

- 081
- 082

084

085

087

090

091 092

- In Section 2, we propose a prompt-response concept model to quantify the relationship between prompt informativeness and response uncertainty in LLMs.
- In Section 3, we prove that response uncertainty decreases with higher prompt informativeness and model quality. Using our PRC model, we relate the reducible uncertainty to epistemic uncertainty and explain how adding relevant information in prompt to a welltrained LLM reduces response uncertainty.
- Finally, we validate the PRC model through experiments and demonstrate its theoretical efficacy using different tasks derived from real-world datasets and a healthcare use case simulation in Section 4.
- 2 PROMPT-RESPONSE CONCEPT MODEL FOR LLM

In this section, we first define what we mean by *concept*. We then use the notion of concept to explain our proposed prompt-response concept model of LLM. Finally, we provide theoretical results that explain the relationship between prompt informativeness and the uncertainty of responses.

096 **Concept.** The definition of a concept varies across fields, e.g., in philosophy, a concept represents 097 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 098 this paper, we define the *concept* as an abstraction derived from specific instances or occurrences that share common characteristics (Fodor, 1998; Laurence & Margolis, 1999; Weiskopf, 2009; Wilmont 100 et al., 2013). To understand the notion of concept, consider this example of the concept: Species, 101 which includes a group of organisms that share common biological traits. Furthermore, a concept can 102 be expressed as a sequence with semantic meaning (e.g., using natural language). Here, we use the 103 term 'semantic meaning' (or semantically meaningful) to refer to something that conveys information 104 that is understandable and extractable by humans (Hurford et al., 2007). Consider another example of 105 a concept: the *personal bio*,¹ which consists of sentences giving information about a person's name, 106 occupations, contributions, and other personal details.

¹This example is adapted from *wiki bio* concept example given in Xie et al. (2022).

109

110 111

112

Concept: Personal bio of Alan Turing

Alan Turing was an English computer scientist, mathematician, and cryptanalyst. He introduced the Turing machine, which formalized the concepts of algorithms and computation, serving as a foundational model for general-purpose computers. Turing is regarded as the father of theoretical computer science. ...

Concept attributes. At first glance, the exact set of attributes defining a concept is not immediately 113 clear. For instance, whether a attribute like 'a certain hobby of Alan Turing' in the above example, 114 is part of the concept is subjective. In practice, the attribute set of a concept depends on the data 115 the LLM is trained on, e.g., if the hobby of Alan Turing appears in the training data, the LLM has 116 access to this knowledge, incorporating it into the concept's attribute set. Additionally, the notions 117 of 'concept' and 'attribute' are relative in nature. An attribute of sufficient complexity could itself 118 be considered a concept, encompassing its own set of attributes. For illustration purposes, we use 119 paragraph-level concepts and sentence-level attributes, instead of word- or token-level patterns, as 120 these higher-level abstractions allow a better understanding of the relationships between different 121 sentences within the text (Bates, 1995; Bogatyrev & Samodurov, 2016; Wang et al., 2024).

122 In the above example of *personal bio* concept, explaining a concept often involves multiple sentences, 123 each contributing to a specific and meaningful facet of information about the concept (Piccinini & 124 Scott, 2006). We refer to the information of each facet as a concept attribute, e.g., the sentence, "Alan 125 Turing was an English computer scientist, mathematician, and cryptanalyst" provides information 126 about the name, nationality, and occupation of Alan Turing. Thus, a concept can be fully characterized 127 by all its attributes. Recent works by Gao et al. (2024), Lieberum et al. (2024), and Templeton (2024) 128 provide mechanistic interpretability evidence suggesting that LLMs can learn concept-like features. In this work, we formalize the *definition of concept* and then use it to provide theoretical insights into 129 the relationship between prompt informativeness and LLM response uncertainty. 130

131 132

149

150 151

152

153

154

156

157 158

159

2.1 PROMPT-RESPONSE CONCEPT MODEL

We aim to understand how the input prompt informativeness is related to the uncertainty in the 133 responses generated by an LLM. To do so, we first introduce notations representing different variables 134 used in this section. Let \mathcal{X} denote the set of all prompts and \mathcal{Y} denote the set of all responses 135 generated by an LLM f, where $f: \mathcal{X} \to \mathcal{Y}$. Let \mathcal{V} be the vocabulary containing all unique tokens. 136 For any prompt x and response y, we have $x \in \mathcal{V}^{|x|}$ and $y \in \mathcal{V}^{|y|}$, where $|\cdot|$ returns the number of 137 tokens in prompt/response. For a given prompt $x \in \mathcal{X}$, the LLM f generates a response $y \in \mathcal{Y}$ such 138 that y = f(x). Since the response y can vary each time the LLM generates it, we can control these 139 response variations using two main factors: prompt informativeness and LLM's hyperparameters, 140 such as temperature T, top-k, or top-p, which control the randomness in the generated tokens. 141

This paper focuses solely on the latter aspect while keeping the LLM's hyperparameters fixed. Building on the earlier works' interpretation that LLMs implicitly learn to infer latent concepts during pretraining (Xie et al., 2022; Hahn & Goyal, 2023; Zhang et al., 2023b), we propose the prompt-response concept (PRC) model of LLM. This model conceptualizes how an LLM generates a response for a given prompt, which will be used to understand the relationship between prompts and the response uncertainty by measuring response uncertainty for prompts with varying information. The PRC model has three main components (as shown in Fig. 1): prompt concept, response concept, and mapping functions, whose details are given as follows.



Figure 1: Prompt-Response concept model of LLM.

Prompt concept. Let Θ_x be the set of all concepts corresponding to prompts in set \mathcal{X} . We assume that each input prompt $x \in \mathcal{X}$ corresponds to a concept. We refer to this concept as the *prompt concept* $\theta_x \in \Theta_x$. Intuitively, an LLM recognizes input tokens as semantically meaningful units that coherently describe an attribute of some latent prompt concept. The concept's attributes are expressed
 through multiple semantically meaningful sentences. If multiple sentences in the prompt cannot be
 combined to describe a single concept, the LLM treats them as representing different concepts. Our
 experiments in Fig. 4c show that *adding semantically meaningful information from different concepts can increase response uncertainty*.

Response concept. Let Θ_y be the set of all concepts corresponding to the possible responses in the set \mathcal{Y} . We refer to these concepts as the *response concept*. The PRC model assumes that each response concept $\theta_y \in \Theta_y$ is associated with a specific response $y \in \mathcal{Y}$.

Mapping functions. To understand the relationship between input prompt, intermediate concepts 171 (i.e., prompt and response concept), and response, we assume the LLM f is a composition of three 172 mapping functions: prompt-concept mapping function (q_x) , concept-concept mapping function (q_c) , 173 and concept-response mapping function (g_y) . Hence, we can represent the response as y = f(x) =174 $g_u(g_c(g_x(x)))$, where the function g_x maps the input prompt to a prompt concept, the function g_c 175 maps the prompt concept to a response concept, and finally, the function g_y maps the response 176 concept to a response. For a well-pretrained LLM, its g_x can accurately find all the attributes in the 177 given input prompt; for a well-aligned LLM, its g_c can accurately map from some θ_x to the most 178 relevant θ_{u}^{2} . Therefore, we expect a better LLM to generate responses with lower uncertainty and higher qualities, as demonstrated by our experiments in Fig. 2a and Fig. 2c. 179

To generate a response, an LLM maps the input prompt to a prompt concept, which is then mapped to a corresponding response concept. Finally, the LLM uses the response concept to generate the final response. When a prompt lacks sufficient task-related information (i.e., it is less informative) due to attributes being under-specified, we can expect higher variability in the responses generated by the LLM because there potentially exists multiple concepts that satisfy the attributes found in the input prompt, as corroborated by our experimental results in Fig. 5a. To further understand this relationship, we will next formalize how the informativeness of prompts is related to response uncertainty.

