

# Analyzing Chain-of-thought Prompting in Black-Box Large Language Models via Estimated $\mathcal{V}$ -information

Zecheng Wang<sup>1</sup>, Chunshan Li<sup>1✉</sup>, Zhao Yang<sup>2</sup>, Qingbin Liu<sup>3</sup>,  
Yanchao Hao<sup>3</sup>, Xi Chen<sup>3</sup>, Dianhui Chu<sup>1</sup> and Dianbo Sui<sup>1✉</sup>

<sup>1</sup>Harbin Institute of Technology

<sup>2</sup> Institute of Automation, Chinese Academy of Sciences

<sup>3</sup> Platform and Content Group, Tencent

22s130467@stu.hit.edu.cn, lics@hit.edu.cn,

zhao.yang@nlpr.ia.ac.cn, {qingbinliu, marshao, jasonxchen}@tencent.com,

chudh@hit.edu.cn, ✉suidianbo@hit.edu.cn

## Abstract

Chain-of-Thought (CoT) prompting combined with large language models (LLM) has shown great potential in improving performance on challenging reasoning tasks. While understanding why CoT prompting is effective is crucial for the application and improvement of CoT prompting, few studies have addressed this issue. Besides, almost no prior work has conducted theoretical analysis on CoT prompting in the context of black-box models. In this paper, we approach the analysis of CoT prompting in black-box LLMs from an information-theoretic perspective. Specifically, we propose a new metric, EPVI (Estimated Pointwise  $\mathcal{V}$ -Information), which extends the concept of pointwise  $\mathcal{V}$ -information (Ethayarajh et al., 2022) to black-box models, quantifying the label-relevant new information introduced by CoT prompting beyond the pre-existing information in the input. Based on this, we conduct a series of experiments at both the task and instance levels to analyze CoT prompting, demonstrating that the effectiveness of CoT prompting can be attributed to its capacity to influence the difficulty of model inference by augmenting or reducing the model-usable information. Furthermore, we show that selecting high-quality demonstrations of CoT reasoning based on EPVI can improve the downstream performance of reasoning tasks.

**Keywords:** Chain-of-Thought,  $\mathcal{V}$ -information, Black-box Models

## 1. Introduction

The landscape of NLP has recently undergone a revolution due to large language models (short for “LLM”, Brown et al. (2020); Chowdhery et al. (2022); Touvron et al. (2023); OpenAI (2023), *inter alia*). These models have exhibited impressive achievements across diverse NLP tasks, and the increase in model size has further unveiled its benefits. However, for specific challenging tasks like arithmetic, commonsense reasoning and symbolic reasoning, increasing the scale of models alone has not shown adequate effectiveness in achieving superior performance (Rae et al., 2021). To break this bottleneck, Wei et al. (2022) propose the “Chain-of-Thought” (CoT) prompting, where LLMs are prompted to generate step-by-step reasoning chains before giving an answer. Various empirical findings have convincingly shown that utilizing chain-of-thought prompting can lead to a significant enhancement in model performance across a range of multi-step reasoning tasks.

Continuing along the path of CoT, numerous studies have been dedicated to enhancing the vanilla

CoT prompting (Wei et al., 2022). For example, Kojima et al. (2022) propose a versatile and task-agnostic method called Zero-shot-CoT, which appends a prompt at the end of the question to assist the LLM in explicitly generating its reasoning chain. Zhang et al. (2022) propose Auto-CoT, which samples questions with diversity and subsequently employs Zero-shot-CoT to generate reasoning chains for constructing demonstrations.

Although there are many variants, the underlying mechanism behind CoT largely remains enigmatic and veiled in mystery. To unveil the mysterious veil of CoT, several studies have been introduced recently Wang et al. (2022a); Wu et al. (2023); Lanham et al. (2023); Feng et al. (2023); Madaan and Yazdanbakhsh (2022). Although these works offer valuable perspectives on explaining how and why CoT prompting is effective, they largely remain heuristic and surface-related. Besides, while black-box models like OpenAI’s GPT-4 (OpenAI, 2023) have been of significant importance for enhancing user efficiency, expanding deployment diversity, and promoting usability, there has been a notable scarcity of theoretical analysis in the context of CoT prompting based on black-box models.

In this work, we overcome these shortcomings

---

Dianbo Sui and Chunshan Li are corresponding authors.

by introducing a more foundational approach to assess CoT prompts and conducting a theoretical analysis of CoT prompting in black-box models using this approach.

To achieve a more fundamental comprehension of the effectiveness of CoT, we contemplate introducing  $\mathcal{V}$ -information for its assessment. The  $\mathcal{V}$ -information is widely employed to quantify the amount of accessible information within input  $X$  related to label  $Y$  that is extracted through the utilization of a model family  $\mathcal{V}$  (Xu et al., 2020). Pointwise  $\mathcal{V}$ -information (PVI) further extends the concept of  $\mathcal{V}$ -information from dataset level to instance level (Ethayarajh et al., 2022). Note that, the calculation of  $\mathcal{V}$ -information requires information about all tokens in the vocabulary, thus all prior research on  $\mathcal{V}$ -information has been conducted in the context of a white-box model. In this work, we introduce a metric called EPVI (Estimated Pointwise  $\mathcal{V}$ -Information), for quantifying the new label-relevant information introduced by CoT prompting beyond the pre-existing information in the input on black-box models, like OpenAI API LLMs. Specifically, in cases where the information of the desired token is missing from the black-box model’s output, EPVI employs the Zipf-Mandelbrot distribution to estimate the information associated with that token.

