Reasoning over Uncertain Text by Generative Large Language Models

Aliakbar Nafar¹, Kristen Brent Venable^{2,3}, Parisa Kordjamshidi¹

¹Michigan State University ²Florida Institute for Human and Machine Cognition ³University of West Florida {nafarali, kordjams}@msu.edu, bvenable@ihmc.org

Abstract

This paper considers the challenges Large Language Models (LLMs) face when reasoning over text that includes information involving uncertainty explicitly quantified via probability values. This type of reasoning is relevant to a variety of contexts ranging from everyday conversations to medical decision-making. Despite improvements in the mathematical reasoning capabilities of LLMs, they still exhibit significant difficulties when it comes to probabilistic reasoning. To deal with this problem, we introduce the Bayesian Linguistic Inference Dataset (BLInD), a new dataset specifically designed to test the probabilistic reasoning capabilities of LLMs. We use BLInD to find out the limitations of LLMs for tasks involving probabilistic reasoning. In addition, we present several prompting strategies that map the problem to different formal representations, including Python code, probabilistic algorithms, and probabilistic logical programming. We conclude by providing an evaluation of our methods on BLInD and an adaptation of a causal reasoning question-answering dataset. Our empirical results highlight the effectiveness of our proposed strategies for multiple LLMs.

Code and Dataset — https://github.com/HLR/BLInD Extended Version — https://arxiv.org/abs/2402.09614

Introduction

Uncertainty in text is communicated in many contexts, ranging from everyday conversations to domain-specific documents, such as those with medical focus (Heritage 2013; Landmark, Gulbrandsen, and Svennevig 2015). Processing this uncertain information is critical. For example, uncertainty in text has been shown to significantly affect decisionmaking in the biomedical domain (Poggi et al. 2019). Reasoning over uncertain text is also closely related to rational reasoning, e.g., if the probabilities of events A and B are low, the probability of both happening simultaneously should also be low. As a result, it is essential for language models to be able to use text with uncertainty and perform inference based on it. While the human intuitive approach to probabilistic reasoning often aligns with Bayesian Rationalism (Oaksford and Chater 2007, 2009), humans usually do not explicitly calculate the probabilities of outcomes.



Figure 1: An example from the BLInD dataset including an underlying Bayesian network, its textual description, and probabilistic queries in natural language form.

Still, probabilistic modeling using Bayesian Networks offers a robust computational approach for dealing with uncertainty. Thus, we tackle the challenge of enabling LLMs to conduct probabilistic reasoning by mapping uncertainty expressed through language to a Bayesian Network. This approach resembles other strategies for enabling mathematical reasoning over text, such as with math word problems (MWPs) (Cobbe et al. 2021; Kim et al. 2023). The common theme of our problem formulation and MWPs is that a formal problem is extracted from the text and solved using external tools (Dries et al. 2017; He-Yueya et al. 2023).

First-generation LLMs were shown to struggle with mathematical reasoning and fail in answering even simple questions (e.g., about summation (Mishra et al. 2022)). With the advent of newer generations of LLMs, their mathematical reasoning capability improved significantly, with GPT4 achieving 92% (OpenAI 2023) and Gemini achieving 94.4% (Google 2023) accuracy on the Grade School Math (GSM8K) dataset (Cobbe et al. 2021). These results have

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

misled to the belief that LLMs are now proficient in mathematical reasoning. However, the LLMs' performance on math problems varies significantly depending on the question types (Kim et al. 2023). Here, we confirm this latter result by showing that LLMs still struggle with the essential task of probabilistic reasoning over text. Furthermore, we illustrate how, depending on the LLM, different limitations and weaknesses hinder their ability to solve these problems. Simply utilizing Chain of Thought (Wei et al. 2022b) or Python code is not always effective (Shi, Zhang, and Lipani 2022; Kim et al. 2023). These observations support the design of customized solutions for each problem and model.

We focus on Bayesian inference over text and introduce a new dataset, Bayesian Linguistic Inference Dataset (BLInD), to evaluate and improve LLMs' probabilistic reasoning. BLInD instances have up to 10 interconnected random variables used to answer a probabilistic query over them. Figure 1 shows a BLInD example, with the Bayesian Network and conditional probability tables at the top. The corresponding natural language explanation, which is given in input to the language models, is shown below. Given the textual context, the models are asked to answer probabilistic queries such as "What is the probability of event P?".