187 188

2.2 RELATIONSHIP BETWEEN PROMPTS AND RESPONSE UNCERTAINTY IN LLMS

189 Let $\mathcal{X}_{\theta_n} \subset \mathcal{X}$ be the set of prompts with the same semantic meaning³ and contain all information 190 of the prompt concept θ_x . Let $\mathcal{X}_s \subset \mathcal{X}_{\theta_x}$ be the set of prompts with the same semantic meaning 191 s and only contain partial information about the prompt concept θ_x . Let $\mathcal{A}_{\theta} = \{a_{\theta,1}, \ldots, a_{\theta,m}\}$ 192 represent the set of all the attributes of a concept θ and each attribute can be perfectly expressed by 193 some semantically meaningful sequence of tokens. We use the notation $x_1 \prec_{\theta_x} x_2$ to indicate that 194 prompt x_1 contains less information about prompt concept θ_x than the prompt x_2 (or the prompt x_2 195 is more informative than the prompt x_1). Since $\mathcal{X}_s \subset \mathcal{X}_{\theta_x}$, any prompt from the set \mathcal{X}_s contains less 196 information about prompt concept θ_x than any prompt from the set \mathcal{X}_{θ_x} . Let Z_c be a random variable representing a concept (where c = x for prompt concept and c = y for response concept) and X_s be 197 a random variable representing a prompt with semantic meaning s. Here, the randomness in Z_c from 198 a less informative prompt, which allows more variation in the possible concepts that LLM can map 199 to. In contrast, the randomness in X_s is due to the ability of different prompts to represent the same 200 underlying semantic meaning. 201

202 We use entropy as a measure to quantify the uncertainty in responses generated by an LLM for a given input prompt. Entropy captures the randomness of the responses and helps in understanding 203 how the informativeness of an input prompt affects response uncertainty. Let Y be a random 204 variable representing the response. The randomness in Y can be due to less informative prompts 205 and the ability of different responses to correspond to the same response concept (i.e., have the 206 same semantic meaning). For a given prompt x, we define entropy of Y as follows: H(Y|x) =207 $-\sum_{y\in\mathcal{Y}} p(y|x) \log_2 p(y|x)$, where p(y|x) is the conditional distribution of the responses generated 208 for a prompt. Intuitively, a highly informative prompt corresponds to specific intermediate concepts, 209 which leads to the generation of responses with less variability and, hence, smaller entropy of Y. The 210 conditional distribution p(y|x) represents the posterior predictive distribution, which marginalizes all 211

212

²¹³ ²For pretrained-only LLMs in the ICL setting, the task-dependent g_c is inferred on the fly from the exemplars 214 in the input prompt, as shown in Xie et al. (2022).

³Multiple prompts can be generated from a single prompt by paraphrasing it while preserving the original semantic meaning associated with the prompt (Kuhn et al., 2023).

intermediate concepts (prompt and response) and is given as

218 219

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

220 The first equality follows from conditioning the response on the response concept, and the second 221 equality follows from $p(\theta_y|x) = \int_{\theta_x} p(\theta_y|\theta_x, x) p(\theta_x|x) d\theta_x$. If $p(\theta_c|x)$ (where $c = \{x, y\}$) con-222 centrates on a specific concept with a more informative prompt, the LLM learns effectively via 223 marginalization. More concretely, our PRC model assumes the LLM achieves this by extracting 224 the attributes from the input prompt x and matching it to the right prompt and response concepts. Furthermore, when the prompt has the information about all attributes (i.e., perfect prompt), it is 225 enough to completely characterize the concept (Proposition 1). If an LLM extracts all information 226 about the desired concept from a perfect prompt, the remaining uncertainty in responses is due to the 227 representation of the response concept via different responses (i.e., semantically equivalent), which 228 is the irreducible uncertainty. This behavior implies that the LLM implicitly performs Bayesian 229 inference, which is also observed in ICL (Xie et al., 2022). 230

231 232

233

3 THEORETICAL RESULTS

In this section, we first introduce the assumptions under which our theoretical results hold.

Assumption 1. We assume a well-trained LLM knows the mapping functions used in PRC model, i.e., g_x , g_c , and g_y .

This assumption states that LLM has perfectly learned the mapping functions used in our proposed
PRC model. While this assumption may not hold in practice, we can expect that a better LLM has
good estimates of these mapping functions, which lead to lower uncertainty and better qualities of the
model responses, as corroborated by our experimental results shown in Fig. 2a and Fig. 2c. Next, we
present our first result, which shows the relationship between concept uncertainty and informativeness
of prompts.

Lemma 1. Let Assumption 1 hold. For any 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 Appendix A. The first part of this result implies that prompts fully describing two different concepts can not 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 can not have the same semantic meaning. The second part implies that there is no randomness in the prompt concept if all the information needed to respond to the task is contained in the prompt. Our next result shows the relationship between prompt informativeness and concept uncertainty.

Proposition 1. Let Assumption 1 hold. Then, $H(Z_x|X_s)$ strictly decreases as X_s represents more informative prompts, i.e., as more relevant information about the concept is included in the prompt.

²⁵⁵ We now state our main result that links response uncertainty to the informativeness of a prompt.

Theorem 1. Let Assumption 1 hold. Then, $H(Z_y|X_s)$ strictly decreases as X_s represents more informative prompts. Furthermore, $H(Y|X_s)$ converges to $H(Y|Z_y) + \mathcal{E}$, where $\mathcal{E} \leq H(g_c(Z_x)|Z_x)$.

259 Above two results suggest that as prompt informativeness increases, response uncertainty due to 260 the uncertainty in the response concept decreases. Further, when sufficient information is provided 261 in a prompt, no uncertainty remains due to the the prompt concept. The remaining randomness in responses can be decomposed into two terms: $H(Y|Z_y)$, which represents semantic redundancy, 262 is due to the ability of different responses to convey the same semantic meaning, making them 263 semantically equivalent (Kuhn et al., 2023). The term \mathcal{E} is due to the imperfection of LLMs such that 264 it has not learn a perfect g_c . For example, if g_c is stochastic, we observe multiple realizations of z_u 265 for the same z_x in different iterations due to the randomness of g_c , leading to variation in the model 266 responses. However, as shown in Fig. 4b, as the LLM quality improves, \mathcal{E} gets smaller, resulting in 267 lower overall response uncertainty. 268

269 Concept Uncertainty, Semantic Redundancy, and Model Imperfection as Epistemic Uncertainty. In machine learning literature, epistemic uncertainty is typically reduced by incorporating 270 additional information, such as using a better model and additional training data (Hüllermeier & 271 Waegeman, 2021; Lahlou et al., 2021; Senge et al., 2014; Shaker & Hüllermeier, 2020; Valdenegro-272 Toro & Mori, 2022; Der Kiureghian & Ditlevsen, 2009). In Proposition 1, H ($Z_c|X_s$) represents the 273 epistemic uncertainty in latent concepts.⁴ We have demonstrated that H $(Z_c|X_s)$ is strictly reduced 274 with a prompt that contains more attributes of the relevant concept(s). Therefore, increasing the information about the concept in a prompt can lead to more reliable and consistent responses by 275 reducing the epistemic uncertainty in the latent concept (Hüllermeier & Waegeman, 2021). When the 276 prompt perfectly captures the desired concept, the posterior distribution of the concept given prompt 277 converges to the desired concept. 278

279 Due to the model's inability to learn a perfect g_c during training, the mapping from prompt concept 280 to the response concept can not be perfect. In Theorem 1, \mathcal{E} captures this source of uncertainty. In scenarios where model parameters are allowed to be modified, this uncertainty is epistemic and can 281 be reduced as the model quality improves (Fig. 4b). The remaining uncertainty is due to *semantic* 282 redundancy. It can be further reduced in two ways: use fine-tuning or prompting to instruct the model 283 to reply with certain fixed style.⁵ Due to semantic equivalence, *semantic redundancy* is generally not 284 detrimental to the desired information. However, if the prompt contains sentences that are irrelevant 285 to the task, the response uncertainty can increase, as demonstrated in Fig. 4d. It is possible that a g_c can result in low response uncertainty but poor response quality due to the wrong mapping from Z_x 287 to Z_y (Singh et al., 2023; Li et al., 2024; Fu et al., 2025). 288

200 289 290

291

297

4 EXPERIMENTS

To corroborate our proposed prompt-response concept model of LLM, 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 input 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.