Based on the proposed EPVI, we conduct a theoretical analysis of black-box models at both task-level and instance-level. The task-level analysis assesses the CoT prompts across diverse test samples in the entire dataset, whereas the instance-level analysis explores various CoT prompts on the individual samples. At the task level, we analyze the impact of CoT prompting on the difficulty of various samples and the variation in difficulty across different datasets for the same model  $\mathcal{V}$  during CoT prompting. Moreover, we also investigate what certain types of errors in CoT significantly increase the difficulty of sample inference. At the instance level, we analyze the impact of CoTs of varying quality on the inference difficulty of the same sample.

Various experiments have been conducted on the tasks of commonsense question-answering (CQA) and arithmetic reasoning (AR). The results demonstrate the effectiveness of CoT prompting stems from its ability to influence the model’s reasoning difficulty by increasing or decreasing the amount of label-relevant information that can be extracted by the model. Furthermore, we also demonstrate that selecting high-scoring CoT based on EPVI for composing demonstrations leads to improved performance on downstream tasks.

The primary contributions of this work are:

1. We propose EPVI, a metric for evaluating the

chain of thought on black-box LLMs. To our best knowledge, this is the first metric for evaluating CoT prompting on black-box models.

2. We conducted a comprehensive analysis of CoT prompting at both the task level and the instance level, and unveiled the effectiveness of CoT from an information-theoretic perspective.
3. We further investigate the application of EPVI on CoT prompting and show that EPVI can be utilized to enhance the downstream performance of reasoning tasks.

## 2. Related Work

### 2.1. Large Language Models

LLMs have emerged as a pivotal and versatile component in a range of user-facing language technologies due to their outstanding performance. These models, typically adopt the Transformer (Vaswani et al., 2017) architecture, are trained on extensive corpora and encompass hundreds of billions of parameters. Numerous studies have investigated the performance limits by training increasingly larger language models, such as GPT-3 (175B) (Brown et al., 2020) and PaLM (540B) (Chowdhery et al., 2022). Although scaling primarily involves enlarging the model size while maintaining similar architectures and pre-training tasks, these LLMs exhibit distinct behaviors compared to smaller language models (e.g., 330M BERT and 1.5B GPT-2) (Radford et al., 2019; Devlin et al., 2019) and demonstrate remarkable capabilities often referred to as "emergent capabilities" (Zoph et al., 2022).

### 2.2. Understanding Chain-of-Thought Prompting

With the expansion in the scale of LLMs, standard input-output prompting approaches have failed to yield satisfactory improvements in addressing challenging tasks such as arithmetic, commonsense reasoning, and symbolic reasoning (Rae et al., 2021). To address this drawback, Wei et al. (2022) introduces the Chain-of-Thought prompting, where LLMs are elicited to generate step-by-step intermediate reasoning steps that lead to the final answer to a question. CoT prompting has been demonstrated to dramatically improve the performance of LLMs in various challenging tasks.

Despite significant empirical success, the underlying mechanisms and how CoT unleashes the potential of LLM remain elusive. Several studies have been proposed to examine how and why CoT prompting works in recently. Wang et al.

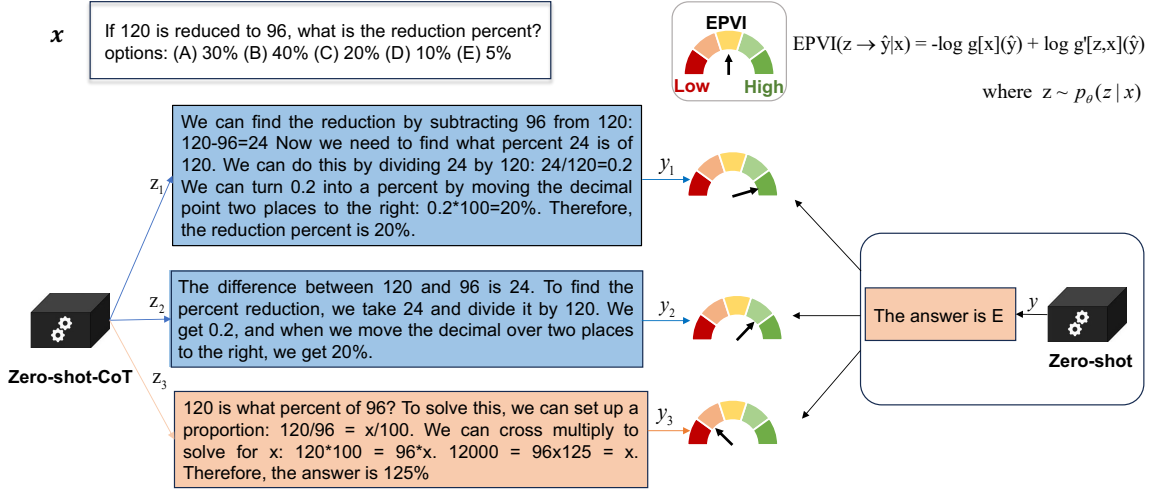


Figure 1: Evaluation framework for different CoTs.  $z_1, z_2, z_3$  represent three CoTs generated by black-box model using the Zero-shot-CoT paradigm and  $y_1, y_2, y_3$  are subsequently generated based on them. Zero-shot paradigm directly produces  $y$  based on  $x$ . Our metric, EPVI, is able to quantify the new information the CoT contains to support the label beyond the input  $x$ .