We design prompts that decompose this complex problem into extracting the probability values and generating the dependency graph prior to probabilistic inference. We investigate solutions that include the above extractions as well as a mapping to symbolic solvers such as pure Python code, probabilistic inference algorithms, and probabilistic logical formalisms. Ultimately, we test our methods on our new challenging dataset and on an adapted version of a causal reasoning question-answering dataset, CLADDER (Jin et al. 2023), further solidifying our results.

In summary, our contributions are as follows: 1) Creating a **new dataset (BLInD)** designed for reasoning over text with uncertainty explicitly quantified as probabilities; 2) **Analyzing the capabilities** of LLMs in solving the complex probabilistic reasoning problems contained in BLInD, highlighting their limitations; 3) Designing **innovative prompt engineering and in-context learning** techniques which leverage mapping to Python code, to inference algorithms, and to a probabilistic logical formalism, leading to improvements in performance across multiple LLMs, proprietary and open-source; 4) **Evaluating** the proposed techniques on our new dataset and an adapted existing benchmark.

Related Work

A few prior works have explored question-answering (QA) involving probabilistic rules. RuleBERT (Saeed et al. 2021) and (Nafar, Venable, and Kordjamshidi 2024) mainly evaluate BERT-based models (Devlin et al. 2019) by fine-tuning them. They use a simple independence structure instead of dealing with arbitrary Bayesian Networks. Their queries are limited to asking the probability of a single variable, and their closed world assumption (CWA) assigns a probability of zero to any event with unspecified probability. In contrast, our queries involve any joint and marginal computation and do not use the CWA. CLADDER (Jin et al. 2023) creates a dataset with probabilistic contexts, but it is mainly designed

to test the causal reasoning capabilities of LLMs with incontext learning and structured prompts. They use a limited number of variables (less than 5), their task setting is limited to binary QA, and their solution is to map the natural language text to a causal reasoning formalism. An adaptation of this dataset for mapping to probabilistic reasoning applies to our problem setting and is used in our experiments.

Looking at reasoning over uncertain text as a form of the math word problem, (Dries et al. 2017; Suster et al. 2021) solve simple probability word problems from introductory mathematics textbooks. However, in most questions, the probabilities are not directly given in the context, and the inference does not necessarily require mapping to Bayesian Networks. These works utilize either classical NLP parsers or fine-tuned LMs instead of our in-context prompting methods. NumGLUE (Mishra et al. 2022) is the first work that analyses Pre-trained and Large Language Models for mathematical reasoning. But, it is limited to questions that require simple arithmetic reasoning. (Bubeck et al. 2023; Frieder et al. 2023; Kim et al. 2023) looks at a broader range of math questions for analyzing LLM's reasoning capabilities. However, none of these works include Bayesian probabilistic questions with complex structures.

In our solutions, we use neuro-symbolic methods to reason over uncertain text. Neuro-symbolic techniques have been used in related research to solve various NLP tasks by integration of symbolic reasoning during training or inference (Rajaby Faghihi et al. 2021, 2023) though not for probabilistic reasoning over uncertain text. In a slightly related work, ThinkSum (Ozturkler et al. 2023) uses probabilistic reasoning by calculating the likelihood of the LLM generating each possible answer and then aggregating these token probabilities. This approach is applied to usual question answering problems that output the final crisp answers. This is fundamentally different from our work that interprets the uncertainty measures that are expressed explicitly in the text and reasons over them to infer the probability of a query.

Problem Definition

The input to the QA task is a textual *context* paired with a probabilistic *query* in a textual form which we refer to as the *question* throughout the paper. The context comprises sentences that describe the probability of random events, which are binary variables, or the conditional probabilities of events. Figure 1 shows a context with five sentences describing the probabilities of random events G, P, and O. The query can be any question that probes the probabilities of these events, such as "What is the probability of G being true and P being false given that O is false?". The output, which is the probability of the query, is a real number that ranges from 0.0 to 1.0.

BLInD Generation

We introduce the Bayesian Language Inference Dataset (BLInD) to investigate the ability of LLMs to perform probabilistic reasoning. Each example in the dataset contains a textual context describing the probability of events and a textual question about the probability of events occurring in the context. Moreover, we provide a Bayesian Network (BN) corresponding to the context with their conditional probability tables (CPTs) and a probability value computed as the answer to the question. In this section, we provide an overview of the generation process and the dataset structure. Details are included in the Appendix of the *arXiv* version of the paper (the link is provided below the abstract).