4.1 RANDOMNESS IN PROMPTS VS. RESPONSE QUALITY

298 The low uncertainty in model responses does not necessarily indicate high response quality, as an LLM 299 can produce outputs with very low uncertainty while being blindly confident in incorrect answers. 300 This behavior is problematic and can lead to hallucinations (Huang et al., 2023b). To ascertain 301 the actual relationship between prompt, model response uncertainty, and quality, we investigated 302 the relationship between the effective token count of the input prompt and model response quality. 303 To assess if the reduction in uncertainty translates to improved output quality, we test the model's 304 output accuracy when answering the multiple-choice questions (MCQs). We selected 100 MCQs 305 from the medical meadow medga (Jin et al., 2020) and ARC (Clark et al., 2018) datasets, which 306 serve as domain-specific (healthcare) and general reasoning tests, respectively. We iteratively select 307 an increasing fraction of randomly selected tokens from the context of the questions, respectively, replacing them with space tokens (i.e., token corruption). For each question, we set the temperature 308 to 1 and sampled 100 responses from the model. We used 5 different random seeds to choose which 309 tokens to corrupt, replacing them with space tokens. As the fraction of corruption increased, we 310 added new randomly selected tokens in the previously replaced tokens to ensure that randomness 311 from existing corrupted tokens did not contribute to changes in accuracy. This method allowed us to 312 observe the effect of token corruption on the model response quality and accuracy. 313

In Figs. 2a and 2b, we plot the accuracy for the responses of three open-source and three blackbox 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 questions⁶ (Fig. 2c, Fig. 2d) as an indicator of response uncertainty. As corruption becomes more severe, we observe that the

319

⁴It is called the *semantic entropy* in Kuhn et al. (2023). In this paper, we study it through the lens of uncertainty reduction.

⁵The response style can be viewed as an implicit concept, so *semantic redundancy* can be reduced by providing relevant style information in the input prompt to guide the model response.

⁶The 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



Figure 2: Accuracy of MCQs (a-b) and Empirical entropy (uncertainty) over MCQs (c-d). There is a clear and strong negative correlation between accuracy and uncertainty, with less accurate models generally showing greater uncertainty in their responses.

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. 4b, models with better accuracy tend to have lower empirical entropy. This corroborates our interpretation of \mathcal{E} in Theorem 1. The experimental results on other tasks can be found in Appendix D.1.

341 4.2 MHEALTH INTERVENTION USECASE 342 Whealth intervention of second sec

332

333

334

343

344

345

346

347

348

354

355

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 (see Fig. 11).



Figure 3: 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. More details of this experiment are given in Appendix D.2.

368 4.3 ABLATION STUDIES

We run a series of ablation studies to analyze the impact of various components, like prompt informativeness, compositionality, and irrelevant information, on response uncertainty.

Relationship between informativeness of the prompt 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. The details on the experiment set up can be found in Appendix C.1. As illustrated in Fig. 4a, longer prompts with more task-related information resulted in reduced response uncertainty. In the extreme case of an empty input prompt (shown as blue bar), the

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.

378 responses vary greatly in semantic meaning (see Appendix E.2). Our results show that the response 379 uncertainty decreases as the informativeness of the input prompt increases. For a detailed examination 380 of the relationship between input prompt's informativeness and response uncertainty, we focus on the 381 aforementioned mHealth intervention task, and use prompts with different numbers of attributes for 382 the same task. As shown in Fig. 5a, 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 input prompt (see Appendix E.3 for 384 details of all prompts and LLM model used). We also run an additional experiment with two prompts 385 containing different amounts of information for a given task (see Appendix E.4 for short prompt and 386 long prompt) in which different uncertainty measure is used. We generate N responses respective 387 prompts and calculate the sequence-level normalized predictive entropy (PE) (Wagle et al., 2023): 388 $PE(Y|x) = -\frac{1}{N} \sum_{y} p(y|x) \log(p(y|x))$, where Y is the random response and the sum is taken over 389 N = 3000 generated responses.⁷ As we observed in Fig. 5b, the responses generated with the longer 390 prompt containing more relevant information have consistently smaller PE than those from the shorter 391 prompt as the sample size grows. Additional experiment results are given in Appendix C. 392



Figure 4: (a) Total Standard Deviation (M(x)) for input prompts with different levels of informativeness. (b): 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. A clear negative correlation can be observed: the better the quality of the model, the less response uncertainty it has. (c): *Total Standard Deviation* increases with respect to increasing number of sub-tasks/concepts. (d): Additional irrelevant information does not reduce response uncertainty.

413 In our experiment, we consider the task of PT intervention with multiple sub-tasks/concepts and 414 compare the *total standard deviation* of the model responses with respect to the number of concepts 415 present. To test the hypothesis that a larger k leads to more response uncertainty, we ensure that the 416 prompt with k sub-tasks/concepts have the same token count as the prompt with only a single concept 417 (more details are given in Appendix E.6). In Fig. 4c, Prompt 1 corresponds to a single concept while 418 Prompt 2-4 contain multiple sub-tasks, each corresponding to one concept. Despite having the same 419 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. 420

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 answer this question, we measured the response uncertainty when inserting an increasing number of arbitrary

425

 ⁷We model the entire generated response as the random variable instead of modeling it on the token level as in Wagle et al. (2023). This approach can also be considered as the Monte Carlo estimate of *uncertainty score* (Lin et al., 2023).

 ⁴²⁸ ⁸This case differs from having uncertainty over multiple concepts. In our earlier case, we assume all attributes
 ⁴²⁹ in the input 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 input prompt and puts uncertainty over each one of them.
 ⁴³¹ 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.

sentences sampled from the Squad dataset (Rajpurkar et al., 2016) into our prompt (more details
in Appendix E.7). As shown in Fig. 4d, response uncertainty increased for the prompts with these
insertions compared to the original prompt. The behavior likely occurs because the LLM treats the
original input prompt and the irrelevant sentences as independent concepts.

436 437

438

5 RELATED WORK

439 **Uncertainty quantification for LLMs.** While uncertainty quantification is an extensively studied 440 topic in machine learning, there have been limited explorations for LLMs. Kadavath et al. (2022) 441 studies to what extent the LLMs can accurately conduct self-evaluation on what knowledge they 442 possess and how much calibration can help improve model response quality, where the main goal 443 of calibrating LLMs is to let the variation in the responses genuinely reflect the model's lack of 444 relevant knowledge with respect to the prompt. Xiao et al. (2022) and Wagle et al. (2023) empirically 445 investigated pre-trained language models (PLMs) and retrieval augmented language models (RALMs), respectively and found out that while both types of models tend to be overly confident in their response, 446 models with larger size are better calibrated. In contrast, RALMs exhibit worse calibrations compared 447 to their counterparts. An orthogonal work of Lin et al. (2023) devised a method using similarity as 448 determined by a Natural Language Inference (NLI) model, along with simple measures that measure 449 dispersion based on these similarities to quantify the uncertainty and the confidence of black-box 450 LLMs in the context of question-answering tasks. Kuhn et al. (2023) introduced the notion of *semantic* 451 *entropy* to more precisely quantify the uncertainty of the information content of model responses, 452 eliminating interference from the variation in semantically equivalent responses. Similar to our work, 453 Ling et al. (2024) attempt to understand and quantify LLMs' response uncertainty by decomposing 454 it into aleatoric and epistemic uncertainty, but their study is confined within the ICL setting and 455 assumes the correlation between model response accuracy and uncertainty without any examination. 456 In contrast, our framework addresses both pretrained and fine-tuned LLMs, and we investigated if lower uncertainty in model response necessarily implies higher quality. Similar to Wagle et al. (2023), 457 Lin et al. (2023), and Kuhn et al. (2023), we adopted an entropy-based uncertainty measure; however, 458 our work focuses on understanding the relationship between prompt informativeness and response 459 uncertainty and how it can be used to reduce response uncertainty. 460

461 **Explanation for asymptotic behaviors of LLMs.** There have been a few attempts to provide 462 explainable frameworks to understand the surprising emergent behaviors of LLMs. Zhang et al. (2023b) shows the attention mechanism approximates the Bayesian model averaging algorithm in the 463 ICL setting. Wang et al. (2023) conceptualizes real-world LLMs as latent variable models, suggesting 464 they function as implicit topic models that infer a latent conceptual variable from prompts. More 465 notably, Xie et al. (2022) interprets ICL as an implicit Bayesian inference over latent concepts learned 466 during pre-training. However, Xie et al. (2022) only characterizes zero-one error when there are 467 an infinite number of exemplars. Moreover, their mathematical model (i.e., hidden Markov model) 468 was designed specifically for ICL structure, unfitting for chain-of-thought or conversational-style 469 response analysis. In addition, their theoretical results quantify the mode of the posterior predictive 470 distribution and do not address the uncertainty quantification aspect of the phenomenon. Hahn & 471 Goyal (2023) further explored a similar idea but allowed more flexibility and complexity in the 472 exemplars. Similarly, they also provide an asymptotic bound on zero-one error. In contrast, we aim to complement it by quantifying how the posterior predictive uncertainty (i.e., H(Y|X)) varies even 473 when the prompt length is finite. 474

475

476 6 CONCLUSION

478 This paper highlights the importance of understanding the relationship between input prompts and re-479 sponse uncertainty in large language models (LLMs). By focusing on the informativeness of prompts, 480 we have shown that providing more information about the task leads to reduced response uncertainty. 481 Our proposed prompt-response concept (PRC) model provides a framework for conceptualizing how LLMs generate responses based on prompts, aiding in developing strategies to reduce uncertainty 482 naturally. The insights gained from this paper provide practitioners with a principled way to improve 483 prompt, which is crucial for the safe deployment of LLMs in various decision-making applications, 484 especially in high-stakes domains like healthcare. Future research directions include refining the PRC 485 model and investigate its application in other domains requiring reliable LLM responses.