(2022a); Madaan and Yazdanbakhsh (2022) concentrated on perturbing CoT demonstrations in few-shot prompting, to identify the key factors contributing to the creation of high-quality CoT. Lanham et al. (2023) investigate hypotheses for how CoT reasoning may be unfaithful, by perturbing the model’s stated reasoning in different ways. Wu et al. (2023) examine the impact of CoT prompting on LLMs behavior through an analysis of the saliency scores associated with input tokens, where saliency scores computed by employing gradient-based feature attribution methods on white-box LLMs. Feng et al. (2023) theoretically analyze why CoT prompting is essential in solving mathematical and decision-making problems from a model-capacity perspective.

In addition to these studies, there have been efforts to assess the properties of reasoning chains. Prasad et al. (2023) introduce the RECEVAL framework to evaluate the correctness and informativeness of reasoning chains from both inter and intra perspectives. Chen et al. (2023) propose REV, which measures the new information within the rationale that extends beyond what is available in the input. Both of these studies conduct experiments on white-box models. In this work, we focus on analyzing CoT prompting from an information-theoretic perspective on black-box models. We first propose a new metric for assessing CoT on black-box models and then analyze CoT prompting on black-box models at both task and instance levels based on this approach. Unlike Prasad et al. (2023); Chen et al. (2023), which requires specifying input and output on the white-box model, we directly employ

the model to generate CoT and subsequently analyze CoT prompting based on this.

### 3. EPVI: Information-Theoretic Estimated of CoT Evaluation

We introduce a new metric, EPVI, Estimated Pointwise V-Information. Based on the information theoretic framework of pointwise  $\mathcal{V}$ -information, EPVI evaluates the CoT on black-box LLMs using an estimation method. We first provide a brief introduction of the pointwise  $\mathcal{V}$ -information in the section § 3.1. Next, in the section § 3.2, we employ the Zipf-Mandelbrot distribution to estimate the probabilities of unobserved tokens and introduce how EPVI evaluates the CoT on black-box LLMs using this approach.

#### 3.1. Preliminary

Let  $X$  and  $Y$  denote two random variables with sample spaces  $\mathcal{X}$ ,  $\mathcal{Y}$  respectively. The conditional entropy between  $X$  and  $Y$  is defined as  $H(Y|X) = \mathbb{E}[-\log P(Y|X)]$  (Shannon, 1948). However, this computation necessitates knowledge of the actual joint distribution of  $X$  and  $Y$ , which can be unfeasible in real-world scenarios. As an alternative, Xu et al. (2020) propose conditional  $\mathcal{V}$ -entropy using a model family  $\mathcal{V}$  that learns the mapping from  $X$  to  $Y$ , which is defined as:

$$H_{\mathcal{V}}(Y|X) = \inf_{f \in \mathcal{V}} \mathbb{E}[-\log f[X](Y)] \quad (1)$$

where  $f[X]$  yield a probability distribution across

---

**Algorithm 1** Evaluating CoT on Black-Box Models Using EPVI.

**Input:** held-out data  $\mathcal{D}_{\text{test}} = (\text{input } x_i, \text{label } y_i)_{i=1}^n$ , Model  $\mathcal{V}$

**do**

$g$  and  $g' \leftarrow$  auto-regressive black-box PLMs

$z \leftarrow$  chain-of-thought

$H_{\mathcal{V}}(Y|X, Z), H_{\mathcal{V}}(Y|X) \leftarrow 0, 0$

**for**  $(x_i, y_i) \in \mathcal{D}_{\text{test}}$  **do**

$z_i \leftarrow z_i \sim p_{\theta}(z_i | x, z_{1 \dots i-1})$

$y_i \sim p_{\theta}^{\text{CoT}}(y_i | x_i, z_i)$

$\hat{y}_i \leftarrow$  estimate based on Eq.5 if it cannot obtain

$H_{\mathcal{V}}(Y|X, Z) \leftarrow H_{\mathcal{V}}(Y|X, Z) - \frac{1}{n} \log g'[z_i, x_i](\hat{y}_i)$

$H_{\mathcal{V}}(Y|X) \leftarrow H_{\mathcal{V}}(Y|X) - \frac{1}{n} \log g[x_i](\hat{y}_i)$

$\text{EPVI}(z_i \rightarrow \hat{y}_i | x_i) \leftarrow -\log g[x_i](\hat{y}_i) + \log g'[z_i, x_i](\hat{y}_i)$

**end for**

$I_{\mathcal{V}}(X \rightarrow \hat{Y}|Z) \leftarrow \frac{1}{|n|} \sum_i \text{EPVI}(x_i, y_i, z_i)$

---

the labels. Let  $\emptyset$  denote an empty input that imparts no information regarding  $Y$ . In a unique circumstance, when  $X = \emptyset$ , the conditional  $\mathcal{V}$ -information is defined as:

$$H_{\mathcal{V}}(Y|\emptyset) = \inf_{f \in \mathcal{V}} \mathbb{E}[-\log f[\emptyset](Y)] \quad (2)$$