Bayesian Network

In the first step of our dataset creation, we generate all isomorphic graphs that would serve as our Bayesian Networks, each including up to ten random variables. We generate these graphs with the following properties: 1) all graphs are Directed Acyclic (DAGs), 2) all graphs are Weakly Connected, and 3) each node has at most one parent (arbores*cence*). This results in dataset splits denoted as V_i for $i \in$ $\{2, 3, \ldots, 10\}$, each including a set of graphs with *i* nodes (random variables). Properties 2 and 3 are necessary to control the complexity of the splits. The complexity increases as *i* increases. To clarify, property 2 prevents the breakdown of a graph into smaller, independent, and subsequently simpler components. Assumption 3 results in 2 + (V - 1) * 4probability entries in a BN's CPTs with V variables (2 probabilities for the root node and 4 for other nodes). For example, in Figure 1, we depict a BN over nodes G (Green), O (Orange), and P (Pink) with corresponding CPTs and a total of 10 probability entries. While property 3 might restrict the networks' coverage, it enables us to analyze the examples with better control over their complexity. Further, we assume each random variable is binary and fill their associated conditional probability tables with uniformly random generated probabilities ranging from 0.01 to 0.99.

Query

We generate only *complex* queries for a given Bayesian Network. By *complex*, we mean those which require all variables in the BN for inference. For example, for a size 2 BN with variables A and B where A is the parent of B, all possible queries are P(A), P(B), P(A, B), P(A|B), and P(B|A). Among these, P(A) is the only query that is not *complex* since it can be answered only with the CPT of A and, therefore, is not selected. We assign true/false values randomly to the query variables.

Textual Context and Question

After generating the BNs, CPTs, and queries, we create the textual context and question that describes the BN (mapping every entry in the CPTs to natural language) and the query, respectively. For the context, sentences follow two templates: 1) For explaining prior probabilities, we use the template "{node name} is True/False with ##% Probability". 2) For explaining dependent nodes, we use the template "{node name} is True/False with ##% Probability". 2) For explaining dependent nodes, we use the template "{node name} is True/False with ##% Probability, if {parent node name} is True/False". Figure 1 shows the context for a given BN and CPTs. The template for textual questions is: "What is the probability that {a node name} is True/False and ... given that {a node name} is True/False and ...?". For a query without evidence variables, the text after "given" is omitted. For example, $P(A| \sim B)$ would be translated to the textual query (question), "What is the probability that A event is True given that B event is False?".

Verification and Inference

At the final step of our dataset generation, we use the Python library pgmpy (Ankan and Panda 2015) to verify the soundness of our BN, probabilities, and queries and to infer the answer to the queries. This library, which is specifically designed for creating and working with Bayesian Networks, takes our generated CPTs as input, verifies their soundness and answers the queries via an exact inference method.

Methodology

Here, we introduce our approach to probabilistic reasoning with LLMs. We use a basic QA and a Chain of Thought prompting as baselines and propose new strategies for mapping to symbolic representations.

Baselines

Basic QA Prompting In this prompting approach, we ask the model to generate a single numerical answer to the probabilistic question. We experiment with zero-shot and fewshot settings. In the zero-shot setting, only the instruction, context, and question are in the prompt. Figure 2 shows the few-shot setting with an in-context example included.

Chain of Thought (COT) Following COT (Wei et al. 2022a), we prompt the LLM to explain its reasoning process while refraining completely (in zero-shot setting) or partially (in few-shot setting) from imposing a strict solution structure. COT's prompting structure is similar to Basic QA's, except that the requested answer should explain the mathematical reasoning, calculate the final answer, and generate the target output probability at the end, as shown in Figure 2.

Structured Prompting with Subtasks

Given the complexity of the probabilistic inference, we propose to divide the problem into multiple steps and demonstrate the steps to the LLMs in one prompt. This approach has shown to be effective in other similar research (Jin et al. 2023; Poesia et al. 2023). The most intuitive step for our problem is identifying the probabilities of the single events and the conditional probabilities. Another important step is to recognize the variables' dependencies, which are the edges in the corresponding BN. Consequently, we use the extraction of probability values from text, named *Number Extraction* subtask, and probabilistic dependencies *Graph Generation* subtask as the prior steps to final reasoning.