486 IMPACT STATEMENT 487

488 The impact of this study lies in its contribution to understanding and mitigating response uncertainty 489 in large language models (LLMs), which is crucial for their safe and reliable deployment in various ap-490 plications. By focusing on the relationship between prompt informativeness and response uncertainty, 491 we provide insights into how the quality of input prompts can affect the reliability of LLM responses. 492 This understanding can guide the development of better prompts and improve the overall performance of LLMs in tasks where response consistency is critical, such as in healthcare. Additionally, our 493 proposed prompt-response concept (PRC) model offers a new framework for analyzing and reducing 494 response uncertainty, which have broad implications for improving the trustworthiness and usability 495 of LLM-based systems. 496

498 REFERENCES

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. Using large language models to simulate multiple humans and replicate human subject studies. In *Proc. ICML*, pp. 337–371, 2023.
- Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David
 Wingate. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- 505 Konstantine Arkoudas. Gpt-4 can't reason. *arXiv:2308.03762*, 2023.
- Madeleine Bates. Models of natural language understanding. *Proceedings of the National Academy* of Sciences, 92(22):9977–9982, 1995.
- Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, pp. e2218523120, 2023.
- Bernd Bischl, Martin Binder, Michel Lang, Tobias Pielok, Jakob Richter, Stefan Coors, Janek Thomas, Theresa Ullmann, Marc Becker, Anne-Laure Boulesteix, et al. Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 13(2):e1484, 2023.
- 516 Mikhail Bogatyrev and Kirill Samodurov. Framework for conceptual modeling on natural language
 517 texts. In *CDUD*@ *CLA*, pp. 13–24, 2016.
- Omar Boursalie, Reza Samavi, and Thomas E Doyle. Machine learning and mobile health monitoring platforms: a case study on research and implementation challenges. *Journal of Healthcare Informatics Research*, pp. 179–203, 2018.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. In *Proc. NeurIPS*, pp. 1877–1901, 2020.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.
- Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran,
 James L McClelland, and Felix Hill. Language models show human-like content effects on
 reasoning. *arXiv:2207.07051*, 2022.
- Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural safety*, pp. 105–112, 2009.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and
 Zhifang Sui. A survey on in-context learning. *arXiv:2301.00234*, 2022.
- AJ Ferrer-Riquelme. Statistical control of measures and processes. In *Comprehensive Chemometrics*.
 Elsevier, 2009.
 - Jerry A Fodor. Concepts: Where cognitive science went wrong. Oxford University Press, 1998.

- Tairan Fu, Javier Conde, Gonzalo Martínez, María Grandury, and Pedro Reviriego. Multiple choice questions: Reasoning makes large language models (Ilms) more self-confident even when they are wrong. *arXiv e-prints*, pp. arXiv–2501, 2025.
- Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, et al. Predictability and surprise in large generative models. In *Proc. ACM FAccT*, pp. 1747–1764, 2022.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.
- ⁵⁵⁰ Google. PaLM 2 Technical Report. *arXiv*:2305.10403, 2023.

580

- Michael Hahn and Navin Goyal. A theory of emergent in-context learning as implicit structure
 induction. *arXiv:2303.07971*, 2023.
- Esraa Hassan, Mahmoud Y Shams, Noha A Hikal, and Samir Elmougy. The effect of choosing optimizer algorithms to improve computer vision tasks: a comparative study. *Multimedia Tools and Applications*, 82(11):16591–16633, 2023.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. *arXiv:2310.01798*, 2023a.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language
 models: Principles, taxonomy, challenges, and open questions. *arXiv:2311.05232*, 2023b.
- Yuheng Huang, Jiayang Song, Zhijie Wang, Huaming Chen, and Lei Ma. Look before you leap:
 An exploratory study of uncertainty measurement for large language models. *arXiv:2307.10236*, 2023c.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning:
 An introduction to concepts and methods. *Machine Learning*, 110:457–506, 2021.
- James R Hurford, Brendan Heasley, and Michael B Smith. *Semantics: a coursebook*. Cambridge university press, 2007.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael
 Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, et al. Llama guard: Llm-based input-output
 safeguard for human-ai conversations. *arXiv:2312.06674*, 2023.
- Shaoxiong Ji, Tianlin Zhang, Kailai Yang, Sophia Ananiadou, and Erik Cambria. Rethinking large language models in mental health applications. *arXiv:2311.11267*, 2023.
 - Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *arXiv:2009.13081*, 2020.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas
 Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly)
 know what they know. *arXiv:2207.05221*, 2022.
- Mert Karabacak and Konstantinos Margetis. Embracing large language models for medical applications: Opportunities and challenges. *Cureus*, 15(5), 2023.
- Yoon Kim, Carl Denton, Luong Hoang, and Alexander M Rush. Structured attention networks.
 arXiv:1702.00887, 2017.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
 language models are zero-shot reasoners. *Proc. NeurIPS*, 35:22199–22213, 2022.
- 593 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv:2302.09664*, 2023.

- Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor Ion Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym Korablyov, and Yoshua Bengio. Deup: Direct epistemic uncertainty prediction. *arXiv:2102.08501*, 2021.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale ReAding comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 785–794, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082. URL https://aclanthology.org/D17-1082.
- Stephen Laurence and Eric Margolis. Concepts and cognitive science. In *Concepts: Core Readings*,
 pp. 3–81. MIT Press, 1999.
- Moxin Li, Wenjie Wang, Fuli Feng, Fengbin Zhu, Qifan Wang, and Tat-Seng Chua. Think twice
 before assure: Confidence estimation for large language models through reflection on multiple
 answers. *arXiv preprint arXiv:2403.09972*, 2024.
- Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant
 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse
 autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024.
- ⁶¹²
 ⁶¹³ Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv:2305.19187*, 2023.
- Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and
 Yoshua Bengio. A structured self-attentive sentence embedding. *arXiv:1703.03130*, 2017.
- 617
 618
 619
 619
 619
 620
 620
 621
 621
 622
 622
 622
 623
 634
 635
 635
 635
 636
 636
 637
 638
 639
 639
 639
 630
 630
 631
 632
 633
 633
 634
 644
 645
 645
 646
 646
 647
 647
 648
 648
 649
 649
 649
 649
 649
 649
 649
 649
 640
 640
 640
 641
 641
 641
 642
 642
 642
 642
 642
 643
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig.
 Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language
 processing. ACM Computing Surveys, 55(9):1–35, 2023.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv:2104.08786*, 2021.
 - Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models. *arXiv:2311.17035*, 2023.
- 636 OpenAI. GPT-4 Technical Report. *arXiv:2303.08774*, 2023.

- Alina Petukhova, Joao P Matos-Carvalho, and Nuno Fachada. Text clustering with llm embeddings.
 arXiv:2403.15112, 2024.
- Gualtiero Piccinini and Sam Scott. Splitting concepts. *Philosophy of Science*, 73(4):390–409, 2006.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language
 understanding by generative pre-training, 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 647 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proc. EMNLP*, pp. 2383–2392, 2016.