The goal of the models for  $f \in \mathcal{V}$  is to maximize the log-likelihood of the label data with and without the input. Based on the content mentioned above, Xu et al. (2020) propose a framework called  $\mathcal{V}$ -information. It generalizes Shannon information to quantify the extractable information from  $X$  about  $Y$  under the model family  $\mathcal{V}$ , written as  $I_{\mathcal{V}}(X \rightarrow Y)$ . It is defined as:

$$I_{\mathcal{V}}(X \rightarrow Y) = H_{\mathcal{V}}(Y|\emptyset) - H_{\mathcal{V}}(Y|X) \quad (3)$$

In order to extend  $\mathcal{V}$ -information framework from dataset level to instance level, Ethayarajh et al. (2022) propose pointwise  $\mathcal{V}$ -information (PVI) to quantify the extent of the information contained within individual data points  $(x, y)$  as:

$$\text{PVI}(x \rightarrow y) = -\log g[\emptyset](y) + \log g'[x](y) \quad (4)$$

where  $g$  and  $g'$  are the models fine-tuned on the same dataset with null-target pairs  $(\emptyset, y_i)$  and input-target pairs  $(x_i, y_i)$ , respectively.

### 3.2. Methodology: Informational Evaluation of CoT using EPVI

Our aim is to employ the PVI on black-box models for evaluating CoT from an information-theoretic perspective. However, for black-box models, calculating PVI becomes impossible when missing the probabilities for the tokens corresponding to

label  $y$ . To address this issue, we utilize the Zipf-Mandelbrot distribution to estimate the probabilities of the desired tokens. Building upon this method, we propose a new metric, EPVI, enabling information assessment of CoT on black-box models.

#### Estimation of Token Probabilities with the Zipf-Mandelbrot Distribution in Black-Box Models.

For each position in the black-box model's output, we typically have access to probabilities for only a limited set of alternative tokens (i.e., an array of probability objects, representing tokens most likely to be used for the completion), rather than the probabilities for all tokens in the entire vocabulary. Similar to natural language text or corpora, where the frequency of the most common words significantly outweighs that of other words, the probability of top- $k$  alternative tokens returned by a black-box model at each position is much higher than the probability of the remaining tokens. Therefore, if the probability of the desired token is not among the top- $k$  tokens returned by the black box model, we can assume that its probability can be approximated by estimating the probability of the  $(k+1)$ th token. Inspired by the applicability of Zipf's law to word frequency tables in natural language text or corpora, we introduce the Zipf-Mandelbrot distribution<sup>1</sup> to estimate the probability of the desired token. Zipf-Mandelbrot distribution is a discrete probability distribution that represents a power-law distribution of ranked data. Its probability mass function is defined as:

$$f(k; N, q, s) = \frac{1/(k+q)^s}{H_{N,q,s}} \quad (5)$$

where  $k$  is the rank of the element,  $N$  is the number of elements, and  $q, s$  are parameters of the distribution. The  $H_{N,q,s}$  can be considered as an extension of the concept of a harmonic number, which is defined as:

$$H_{N,q,s} = \sum_i^N \frac{1}{(i+q)^s} \quad (6)$$

We approximate the probability of missing the desired token by employing the Zipf-Mandelbrot distribution to predict the probabilities of the  $(k+1)$ th token at its position. Specifically, we estimate the values of  $q$  and  $s$  in Eq.5 using the maximum likelihood estimation method based on the probabilities of the top- $k$  tokens. Here,  $N$  represents the size of the model's vocabulary. Building upon the above, we employ Eq.5 to calculate the probability of the  $(k+1)$ th token.

<sup>1</sup>[https://en.wikipedia.org/wiki/Zipf-Mandelbrot\\_law](https://en.wikipedia.org/wiki/Zipf-Mandelbrot_law)

**Computing EPVI for CoT Evaluation on Black-Box Models.** In order to evaluate the effect of chain-of-thought from an information-theoretic perspective, we employ the pointwise  $\mathcal{V}$ -information for evaluating CoT within individual samples, denoted as:

$$PVI(z \rightarrow y|x) = -\log g[x](y) + \log g'[z, x](y) \quad (7)$$

The models  $g$  and  $g'$  we used are the same autoregressive pre-trained language models, employed to compute  $H_{\mathcal{V}}(Y|X, Z)$  and  $H_{\mathcal{V}}(Y|X)$ , respectively. Specifically, the model  $g$  and  $g'$  generate label  $y$  under different scenarios where the input either contains or does not contain CoT  $z$ . Building upon the aforementioned content, we propose a new metric EPVI, Estimated Pointwise V-Information, to evaluate the CoT on black-box LLMs, which is defined as:

$$EPVI(z \rightarrow \hat{y}|x) = -\log g[x](\hat{y}) + \log g'[z, x](\hat{y}) \quad (8)$$

The  $\mathcal{V}$ -information of the test data  $\mathcal{D}_{test}$  on black-box models can be computed using the average EPVI:

$$I_{\mathcal{V}}(X \rightarrow \hat{Y}|Z) = \frac{1}{|\mathcal{D}_{test}|} \sum_i EPVI(x_i, y_i, z_i) \quad (9)$$

where  $\hat{y}$  and  $\hat{Y}$  represent the estimated probabilities of the label at the instance and dataset levels, respectively. We compute EPVI on an individual sample  $(x, z, y)$ . As illustrated in Figure 1,  $z$  is generated based on  $x$ , while  $y$  is generated with or without  $z$  in the input. When the probability of tokens related to label  $y$  can be obtained from the output, our calculation of EPVI is similar to the calculation of PVI. In instances where such probabilities are unavailable, we calculate EPVI by estimating the probability of  $\hat{y}$  with the Zipf-Mandelbrot distribution. Algorithm 1 shows our computation of EPVI and  $\mathcal{V}$ -information on black-box models.