Number Extraction In this subtask, the LLM should extract the CPT probabilities from the input context and output them in a structured format as Python-compatible variable assignments. Each line of the output represents either a probability of an event or a conditional probability. This format is shown in the "PAL" column of Figure 2. To add this subtask, we prepend its corresponding instruction to the prompt, i.e., "Extract the probabilities..." and its answers to the in-context examples. This subtask aims to facilitate the



Figure 2: This figure shows our main prompting approaches, PAL, Monte Carlo, and ProbLog, alongside the baseline approaches, Basic QA and COT. Each prompt begins with an instruction (purple boxes) that describes the problem and the answer format. Then, the context, question, and answer are demonstrated depending on the approach. We display only the first incontext example here but use 3 in our experiments. When we require the use of our designed subtasks in the prompt, their instructions and answers are prepended to the main approach, as shown in the PAL method for *Number Extraction* and the Monte Carlo method for *Graph Generation*.

correct extraction of probabilities before reasoning to avoid hallucination (Ouyang et al. 2022) of incorrect numbers.

Graph Generation In this subtask, we want the LLM to generate the underlying Bayesian Network of the given context as a list of edges, each indicating the direct dependency between two variables. For a BN with V variables, the output should consist of V - 1 edges in the format of $v_i - > v_j$ separated by ','. This will help the model capture the random variables' dependency structure and utilize it in mapping to symbolic solutions. Similar to *Number Extraction*, to use *Graph Generation*, its instruction and answers are added to the prompts. An example of a simple Bayesian Network with two nodes, Green and Pink, is shown in the "Monte Carlo Inference Algorithm" column of Figure 2.

Mapping to Symbolic Computations

Program-aided Language Models (PAL) PAL (Gao et al. 2023) is the first study to analyze the use of Python in various mathematical reasoning QA datasets. However, most problems tested in PAL require only a few lines of code, unlike BLInD, which may need complex, multi-line calculations depending on the BN. A benefit of the PAL method is it bypasses the challenge of mathematical calculations by the LLM itself. Similar to the original PAL paper, we instruct the LLM to solve the problem by explaining the mathematical reasoning process and mapping to the basic arithmetic calculations in Python code leading to the

answer. Figure 2 shows an in-context example of this approach, where the *Number Extraction* subtask is first used, followed by mapping the reasoning solution to Python code.

Monte Carlo Inference Algorithm Given the efficiency and popularity of Monte Carlo Algorithms (Koller and Friedman 2009) for approximate inference, we try to map our problem the *Direct Sampling* technique for inference. An LLM can use this method by generating a Python function, we call *simulate*, that samples all events according to the probabilistic dependencies expressed in the BN. Here, all parent variables must be sampled before their children. Keeping this order is the main challenge in mapping to this algorithm with LLMs. Figure 2 shows an in-context example of the Monte Carlo method with the *simulate* function defined as a part of the answer. The LLM is also instructed to call this function in the generated Python code which leads to computing the answer to the probabilistic question.

Probabilistic Logical Solver (ProbLog) In our neurosymbolic method, we employ a technique that involves mapping the context and question to a probabilistic logical formalism. We use ProbLog (De Raedt, Kimmig, and Toivonen 2007), a probabilistic programming language that extends Prolog (Bratko 2000) to incorporate probabilistic logical reasoning. Here, the LLM is asked to generate a ProbLog code corresponding to the probabilities given in the context and to create the formal *ProbLog query* based on the question. We subsequently execute the ProbLog code and extract the final answer. An in-context example of this code is shown in Figure 2. The LLM does not need extensive ProbLog programming knowledge; The three in-context examples we supply are sufficiently complex and encompass all the necessary ProbLog syntax information to enable the generation of code for our contexts and questions.

Experiments

Here, we present the results of our experiments on our baselines and proposed prompting techniques, which we evaluate on BLInD and an adapted version of the CLADDER dataset. Refer to the Appendix of the *arXiv* version of the paper for additional information, including our models' hyperparameters (the link is provided below the abstract).

LLM Models In our experiments, we employ three LLMs: Llama3 (AI@Meta 2024), specifically the *meta-llama-3-70b-instruct* variant; GPT3.5 (Brown et al. 2020), using the *gpt-3.5-turbo-0613* release; and GPT4 (OpenAI 2023), with the *gpt-4-0613* version. These models are evaluated in zero-shot and few-shot settings without any fine-tuning.

Few-shot Example Selection We selected a set of shots from a development dataset and manually crafted their solutions. After evaluating these shots on the same development dataset, we identified the three most effective examples through an iterative, trial-and-error process. To ensure a fair comparison, we consistently use these three examples across all models and methods rather than tailoring the examples to each specific approach or model.