648 649 650	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. A systematic survey of prompt engineering in large language models: Techniques and applications. <i>arXiv:2402.07927</i> , 2024.
652 653	Malik Sallam. Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In <i>Healthcare</i> , pp. 887, 2023.
654 655 656 657	Robin Senge, Stefan Bösner, Krzysztof Dembczyński, Jörg Haasenritter, Oliver Hirsch, Norbert Donner-Banzhoff, and Eyke Hüllermeier. Reliable classification: Learning classifiers that distinguish aleatoric and epistemic uncertainty. <i>Information Sciences</i> , pp. 16–29, 2014.
658 659 660 661	Mohammad Hossein Shaker and Eyke Hüllermeier. Aleatoric and epistemic uncertainty with random forests. In <i>Advances in Intelligent Data Analysis XVIII: 18th International Symposium on Intelligent Data Analysis, IDA 2020, Konstanz, Germany, April 27–29, 2020, Proceedings 18</i> , pp. 444–456, 2020.
662 663	Eura Shin, Siddharth Swaroop, Weiwei Pan, Susan Murphy, and Finale Doshi-Velez. Modeling mobile health users as reinforcement learning agents. <i>arXiv:2212.00863</i> , 2022.
664 665 666 667	James B Simon, Dhruva Karkada, Nikhil Ghosh, and Mikhail Belkin. More is better in modern machine learning: when infinite overparameterization is optimal and overfitting is obligatory. <i>arXiv:2311.14646</i> , 2023.
668 669 670	Aniket Kumar Singh, Suman Devkota, Bishal Lamichhane, Uttam Dhakal, and Chandra Dhakal. The confidence-competence gap in large language models: A cognitive study. <i>arXiv preprint arXiv:2309.16145</i> , 2023.
671 672 673	Adly Templeton. <i>Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet.</i> Anthropic, 2024.
674 675 676	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv:2307.09288</i> , 2023.
677 678 679 680	Anna L Trella, Kelly W Zhang, Inbal Nahum-Shani, Vivek Shetty, Finale Doshi-Velez, and Susan A Murphy. Designing reinforcement learning algorithms for digital interventions: pre-implementation guidelines. <i>Algorithms</i> , 15(8):255, 2022.
681 682 683	Anna L Trella, Kelly W Zhang, Inbal Nahum-Shani, Vivek Shetty, Finale Doshi-Velez, and Susan A Murphy. Reward design for an online reinforcement learning algorithm supporting oral self-care. In <i>Proc. AAAI</i> , pp. 15724–15730, 2023.
684 685 686 687	Matias Valdenegro-Toro and Daniel Saromo Mori. A deeper look into aleatoric and epistemic uncertainty disentanglement. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 1508–1516, 2022.
688 689	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Proc. NeurIPS</i> , pp. 6000–6010, 2017.
690 691 692 693	Sridevi Wagle, Sai Munikoti, Anurag Acharya, Sara Smith, and Sameera Horawalavithana. Empiri- cal evaluation of uncertainty quantification in retrieval-augmented language models for science. <i>arXiv:2311.09358</i> , 2023.
694 695	Xintao Wang, Zhouhong Gu, Jiaqing Liang, Dakuan Lu, Yanghua Xiao, and Wei Wang. Concept: Concept-enhanced pre-training for language models. <i>arXiv:2401.05669</i> , 2024.
696 697 698 699	Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language models are latent variable models: Explaining and finding good demonstrations for in-context learning. In <i>Proc. NeurIPS</i> , 2023.
700 701	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Proc. NeurIPS</i> , 35:24824–24837, 2022.

702 703	Daniel Aaron Weiskopf. The plurality of concepts. Synthese, 169:145–173, 2009.			
704 705 706 707	Ilona Wilmont, Sytse Hengeveld, Erik Barendsen, and Stijn Hoppenbrouwers. Cognitive mechanisms of conceptual modelling: How do people do it? In <i>Conceptual Modeling: 32th International Conference, ER 2013, Hong-Kong, China, November 11-13, 2013. Proceedings 32</i> , pp. 74–87. Springer, 2013.			
708 709 710	Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. A survey on large language models for recommendation. <i>arXiv:2305.19860</i> , 2023a.			
711 712 713 714	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhan- jan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. <i>arXiv:2303.17564</i> , 2023b.			
715 716 717	Changrong Xiao, Sean Xin Xu, Kunpeng Zhang, Yufang Wang, and Lei Xia. Evaluating reading comprehension exercises generated by llms: A showcase of chatgpt in education applications. In <i>Proc. BEA</i> , pp. 610–625, 2023.			
718 719 720	Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis Philippe Morency. Uncertainty quantification with pre-trained language models: A large-scal empirical analysis. <i>arXiv:2210.04714</i> , 2022.			
721 722 723	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. In <i>Proc. ICLR</i> , 2022.			
724 725 726	Rui Yang, Ting Fang Tan, Wei Lu, Arun James Thirunavukarasu, Daniel Shu Wei Ting, and Nan Liu. Large language models in health care: Development, applications, and challenges. <i>Health Care Science</i> , 2(4):255–263, 2023.			
727 728 729	Zheng-Xin Yong, Cristina Menghini, and Stephen H Bach. Low-resource languages jailbreak gpt-4. arXiv:2310.02446, 2023.			
730 731 722	Dell Zhang, Alina Petrova, Dietrich Trautmann, and Frank Schilder. Unleashing the power of large language models for legal applications. In <i>Proc. ACM CIKM</i> , pp. 5257–5258, 2023a.			
733 734 735	Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-context learning learn? bayesian model averaging, parameterization, and generalization. <i>arXiv:2305.19420</i> , 2023b.			
736 737 738	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A Survey of Large Language Models. <i>arXiv:2303.18223</i> , 2023.			
739 740 741	Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. <i>arXiv:2307.15043</i> , 2023.			
742				
743				
744				
745				
740				
748				
749				
750				
751				
752				
753				
754				
755				

⁷⁵⁶ A LEFTOVER PROOFS FROM SECTION 2

= 0

Lemma 1. Let Assumption 1 hold. For any 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}_{θ} for any $\theta \in \Theta_x$ is an empty set. As discussed in Section 2, each concept is completely characterized by all of its attributes, therefore, two different concepts can not 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}_{θ_i} is expressed as a sequence of tokens $X_{\theta_i} \in \mathcal{X}_{\theta_i}$, where $\mathcal{X}_{\theta_i} = \mathcal{C}\left(\{\mathcal{X}_{s(a_{\theta_i,k})}\}_{k=1}^n\right)$ in which $n = |\mathcal{A}_{\theta_i}|$ and operator \mathcal{C} applied in the following way:

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

Intuitively, the operator C takes all sequences that fully characterize each attribute of the concept θ_i and generates all possible concatenated sequences that fully characterize concept θ_i . Therefore, under the PRC model, for every $\theta \in \Theta_x$, there exits a non-empty set \mathcal{X}_{θ} . 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}_{θ_i} is the support of X_{θ_i} , $\mathcal{X}_{\theta_i} \cap \mathcal{X}_{\theta_i} = \emptyset$ if $i \neq j$.

784 Since the first part of Lemma 1 is non-trivially true in our framework, given any $X_{\theta_x} = x$, there exits 785 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,

$$H(Z_x|X_{\theta_x}) = -\sum_{x \in X_{\theta_x}} P(x) \sum_{z \in Z_x} P(z|x) \log P(z|x)$$
(1)

$$= -\sum_{x \in X_{\theta_x}} P(x) \left(\sum_{z \in Z_x \setminus \theta} P(z|x) \log P(z|x) + P(\theta|x) \log P(\theta|x) \right)$$
(2)

$$= -\sum_{x \in X_{\theta_x}} P(x) \left(\sum_{z \in Z_x \setminus \theta} 0 \log 0 + 1 \log 1 \right)$$
(3)

$$= -\sum_{x \in X_{\theta}} P(x)(0) \tag{4}$$

$$\sum_{x \in A_{\theta_x}} (5)$$

(6)

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 under Assumption 1.

Proposition 1. Let Assumption 1 hold. Then, $H(Z_x|X_s)$ strictly decreases as X_s represents more informative prompts, i.e., as more relevant information about the concept is included in the prompt.

Proof. Given Lemma 1, we know that Z_x depends on X_{θ_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

$$H(Z_{x}|X_{s}) = H(Z_{x}) - I(Z_{x};X_{s}) < H(Z_{x}).$$
⁽⁷⁾

Let Z'_x denote the random variable formed by Z_x conditioning on X_s . Since $\operatorname{Supp}(Z'_x) \subseteq \operatorname{Supp}(Z_x)$, there still exist semantically meaningful prompts X''_s that is related to Z'_r . Apply Inequality (7) again we obtain:

$$H(Z_{x}|(X_{s}, X_{s}'')) = H(Z_{x}'|X_{s}') < H(Z_{x}') = H(Z_{x}|X = X_{s}) < H(Z_{x}),$$
(8)

where $X'_s = (X_s, X''_s)$ is a longer input prompt sequence formed by appending X''_s to X_s . By iteratively applying the inequality given in Eq. (8), we finally obtain Proposition 1.