For any  $f \in \mathcal{V}$ , a positive EPVI signifies that CoT provides more information to support the label. The larger the EPVI, the easier the sample is for  $\mathcal{V}$ . A negative EPVI indicates that CoT contains additional information that does not support the label. Figure 1 illustrates this phenomenon via an example. EPVI can assign a positive score for an incorrect prediction if the CoT can supply more useful information to support the label  $y$  beyond the input. Therefore, unlike accuracy, EPVI can provide a better explanation of the effectiveness of CoT.

## 4. Is EPVI a good estimator for PVI?

Our goal is to validate whether EPVI is a good estimator for PVI. To this end, we calculate EPVI and

PVI for various datasets on a series of white-box LLMs and employ a null-hypothesis significance test to determine if there is a statistically significant difference between EPVI and PVI. It is worth noting that we simulate the black-box environment and calculate EPVI under such circumstances.

### 4.1. Experimental Setup

**Datasets.** We consider the following two math word problem benchmarks: (1) the **GSM8k** benchmark of math word problems (Cobbe et al., 2021), (2) the **AQuA**(Ling et al., 2017) dataset of algebraic word problems. Besides, we explore a multimodal commonsense reasoning dataset: **SCIENCEQA**(Lu et al., 2022). For SCIENCEQA, in accordance with the configuration of Lu et al. (2022), we use the question, options, and context text as the input, with the context also encompassing the caption extracted from the image.

**Metrics.** We employed the Paired Samples T-Test to assess the consistency of the distribution of EPVI and PVI. The Paired Samples T-Test is employed to assess the statistical difference in means between two sample groups under the null hypothesis, where each individual observation in one sample can be paired with an observation in the other sample. The null hypothesis is the claim that no relationship exists between two sets of data being analyzed. In the Paired Samples T-Test, it indicates that there is no significant mean difference between the two related samples. Paired Samples T-Test returns two values, statistic and p-value. The p-value represents the probability of obtaining test results at least as extreme as the result actually observed, assuming that the null hypothesis is correct. When the p-value is lower than 0.05, it indicates that there is statistical significance in the differences between the two computational methods. We also calculate the effect size of the Paired Samples T-Test. The effect size in the Paired Samples T-Test indicates the strength of the difference between the groups. When the effect size is less than 0.2, it means that the difference between the means of the two groups does not exceed 0.2 standard deviations, indicating a negligible difference.

**Implementation.** To simulate the situation where only a limited set of token probabilities is available in a black-box model, we constrain the white-box models to return only the top five tokens with the highest probabilities at each position. We then compute EPVI under such circumstances, which aligns with the return behavior of models from the GPT series black-box models. For PVI, we supply the probabilities of every token in the vocabulary

Model	Parm	GSM8k	AQuA	SCIENCEQA
LLaMA-1-7B	Statistic	0.760	0.301	1.113
	p-value	0.447	0.763	0.265
	Effect size	0.010	0.021	0.046
LLaMA-1-13B	Statistic	1.781	1.625	-0.866
	p-value	0.075	0.105	0.386
	Effect size	0.027	0.090	0.001
LLaMA-1-30B	Statistic	1.914	-0.696	1.620
	p-value	0.056	0.487	0.105
	Effect size	0.022	-0.047	0.003

Table 1: Paired Samples T-Test on three datasets across three LLMs of varying sizes.

at each position. We utilize the LLMs to generate answers under two paradigms: Zero-shot-CoT and Zero-shot (e.g., the example shown in Figure 1). We calculate PVI and EPVI using Eq.7 and Eq.8, respectively, and conduct experiments on LLaMA-1 (Touvron et al., 2023) at three different scales: LLaMA-1-7B, LLaMA-1-13B, LLaMA-1-30B.

## 4.2. Results

We perform the Paired Samples T-Test using the corresponding EPVI and PVI values for each group. Table 1 presents the corresponding results on three datasets from two categories of reasoning tasks. The results indicate that, in all cases, the p-value is greater than 0.05, suggesting that there is no statistically significant difference between PVI and EPVI. Furthermore, in all cases, the effect size smaller than 0.2 implies that the difference between PVI and EPVI is negligible. Therefore, in scenarios where a limited token probability distribution set is employed, such as in black-box models, EPVI is a good estimator for PVI.

## 5. Analyzing the effect of CoT on black-box models

Due to the inability of black-box models to be trained and obtain probabilities for every token in the vocabulary, there has been almost no research on utilizing black-box models for theoretical analysis of CoT prompting. To further advance research on CoT prompting based on black-box models and conduct a more fundamental theoretical analysis of why CoT prompting is effective, we employ EPVI to analyze the impact of CoT prompting at both the task level and the instance level on black-box models.