Evaluation Metrics. Given a context and a question, we consider an answer probability to be correct if it is within the ± 0.01 range of the ground truth probability (ex. any answer within [0.30-0.32] is correct for a ground truth of 0.31). We chose this threshold because we found that the outputs were either correct or wrong with a large margin since the exact line of computations is not followed in those cases. This bimodal behavior, which differs from traditional regression models, renders evaluation metrics such as MSE and L1 ineffective. In this context, correct predictions contribute minimally to the error, while incorrect predictions dominate the error metric in a way that lacks relevance. Additionally, this behavior rendered larger thresholds useless, as the accuracy at a threshold of 0.01 was nearly identical to that at 0.05. A narrower threshold would cause the challenge of number precision which is not in our interest due to the nature of our task and has been previously highlighted in (Gao et al. 2023) for other mathematical reasoning problems. For the evaluation of the subtasks, we count an output as accurate if **all** the numbers in Number Extraction and all the edges in Graph Generation are correctly generated without redundancy. As a result, the numbers in all Tables are accuracy values in percentages based on these criteria.

Evaluation Splits of the Dataset. To assess our methods, we randomly select 100 instances from each data split V_i , resulting in a total of 900 instances. This test set remains consistent across all of our LLMs.

Solving Probabilistic Questions Directly

Here, we apply the baseline methods of Basic QA and COT, focusing on answering probabilistic questions directly. Their performance is detailed in Table 1. In Basic QA, overall, the results are very low, and only GPT4 achieves meaningful results for some of the dataset splits with **smaller BNs**, i.e. V_i with $i \leq 5$. In the few-shot setting of Basic QA, the additional examples, which do not explain their solutions, worsen the results for all the LLMs. Using COT improved the results for all models. However, even with COT, these models struggle particularly in dataset splits with **larger BNs**, i.e. V_i with i > 5. We will use these baselines to compare with our symbolic methods.

Subtasks

Before reporting the results of the final symbolic solvers, we discuss the results of *Number Extraction* and *Graph Generation*. The results of *Number Extraction* are shown in Table 2, which indicates this subtask is quite straightforward to solve. Llama3 and GPT4 extract all numbers correctly, achieving 100% accuracy in all V_i s. For GPT3.5, although the accuracy drops as the number of variables increases, it remains overall very high above 90%.

The results of the more challenging *Graph Generation* subtask are shown in Table 3. Mirroring the pattern observed in Table 2, we notice a decline in accuracy as the number of variables increases. However, the drop in accuracy is more notable and goes from 100% in V_2 to as low as 73% in V_9 for GPT3.5. GPT4 generated all graphs correctly in all V_i s. Looking at the results in Table 3, we observe minor inconsistencies such as a lower accuracy for V_9 compared to V_{10} . These inconsistencies stem from the inherent randomness in the output generation of LLMs and our random selection dataset instances. These small inconsistencies happen in some other parts of our experiments but they do not detract from the core message and pattern of our findings.

Note that the accuracies reported here are calculated when each subtask is prompted to the LLM as a standalone problem. When integrating these subtasks within our solutions, we prompt the LLM to generate both the subtask and the problem solution together. This affects the subtasks' accuracy and, consequently, usefulness depending on the symbolic method, as discussed in the next section.

Evaluation of Proposed Methods

Here, we assess our three proposed approaches, PAL, Monte Carlo, and ProbLog combined with *Number Extraction* and *Graph Generation* in three LLMs. We discuss how effective the subtasks are with each method and analyze their impact based on factors like the number of variables and the employed LLM. The results of these experiments are presented in Table 4. Not all of the combinations are useful; some lead to lower final accuracy. These underperforming configurations are not presented in the table.

PAL, Monte Carlo, and ProbLog As seen in Table 4, there is a significant improvement in the performance of all of these methods, compared to Basic QA and COT (previously shown in Table 1). Additionally, accuracy consistently