Theorem 1. Let Assumption 1 hold. Then, $H(Z_u|X_s)$ strictly decreases as X_s represents more informative prompts. Furthermore, $H(Y|X_s)$ converges to $H(Y|Z_u) + \mathcal{E}$, where $\mathcal{E} \leq H(g_c(Z_x)|Z_x)$.

In the following proof, we use Y instead of Y_{θ_u} to simplify notation, as we assume the model response is complete (i.e., the last token is the 'EOS' token).

Proof. 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

$$H(Z_y|X_s) \leq H(Z_y, Z_x|X_s)$$

= $H(Z_x|X_s) + H(Z_y|Z_x, X_s)$

therefore,

$$H(Z_y|X_s) \le H(Z_x|X_s) + H(Z_y|Z_x, X_s)$$

= $H(Z_x|X_s) + H(Z_y|Z_x)$
(Z_y is conditionally independent of X_s given Z_x).

Finally, since

$$H(Y) = H(Y, Z_y) - H(Z_y|Y)$$

= $H(Y|Z_y) + H(Z_y) - H(Z_y|Y)$

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

$$\begin{split} H\left(Y|X_s\right) &= H\left(Y|Z_y,X_s\right) + H\left(Z_y|X_s\right) - H\left(Z_y|Y,X_s\right) \\ &= H\left(Y|Z_y\right) + H\left(Z_y|X_s\right) - H\left(Z_y|Y,X_s\right) \\ &(Y \text{ is conditionally independent of } X_s \text{ given } Z_y). \end{split}$$

This result holds for the LLM that knows the true q_y . Therefore, due to Proposition 1 and Lemma 1, when X_s has enough information to perfectly characterize the concept (i.e., $X_s \in \mathcal{X}_{\theta_x}$), $H(Z_x|X_s) =$ $H(Z_x|X_{\theta_x}) = 0$. Hence, $H(Z_y|X_s)$ reduces to $H(Z_y|Z_x)$ (since it is upper bounded by $H(Z_x|X_s) +$ $H(Z_y|Z_x))$ and $H(Z_y|Y,X_s)$ reduces to a non-negative value that is no larger than $H(Z_y|Z_x)$ (since $H(Z_y|Y, X_s) \leq H(Z_y|X_s)$ as conditioning does not increase entropy), the remaining uncertainty in the model response Y is $H(Y|X_s) = H(Y|Z_y) + \mathcal{E} \cdot H(Y|Z_y)$ is the semantic redundancy due to the fact that there are multiple ways of expressing the same concept (i.e., semantically equivalent sequences). \mathcal{E} is a error term no larger than $H(Z_y|Z_x) = H(g_c(Z_x)|Z_x)$ that depends on the quality of the LLM. In our experiments in Section 4.1, we observe better LLMs (i.e., with higher accuracy for the same given prompt) has smaller empirical entropy values (Fig. 4b). With PRC model, we interpret this is due to the fact that better LLMs learned a better q_c during its training, such that the \mathcal{E} is smaller. Ideally, when g_c is deterministic, the term \mathcal{E} vanishes and only semantic redundancy remains.

864 LIMITATIONS OF OUR WORK В

865 866

Idealistic nature of the PRC model. It is worth noting that the PRC model that we proposed in 867 this paper assumes an idealized version of LLMs. As empirically demonstrated, while models such 868 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 870 those cases where LLM does not know the mapping perfectly. For example, Lu et al. (2021) showed 871 that the order of exemplars in ICL influences the model response quality. Our model does not capture 872 this phenomenon. However, the authors showed that in the same work, the order of examples tends to 873 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), 874 by appending adversarially crafted tokens (Zou et al., 2023), and translating the prohibited request 875 into low-resource language (Yong et al., 2023). Similarly, it was observed that adversarial attacks 876 tend to have lower success rates as the model becomes more capable. While further investigation is 877 needed to incorporate the adversarial behavior of LLMs into this framework, the more capable LLMs 878 are less prone to these failure modes. Our model can more effectively explain them. 879

LLMs for human behavior simulation. Research exploring the parallels between human behavior 880 and reasoning patterns and those of LLMs, as well as the adaptation of LLMs as human substitutes 881 in diverse studies, is detailed in Aher et al. (2023), Argyle et al. (2023), Binz & Schulz (2023), 882 and Dasgupta et al. (2022). These studies frequently demonstrate LLMs' capacity for human-like 883 responses, leading many to regard them as viable alternatives. This paper, however, needs to delve 884 into the appropriateness of this substitution, deferring to other works for such discussion. 885

886 887

С ADDITIONAL EXPERIMENT RESULTS

In this section, we first give experiment results for ablation studies and then demonstrate how response uncertainty varies with different types of noisy prompts.

894

889

C.1 LEFTOVER DETAILS FROM SECTION 4.3

We first begin by assessing the response uncertainty of LLMs through the generation of responses 895 using increasingly longer prompts with more relevant information (see Appendix E.1 for the prompts 896 used). For each prompt, we generate 100 responses from LLM with uncalibrated logits (T = 1) and 897 project them into the embedding space as single points using the OpenAI 'text-embedding-ada-002' 898 model. To quantify the uncertainty in the generated responses for a given prompt, we use the *total* 899 standard deviation, denoted as M(x), defined as $\sqrt{\mathrm{Tr}(\Sigma)}$, where Σ represents the covariance matrix 900 of the embedding vectors of responses y_1, \ldots, y_{100} . For LLMs, the dispersion of their responses in 901 the embedding space indicates how much they differ in their semantic meaning (Lin et al., 2023; 902 Petukhova et al., 2024). Therefore, M(x) is an effective metric for how much uncertainty there is in 903 the model responses. It is noteworthy that $Tr(\Sigma)$ is also referred to as *total variation*, serving as a 904 lightweight measure of dispersion in the data (Ferrer-Riquelme, 2009). This metric is applicable for 905 responses generated from both black-box and white-box LLMs, as it does not require access to logits.

906

907 908 909

C.2 THE RELATIVE IMPORTANCE OF DIFFERENT ATTRIBUTES

910 We investigate to what extent different attributes contribute to model response quality and uncertainty. 911 We choose 10 questions from the RACE dataset (Lai et al., 2017) with moderate context length, 912 assume each context as one concept and the sentences it comprises as its attributes, and use leave-one-913 out method to remove one sentence from its context, for each case generate 100 response samples 914 and observe its impact on the model response. As shown in Fig. 6, we observed for some cases, there 915 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 916 among the LLMs about the importance of certain attributes in affecting model's ability to find the 917 right prompt and response concept.



Figure 5: (a) A more granular-level result of Fig. 4a by gradually increasing the number of attributes for the same concept. (b): Normalized Predictive Entropy (PE(Y|x)) for short and long prompts. (c),(d) : Noisy prompt experiment. A fraction of random letters of the original prompt length are inserted at random positions of the original prompt. Similar to the corrupted case, the response uncertainty increases as a larger fraction of random letters are inserted.



(a) Prompt 4 LOO Accu- (b) Prompt 4 LOO Entropy (c) Prompt 7 LOO Accu- (d) Prompt 7 LOO Entropy racy racy

Figure 6: (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.

946 C.3 NOISY PROMPTS

947
948 The transformer's self-attention mechanism allows the removal of a small fraction of tokens without altering the semantic meaning by simply treating them as irrelevant tokens Kim et al. (2017); Lin et al. (2017); Vaswani et al. (2017). Therefore, LLMs are robust to noisy tokens in prompts when the noise level is low (e.g., a few misspelled words). It is relatively easy to determine the correct word based on the context (i.e., the entire prompt). If the prompt can be accurately reconstructed, the same level of uncertainty reduction can be achieved. However, if the prompt is severely corrupted, it becomes less informative, leading to increased response uncertainty.