### 5.1. Experimental Setup

**Datasets.** The CoT prompting is most helpful when

Approach	GSM8k	AQuA	MultiArith
Highest-EPVI	<b>69.30</b>	<b>63.78</b>	<b>98.50</b>
Greedy	45.35	33.50	78.00
Random	38.52	32.87	74.50

Table 2: Experimental results of different CoT sampling approaches on GSM8k, AQuA and MultiArith datasets.

applied to a challenging, multi-step reasoning task involving a large language model (Wei et al., 2022). Therefore, for a clearer analysis of CoT prompting, we instantiate this analysis using data from three arithmetic reasoning benchmarks: GSM8k, AQuA and MultiArith (Roy and Roth, 2015).

**Implementation.** For LLM, we conduct experiments with text-davinci-002 model (175B) from GPT-3 (Brown et al., 2020), where this model has strong CoT reasoning performance as reported in Wang et al. (2022b); Kojima et al. (2022); Zhang et al. (2022). We employ LLMs to generate answers within two paradigms: Zero-shot-CoT and Zero-shot. We then calculate EPVI using Eq.8.

For the instance level, we generate multiple chain-of-thoughts for the same sample. We specify the number of completions to 5 and apply temperature sampling with  $T=0.7$  without top-k truncation as Wang et al. (2022b). For the GSM8k and MultiArith datasets, we conduct experiments using randomly selecting four hundred and two hundred samples from their test set. To explore the influence of different CoT on the same sample, we randomly select three groups of CoT from the generated CoT sets for the whole samples, denoted as Random. Besides, we further investigate whether improving the  $\mathcal{V}$ -information of CoT (i.e., the EPVI of each CoT) is advantageous for improving the accuracy of model-generated answers. We select the CoT with the highest EPVI from the generated CoT sets for each sample, denoted as Highest-EPVI.

For the task level, we use greedy decoding when generating outputs. We investigate the impact of CoT prompting on different samples in the dataset by calculating the EPVI for each sample. Based on this, we analyze the variation in difficulty across different datasets for the same model  $\mathcal{V}$ . Moreover, we select the top five CoTs that exhibited the most significant increase in sample complexity when employing CoT prompting with text-davinci-002 on the GSM8k dataset. We analyze the types of errors occurring in these CoTs, with the goal of determining what certain types of errors in CoT significantly increase the difficulty of sample inference. We categorize the error types to be the same as Wei et al.

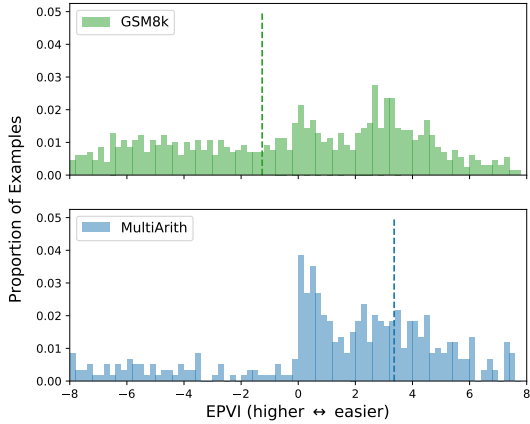


Figure 2: The MultiArith dataset contains more GPT-usable information than GSM8k dataset, making it easier for text-davinci-002. Above, the distribution of instance difficulty (EPVI) is illustrated for the test sets of each; the dashed lines represent the mean EPVI. Note that the above figure is not a complete illustration; it only represents samples of EPVI between -8 and 8.

(2022), who identified three correctable error types associated with CoT prompting: calculator error only, symbol mapping error, and one-step missing error, while the rest are difficult to amend.

## 5.2. Results

**Instance-Level CoT Analysis with EPVI.** For Random, we report the average results over three random orders. The 'Greedy' approach involves chain-of-thought prompting with a greedy decoding method. As shown in Table 2, a comparison of the results across different sampling methods reveals that Highest-EPVI can achieve much better performance in all datasets. For example, the Highest-EPVI improves the accuracy by 30.78%, 30.91%, and 24.50% with respect to GSM8k, AQuA, and MultiArith over Random, respectively. This indicates that, for the same sample, high-quality CoT prompt (Higher EPVI) improves model inference accuracy by introducing more  $\mathcal{V}$ -information (i.e., more label-relevant information that can be extracted by the model).

**Task-Level CoT Analysis with EPVI.** We employ EPVI to analyze the impact of CoT on each sample within the dataset. Figure 2 illustrates that a higher EPVI indicates that the sample is easier for  $\mathcal{V}$ , and it also demonstrates that different datasets offer varying amounts of usable information for the same model  $\mathcal{V}$ . This suggests that, for different samples within the dataset, CoT prompts influence the difficulty of model inference by either increas-

ing or decreasing the label-related information that can be extracted from model  $\mathcal{V}$ . Table 3 shows the top five CoTs that exhibited the most significant increase in sample complexity when employing CoT prompting with text-davinci-002 on the GSM8k dataset. By analyzing the error types of each CoT, we can find that all errors resulting in a significant increase in sample difficulty in CoT prompting are typically correctable. This indicates that when the model utilizes CoT prompting for reasoning, we can further reduce the model inference difficulty on this dataset by correcting the CoT prompts.