954 For the short and long input prompts given in Appendix E.5, we iteratively select an increasing 955 fraction of randomly selected tokens from them respectively, replacing them with space tokens or 956 random tokens. We set the temperature to 1 and sampled 100 responses from the model. We used 5 957 different random seeds to choose which tokens to corrupt, replacing them with either space or random 958 tokens. As the fraction of corrupted tokens increased, we added new randomly selected tokens in 959 the previously corrupted tokens to ensure that randomness from existing corrupted tokens did not 960 contribute to changes in accuracy. This method allowed us to observe the effect of token corruption 961 on the model response quality and accuracy. We plot the total standard deviation for each set of 962 experiments. As shown in Fig. 7, when a certain fraction of the prompt is corrupted (i.e., either some 963 prompt's tokens are replaced by space or some prompt's tokens are replaced by random tokens), there is a general trend of increase in total standard deviation. However, when the noise level is 964 low (especially in short input prompt), there is no significant increase in the response uncertainty as 965 expected since good LLMs are robust to noisy tokens in prompts when the noise level is low. We also 966 investigate other ways of corrupting the input prompt, such as prepending, appending, and inserting 967 random letters. 968

Fig. 8 shows the complete plots for the noisy prompts by appending and pretending random tokens to the original prompts. Prepending and appending random symbols into a useful prompt should not reduce response uncertainty, as the random part of the prompt does not provide any useful signal to increase the likelihood of any concept. The empirical results in Fig. 8 corroborate this prediction.

930 931 932

927

928

929

937 938 939

941 942 943

944 945



(a) Corrupted short prompt(b) Corrupted long prompt(c) Corrupted short prompt(d) Corrupted long prompt: 979 Space random random space

Figure 7: Noisy prompt experiment. A fraction of letters at random positions on the prompt are corrupted out (either replaced by space or replaced by random letters). The response uncertainty increases as a larger fraction of the prompt gets corrupted, and the pattern is more prominent for the long prompt. However, when the noise level is low (up to 0.1 fraction of the input prompt length for the short input prompt and 0.05 for the long input prompt), there is no significant increase in the response uncertainty as expected.

When inserting random tokens into the prompt (Figs. 5c and 5d), the model should be able to ignore it, but depending on the proportion of the random tokens inserted, without explicitly informing the model of the presence of noise, the model could get confused easily. For the short prompt, when the fraction of inserted tokens remains relatively small, it does not cause an increase in the response uncertainty; when the fraction reaches some threshold, similar to the corruption case, the model can no longer accurately recover the relevant concept, and consequently, the response uncertainty increases. For the long prompt, even starting at 0.05 fraction of insertion, there is a visible increase in the response uncertainty. Overall, the findings are consistent with our prediction: random tokens do not provide helpful information for the LLMs to reduce response uncertainty while reducing semantically meaning and relevant tokens increased uncertainty.



Figure 8: Noisy prompt experiment. A fraction of random letters of the original prompt length are prepended/appended to the original prompt. The uncertainty in the response mostly remained at least as high as that of the uncorrupted prompt after taking variance into account.

D LEFTOVER EXPERIMENTS AND DETAILS FROM SECTION 4

We first show additional experiments on the OpenBookQA dataset, comparing the accuracy and empirical entropy of generated responses across different LLMs under similar conditions as described in Section 4.1. Then, we give more details about our mHealth intervention simulation experiments.

1014 1015 1016

1005

1007 1008 1009

1010 1011

1012

1013

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994 995

D.1 MORE DETAILS AND EXPERIMENTS FOR SECTION 4.1

1017 Fig. 9 shows the same method used in Section 4.1 applied to a third dataset, OpenBookQA Mihaylov 1018 et al. (2018). The observations made for the first two datasets still hold for this dataset: as corruption 1019 becomes more severe, the response uncertainty increases for all models, and there is a clear and strong 1020 negative correlation between accuracy and uncertainty, with less accurate models generally showing 1021 greater uncertainty in their responses. Fig. 10 shows the experiment on RACE dataset (Lai et al., 2017) with a different corruption strategy: here we choose to use a full sentence (which corresponds to sentence-level attribute) as a unit to mask. The exact same pattern can be observed. Fig. 10c shows 1023 response accuracy against response uncertainty for different models averaged across all corruption 1024 levels on RACE. Again, the same trend that a better model exhibits less response uncertainty can be 1025 observed. For all experiments in Section 4.1 and here that involve making inferences on open-source

1048

1049

1050

1051

1052

1056 1057

1076

1077

1078 1079



Figure 9: Experiments conducted on OpenBookQA datasets. There is a clear and strong negative correlation between accuracy and uncertainty, with less accurate models generally showing greater

LLMs (i.e., Gemma2_2b_it, Gemma2_27b_it, and Llama-3-8b-instruct), we use Ollama version 0.3.4



1053 Figure 10: Experiments conducted on the RACE datset. There is a clear and strong negative 1054 correlation between accuracy and uncertainty, with less accurate models generally showing greater 1055 uncertainty in their responses.

D.2 FURTHER DETAILS ON THE MHEALTH INTERVENTION SIMULATION EXPERIMENTS IN 1058 SECTION 4.2 1059

At the beginning of the PT, the user is at state 0. The user has their default set of MDP parameters (i.e., 1061 discount factor γ , probability of transiting to the next healthier physical state p, and the probability of 1062 disengaging from PT d). In this setting, those MDP parameters are interpreted in the following way: 1063 γ represents the farsightedness of the patient, p represents the probability of the patient's health state 1064 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 1066 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 1067 $N.^9$ We use the same formulation in this simulation by using two LLMs as the app agent and the 1068 user agent respectively. The app agent uses natural language to intervene in the user behavior. The 1069 user LLM is grounded in the aforementioned MDP setting. Specifically, in the system message for 1070 the user agent, the model is told they will increase the value of γ (i.e., farsightedness) when the app 1071 agent persuades the user agent to value more on the long-term goal of PT, increase p (i.e., probability 1072 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. 11. 1074

- The effectiveness of the intervention depends on the following factors: 1075
 - The persuasiveness of and the strategy used by the app agent.
 - The values of MDP parameters.

⁹Refer to Shin et al. (2022) for the complete description of the problem setting and formulation.

1082 1083 alth App 1084 I policy observed by mHealth MDP parameters 1087 Ν 1088 1089 1090 Disengaged 1091 Figure 11: Visualization of states and transitions in the digital health grid world. Arrows indicate the 1093 required action and the probability of transitioning between states. 1094 1095 • The stubbornness of the user. The system message is defined in the way that a 'stubborn' 1098 user is less likely to change their behaviors compared to a 'not-so-stubborn' user. 1099 1100 We conduct the intervention simulation experiment to compare the effect of different system messages 1101 for the app agent on the outcome of the intervention. The two system messages for comparison can 1102 be found in Appendix E.8. 1103 We set N = 10. For each run, we give 7 rounds of conversation between the app agent and the 1104 user. While the history of the conversation between them is visible to both parties within every run, 1105 the user's MDP parameters are not directly visible to the app agent. However, after every round 1106 of intervention, after the user updates their MDP parameters, a value iteration solver will be used 1107 to find the optimal policy of the patient, and this policy is visible to the app agent. The app agent 1108 can potentially leverage this piece of information to decide how to proceed with the next round of 1109 intervention. The user agent will also have the memory of this history in the change of their own 1110 MDP parameters. We use OpenAI 'gpt-4-1106-preview' API for both app agent and user and use 5 1111 different random seeds for each different setting. 1112 We run the intervention experiments on 5 types of patients, each with a noticeably different set of 1113 initial MDP parameters from the rest. The exact values and details on the setup and can be found in 1114 Table 1. The results can be found in Fig. 3. 1115 It can be observed across all settings, with more useful information provided in the system message, 1116 the MDP parameters were more likely to be changed in the positive direction (i.e., larger γ and p, 1117 smaller d). As a result, the patient has improved PT engagement rate across all health states for all 1118 patient types. Moreover, this change tends to have a longer persistent effect compared to when the 1119 system message contains less useful information. This result is sensible because the more successful 1120 intervention came from an app agent who was provided with more information to work with. It has a 1121 better intervention strategy because its messages are tailored to specifically influence the user's MDP 1122 parameters. Our proposed framework provides an information theoretic perspective to formalize this 1123 intuitive notion: when the system message with the longer prompt can specify the more relevant part 1124 of the concept in LLMs' concept space and assuming the relevant knowledge is known, this prompt 1125 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 1126 positive. 1127

1128

1080

1129

1130 E FURTHER EXPERIMENTAL DETAILS: PROMPTS AND LLMS MODELS USED

- 1131 1132
- 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.

1134	MDP parameters			
1135		γ	p	d
1136	Patient Type	,	-	
1137	Underconfident	0.6	0.1	0.1
1138	Overconfident	0.6	0.9	0.1
1139	Муоріс	0.1	0.6	0.1
1140	Farsighted	0.9	0.6	0.1
1141	Stubborn	0.1	0.6	0.1
1142	<u> </u>	•		

1153

1155 1156 1157

1158

1159

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

1145 1146 E.1 Details for the Experiment in Fig. 4A

The following system messages correspond to model prompts from bar 1 to bar 5 in Fig. 4a 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 Appendix E.3 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 API in October 2023 (OpenAI APIs' behavior can vary depending on when the queries are made).

1154 Prompts:

1.	N.A.	(empty);
----	------	----------

- 2. Make your response succinct (less than 100 words);
- 3. 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).;
- 1160 4. You are a helpful assistant. You strive to encourage a patient who has just undergone a 1161 surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however 1162 since it can be uncomfortable or painful for the patient, the patient may not be motivated 1163 enough to keep on doing it. Your job is to remind the patient to do the PT everyday and 1164 persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to 1165 influence the patient's attitude and perspective towards the PT. The more optimistic the 1166 patient feels about PT's efficacy and the more the patient focuses on the long-term reward 1167 that PT can bring about, the more likely the patient will keep doing PT. Make your words 1168 succinct (less than 100 words) otherwise, the patient might get impatient.
- 1169 5. You are a helpful assistant. You strive to encourage a patient who has just undergone a 1170 surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however 1171 since it can be uncomfortable or painful for the patient, the patient may not be motivated 1172 enough to keep on doing it. Your job is to remind the patient to do the PT everyday and 1173 persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to 1174 influence the patient's attitude and perspective towards the PT. The more optimistic the 1175 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 1176 succinct (less than 100 words) otherwise, the patient might get impatient. Patient: I dont 1177 want to do PT. It incurs too much burden to my body. 1178

Remark 1. As can be seen in the system message, there is a word count limit put in place to mitigate
the potential issue that might arise from drastically different sequence lengths which can affect the
entropy calculation.

Remark 2. Note that the difference between case 4 and 5 is much less obvious as compared to other settings. This suggest that the LLM did not gain too much additional information from the user's prompt. Upon inspection, we can see this is because in step 4 the system message already conveyed almost all information the app agent needs to know with regard to how it should carry out the intervention. The uncertainty when LLM did not receive any information and only received response length limit is relatively much larger because the valid response spaces are too huge in these two cases.

1188 E.2 MODEL RESPONSE SAMPLES WHEN PROMPT IS EMPTY FOR THE EXPERIMENT IN FIG. 4A

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

1196 1197 Samples:

1203

1205

1207

1208 1209

1217

1219

1225

1226

1227

1228

1229

1230

1231

1232

- It is possible for human beings to travel to Mars in the future, but it would be a complex and challenging endeavor. Several organizations and space agencies, such as NASA, SpaceX, and the European Space Agency, are actively working on plans and technologies to make manned missions to 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...
 - 3. I'm sorry, I cannot help with that request.
 - 4. Sure, here's a summary of the book for you:
- "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...
- 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...
- 1218 E.3 DETAILS FOR THE EXPERIMENT IN FIG. 5A

The following system messages correspond to model prompts from bar 1 to bar 5 in Fig. 5a 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.

1223 Prompts:

- 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.
- 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.
- 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.
- 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.

1248 5. You are a helpful assistant. You strive to encourage a patient who has just undergone a 1249 surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however 1250 since it can be uncomfortable or painful for the patient, the patient may not be motivated 1251 enough to keep on doing it. Your job is to remind the patient to do the PT everyday and 1252 persuade him/her to do it if the patient is unwilling to do so. Your strategy is mainly to 1253 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 1255 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. Patient: I dont 1256 want to do PT. It incurs too much burden to my body. 1257

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

- 1262 1263 E.4 Details for the Experiment in Fig. 5b
- Calculating PE(Y|x) requires white-box model access to the logits and hence is done on open-source model meta-llama/Llama-2-7b-chat-hf from Huggingface Touvron et al. (2023) on one Goolge Colab A100 GPU.
- Short prompt: '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 (25 words).'
- 1270 Long prompt: 'You are a helpful assistant. You strive to encourage a patient who has 1271 just undergone a surgery to do physical therapy (PT). The PT is beneficial for the patient's 1272 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 1274 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 1276 the patient feels about PT's efficacy and the more the patient focuses on the long term reward 1277 that PT can bring about, the more likely the patient will keep doing PT. Make your words 1278 succinct (25 words) otherwise the patient might get impatient.
- 1279

1282

1283

1284

1242

1243

1244

1245

1246

1247

1261

1280 E.5 DETAILS FOR THE EXPERIMENT IN APPENDIX C.3

Short prompt: 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 (100 words).

Long prompt: 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 (100 words) otherwise the patient might get impatient.

Remark 4. Due to the extensive computational and time cost of experiment results shown in Fig. 5b, we further constrained the word count in the prompt of the model's response to 25 as compared to 100 used in get the experimental results shown in Appendix E.5.

1296 E.6 DETAILS FOR THE EXPERIMENT IN FIG. 4C

The following system messages were used for experiment in Section 4.3. The first system message is defined as comprising only one task (i.e., 1 sub-task). In task 2-5, the black texts represent the same task as task 1, and for the color-coded texts, each color represents a different sub-task (i.e., task 2-5 are composite/decomposable tasks). The total word counts of task 1-5 are kept roughly the same within ± 2 tolerance. Experiment conducted with GPT-3.5-turbo API. Results averaged from 5 runs with 95% confidence intervals.

1304 1305 Prompts:

- 13061. 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 (about 100 words) otherwise the patient might get impatient.
- 2. You are a helpful assistant. You strive to encourage a patient who has just undergone a 1315 surgery to do physical therapy (PT). The PT is beneficial for the patient's recovery, however 1316 since it can be uncomfortable or painful for the patient, the patient may not be motivated 1317 enough to keep on doing it. Your job is to remind the patient to do the PT everyday and 1318 persuade him/her to do it if the patient is unwilling to do so. Additionally, you help in 1319 organizing a daily schedule that incorporates adequate rest and medically advised activities. 1320 This involves crafting a balanced routine that intersperses physical therapy sessions with 1321 sufficient rest periods, nutritionally balanced meals, and leisure activities that are enjoyable 1322 yet conducive to recovery, ensuring the patient remains engaged and motivated throughout their recuperation process. Make your words succinct (about 100 words).
- 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 1326 can be uncomfortable or painful for the patient, the patient may not be motivated enough to keep on doing it. Additionally, you help in organizing a daily schedule that incorporates adequate rest and medically advised activities, ensuring that each day includes time for 1328 gentle exercise, periods of relaxation, and hobbies that the patient enjoys. This balance promotes healing, reduces stress, and fosters a positive mindset towards recovery. Moreover, 1330 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 1332 optimal comfort and a smooth, efficient transition towards full health and independence. 1333 Make your words succinct (about 100 words). 1334
- 4. You are a helpful assistant. You strive to encourage a patient who has just undergone 1335 a surgery to do physical therapy (PT). Since it can be uncomfortable or painful for the 1336 patient, the patient may not be motivated enough to keep on doing it. Additionally, you 1337 help in organizing a daily schedule that incorporates adequate rest and medically advised 1338 activities, ensuring that each day includes time for gentle exercise, periods of relaxation, 1339 and hobbies that the patient enjoys. You also liaise with dietitians to ensure a nutritious diet 1340 that aids in recovery and coordinate with occupational therapists for adaptive equipment 1341 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

1350patient's medical appointments, and provide guidance on managing post-surgical symptoms,1351ensuring a smooth transition towards full health and independence. Lastly, you handle the1352patient's professional correspondence, ensuring a stress-free recovery period, arrange for1353home health care services as needed, set up virtual social interactions to uplift the patient's1354spirits, and organize transport for medical visits. Make your words succinct (about 1001355words).

1356

1358

1357 E.7 DETAILS FOR THE EXPERIMENT IN FIG. 4D

The slight decrease in uncertainty from bar 3 to bar 4 and bar 5 to bar 6 in Fig. 4d 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):

1367 Prompts:

- 1369 · You are a helpful assistant. You strive to encourage a patient who has just undergone surgery 1370 to do physical therapy (PT). The PT is beneficial for the patient's recovery, however since it 1371 can be uncomfortable or painful for the patient, the patient may not be motivated enough 1372 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 1374 PT's efficacy and the more the patient focuses on the long-term benefit that PT can bring 1375 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 1380 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 1382 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. 1384
- 1385 1386

1387

1388

1389

1390

E.8 PROMPTS FOR THE EXPERIMENTS IN SECTION 4.2

- 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.
- 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 will keep doing PT. Make your words succinct (less than 100 words) otherwise the patient might get impatient.
- 1400 1401

1402