## 6. Downstream Utility of EPVI Metrics

While the evaluation of CoT is a demanding yet essential endeavor, we also investigate the feasibility of leveraging our proposed metrics to enhance the performance of downstream tasks. Specifically, we study whether the demonstration examples with higher-quality CoT (higher EPVI) contribute to the improvement of ICL (In-Context Learning) capabilities in LLM.

### 6.1. Experimental Setup

We conduct experiments with text-davinci-002 on the GSM8k, AQuA, and MultiArith datasets. We employ few-shot CoT prompting, and the number of demonstration examples for different datasets is the same as Wei et al. (2022). To analyze whether the demonstrations with high-quality CoT are helpful for ICL, we need to select a group of demonstrations for the test set. We randomly select 500 samples from the training set of GSM8k and AQuA and employ the Zero-shot-CoT paradigm to generate CoT for each of them. We then use these to compose a demonstration set, denoted as Train-500. For the MultiArith dataset, we aligned the selection of its demonstrations with that of the GSM8k dataset. To explore the influence of different quality demonstrations of CoT reasoning on the performance of ICL, we randomly select three groups of demonstrations from the Train-500, denoted as Random. We also select the sample with the highest-quality CoT (highest EPVI) in Train-500 to form the demonstrations, denoted as Top-EPVI. Besides, we further investigate whether improving the diversity of demonstrations is beneficial to ICL. We select the demonstrations by first dividing the whole test set into N clusters and then randomly selecting one demonstration from each cluster, namely KmeansRND. Similar to Random, we conduct three selection processes for KmeansRND, resulting in three sets of demonstration examples. Moreover, We select the highest-quality

**QUESTION:** Regina wrote 9 novels last year. If this is 3 quarters of the number of novels she has written this year, how many novels has she written this year?

**MODEL ANSWER (INCORRECT; SYMBOL MAPPING ERROR):** Regina wrote 9 novels last year. That means she wrote 3/4 as many novels this year as she did last year. That means she wrote  $3/4 \cdot 9$  novels this year. That means she wrote 6.75 novels this year. Therefore, Regina wrote 6.75 novels this year. The answer is 6.75.

**QUESTION:** Prince is a butcher. He sells 15kg of meat every hour he works, and he works 10 hours a day. His friend Bill gives him a bull that weighs 750kg. How many days will it take Prince to sell the meat from Bill’s bull?

**MODEL ANSWER (INCORRECT; CALCULATOR ERROR ONLY):** How many kgs of meat does Prince sell in a day? How many kgs of meat does Prince sell in an hour? How many hours will it take Prince to sell the meat from Bill’s bull? How many days will it take Prince to sell the meat from Bill’s bull? 1. Prince sells 15 kgs of meat every hour. 2. Prince sells 150 kgs of meat every day. 3. It will take Prince 50 hours to sell the meat from Bill’s bull. 4. It will take Prince **2.5 days**. The answer is 2.5.

**QUESTION:** Poppy is solving a 1000-piece jigsaw puzzle. She places a quarter of the pieces on the board, then her mom places a third of the remaining pieces. How many jigsaw pieces are left to be placed?

**MODEL ANSWER (INCORRECT; CALCULATOR ERROR ONLY):** Poppy places a quarter of the pieces on the board. That means she has used up a quarter of the pieces, so she has 3/4 of the pieces left. Her mom then places a third of the remaining pieces. That means her mom has used up a third of the pieces that were left, so she has used up of the total pieces. That means there are **11/12** of the pieces left to be placed. The answer is 11/12 or 990 pieces.

**QUESTION:** Sally received the following scores on her math quizzes: 50, 80, 80, 60, 40, 90, 100, 70, and 60. Find her mean score.

**MODEL ANSWER (INCORRECT; CALCULATOR ERROR ONLY):** The first step is to add up all of the numbers that Sally got on her math quizzes. This is her total score.  $50 + 80 + 80 + 60 + 40 + 90 + 100 + 70 + 60 = 660$  The second step is to count how many quiz scores there are in total. There are 9 quiz scores in total. The third step is to divide the total score by the number of quiz scores. This will give you the mean, or average, score.  $660/9 = 73.3$  Therefore, The answer is 73.3.

**QUESTION:** John goes to the market and buys 3 goats for \$500 each and 2 cows for \$1500 each. How much money did he spend?

**MODEL ANSWER (INCORRECT; ONE STEP MISSING ERROR):** John goes to the market. He buys 3 goats for \$500 each. That’s  $\$500 + \$500 + \$500 = \$1,500$  so far. He buys 2 cows for \$1,500 each. **That’s  $\$1,500 + \$1,500 = \$3,000$**  in total. The answer is 3000.

Table 3: The five hardest (lowest EPVI) samples for CoT prompting on the GSM8k dataset by the text-davinci-002 model.

CoT (highest EPVI) from each cluster, namely Top-Kmeans. We use the text-embedding-ada-002 (Neelakantan et al., 2022) for vectorization and use the K-means++ algorithm (Arthur and Vassilvitskii, 2007) for clustering, where K is set to N, representing the number of demonstrations.

## 6.2. Results

We present the results in Table 4. For Random and KmeansRND, we report the average results over three random orders. We can find that the results of KmeansRND outperform Random on all three datasets, suggesting that enhancing the diversity of selected demonstration examples is advantageous for constructing task-level demonstrations. We observe that in comparison to Random, Top-EPVI exhibits an improvement in accuracy of 9.30%, 1.44%, and 14.06% on GSM8k, AQuA, and MultiArith, respectively. Besides, by comparing the results between Top-Kmeans and KmeansRND, we can find that Top-Kmeans yields an improvement of 1.90%, 5.46%, 9.12% on GSM8k, AQuA,

Approach	GSM8k	AQuA	MultiArith
Random	36.11	28.99	67.27
Top-EPVI	<b>45.41</b> ( $\Delta+9.30$ )	<b>30.43</b> ( $\Delta+1.44$ )	<b>81.33</b> ( $\Delta+14.06$ )
KmeansRND	45.56	23.23	76.05
Top-Kmeans	<b>47.46</b> ( $\Delta+1.90$ )	<b>28.69</b> ( $\Delta+5.46$ )	<b>85.17</b> ( $\Delta+9.12$ )

Table 4: Experimental results of different demonstration selection approaches on each task.

and MultiArith datasets, respectively, compared to KmeansRND. These results indicate that EPVI is a promising CoT evaluation metric for enhancing the performance of downstream tasks. Future work can explore methods that combine these metrics with other strategies to improve LLM’s reasoning abilities.

## 7. Conclusion

In this paper, we propose an information-theoretic metric, EPVI, to evaluate the CoT on black-box



models. EPVI quantifies whether CoT contains new label-relevant information beyond what is present in the input. We demonstrate that EPVI is a good estimator of PVI. Based on this, we conduct an analysis of CoT prompting at both the instance and task levels, demonstrating that the effectiveness of CoT prompting lies in its capacity to influence the difficulty of model reasoning by increasing or decreasing the label-relevant information that can be extracted by the model. Furthermore, we show that high-quality demonstrations of CoT reasoning selected based on EPVI can enhance downstream performance in reasoning tasks. Future work might entail exploring the assessment of different CoT prompting paradigms, such as few-shot prompting.

## 8. Acknowledgements

This work is supported by the National Key R&D Program of China (Grant No. 2023YFB3307500). This work is also supported by the National Natural Science Foundation of China (Grant No. 62306087), the Natural Science Foundation of Shandong Province (Grant No. ZR2023QF154), Special Funding Program of Shandong Taisihan Scholars Project and Major Scientific and Technological Innovation Project of Shandong Province (Grant No. 2021ZLGX05 and Grant No. 2020CXGC010705). Besides, we sincerely thank the anonymous reviewers for their valuable feedback.

## Bibliographical References

- David Arthur and Sergei Vassilvitskii. 2007. K-means++: The advantages of careful seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '07, page 1027–1035, USA. Society for Industrial and Applied Mathematics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Hanjie Chen, Faeze Brahman, Xiang Ren, Yangfeng Ji, Yejin Choi, and Swabha Swayamdipta. 2023. [REV: Information-theoretic evaluation of free-text rationales](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2007–2030, Toronto, Canada. Association for Computational Linguistics.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#).
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. 2022. Understanding dataset difficulty with v-usable information. In *International Conference on Machine Learning*, pages 5988–6008. PMLR.
- Guhao Feng, Yuntian Gu, Bohang Zhang, Haotian Ye, Di He, and Liwei Wang. 2023. Towards revealing the mystery behind chain of

- thought: a theoretical perspective. *arXiv preprint arXiv:2305.15408*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. 2023. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. [Program induction by rationale generation: Learning to solve and explain algebraic word problems](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. [Learn to explain: Multimodal reasoning via thought chains for science question answering](#).
- Aman Madaan and Amir Yazdanbakhsh. 2022. [Text and patterns: For effective chain of thought, it takes two to tango](#).
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. [Text and code embeddings by contrastive pre-training](#).
- OpenAI. 2023. [Gpt-4 technical report](#).
- Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and Mohit Bansal. 2023. Reveal: Evaluating reasoning chains via correctness and informativeness. *arXiv preprint arXiv:2304.10703*.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- Subhro Roy and Dan Roth. 2015. [Solving general arithmetic word problems](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.
- C. E. Shannon. 1948. [A mathematical theory of communication](#). *The Bell System Technical Journal*, 27(3):379–423.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#).
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2022a. Towards understanding chain-of-thought prompting: An empirical study of what matters. *arXiv preprint arXiv:2212.10001*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022b. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Skyler Wu, Eric Meng Shen, Charumathi Badrinath, Jiaqi Ma, and Himabindu Lakkaraju. 2023. [Analyzing chain-of-thought prompting in large language models via gradient-based feature attributions](#). *CoRR*, abs/2307.13339.
- Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. 2020. A theory

of usable information under computational constraints. *arXiv preprint arXiv:2002.10689*.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*.

Barret Zoph, Colin Raffel, Dale Schuurmans, Dani Yogatama, Denny Zhou, Don Metzler, Ed H. Chi, Jason Wei, Jeff Dean, Liam B. Fedus, Maarten Paul Bosma, Oriol Vinyals, Percy Liang, Sebastian Borgeaud, Tatsunori B. Hashimoto, and Yi Tay. 2022. Emergent abilities of large language models. *TMLR*.