Model	Method	V_2	V_3	V_4	V_5	V_6	V_7	V_8	V_9	V10	V_{2-5}	V_{6-10}	V_{2-10}
	Basic QA ZS	33	13	5	4	6	2	3	1	2	13	2	7
CDT2 5	Basic QA FS	3	0	1	1	2	2	1	1	0	1	1	1
0F15.5	COT ZS	53	8	4	5	10	5	2	2	0	17	3	9
	COT FS	52	23	12	5	8	4	1	4	2	23	3	12
	Basic QA ZS	31	21	5	6	6	5	1	1	0	16	3	8
L lama 3	Basic QA FS	3	0	1	1	2	2	1	1	0	1	1	1
Liamas	COT ZS	63	45	21	17	18	11	9	4	2	37	9	21
	COT FS	63	46	21	12	20	15	7	8	5	36	11	22
	Basic QA ZS	44	23	9	9	11	11	8	8	2	21	8	14
CDT4	Basic QA FS	3	0	1	1	2	2	1	1	0	1	1	1
OF14	COT ZS	79	63	27	10	17	6	5	7	6	45	8	24
	COT FS	78	64	36	25	22	16	7	7	7	50	12	29

Table 1: Comparison of GPT3.5, Llama3, and GPT4 accuracy results for Basic QA and COT methods, presented as percentages. The columns represent dataset splits V_i , and the average results for smaller BNs V_{2-5} , larger BNs V_{6-10} , and all BNs V_{2-10} . The rows show the methods tested with zero-shot (ZS) or few-shot (FS) settings.

LLM / V_i	V_2	V_3	V_4	V_5	V_6	V_7	V_8	V_9	V_{10}
GPT3.5	100	100	100	100	96	95	98	94	94
Llama3	100	100	100	100	100	100	99	100	100
GPT4	100	100	100	100	100	100	100	100	100

Table 2: *Number Extraction* accuracy of our models, presented as percentages and based on the exact match of all the extracted probabilities of the context.

LLM / V_i	V_2	V_3	V_4	V_5	V_6	V_7	V_8	V_9	V_{10}
GPT3.5	100	95	92	93	84	75	79	73	78
Llama3	99	99	99	100	99	95	94	93	89
GPT4	100	100	100	100	100	100	100	100	100

Table 3: *Graph Generation* accuracy, presented as percentages. The extracted graph should exactly match the correct BN graph to be counted as correct.

increases across all models (closed and open-source) by transitioning from PAL to Monte Carlo and then to ProbLog. This suggests that the proposed methods' effectiveness is independent of the LLMs. All models struggle to generate a solution with PAL for larger BNs. In contrast, when we utilize the Monte Carlo approach, the accuracy of these larger BNs sharply increases, suggesting proficiency of LLMs at mapping the entire BN correctly to a Monte Carlo algorithm code, even for a large number of variables.

ProbLog eliminates the challenge of structural programming and requires only the correct extraction of probabilities (represented declaratively) and generating a corresponding *ProbLog query*. In this case, GPT4 can solve almost every question. GPT3.5 is mainly held back by the challenge of writing probabilistic logical programming code. While Llama3 (like GPT4) featured nearly 100% correct Python syntax in the PAL and Monte Carlo methods, it sometimes fails to create coherent ProbLog code. This leads to somewhat inconsistent performance among smaller BNs.

PAL with Number Extraction This combination, which is shown as "PAL w/NE" in Table 4, shows that the accuracy of both GPT3.5 and Llama3 benefits from the addition

of the *Number Extraction* subtas to the PAL prompt. This appears to reduce the hallucination (Ouyang et al. 2022) of probability values in PAL solutions, as we further confirmed by analyzing several test cases. This subtask was not needed for the more robust GPT4, which can remember the numbers and its addition resulted in marginal improvements.

Monte Carlo with Graph Generation The accuracy of GPT3.5 and GPT4 improved when the Monte Carlo method was combined with the *Graph Generation* subtask, shown as 'Monte Carlo w/GG'' in Table 4. *Graph Generation* subtask further improves the already high performance of Monte Carlo for larger BNs and enables GPT4 to reach a nearperfect average accuracy of 95% in this setting. The improvements caused by the addition of *Graph Generation* are not surprising since the Monte Carlo method generates a Python function with many nested "if" structures tied to the underlying Bayesian Network's graph structure. However, Llama3 is the exception and the only model that does not benefit from this configuration, as discussed further below.

Discussion

Practical Use of Subtasks While our proposed approaches proved effective, independently of the LLMs, this was not true for our subtasks. As mentioned earlier, we comprehensively tested all the configurations, but not all of them improved our models. This raises a few questions: Q1) Why do the added subtasks not always help? For example, in principle, the Monte Carlo method could use *Number Extraction* in its code to improve. Q2) Why does ProbLog not improve with any subtasks? Q3) Why does Llama3 not improve with *Graph Generation* in its Monte Carlo method like GPT3.5 and GPT4? Q4) Why does no method improve by adding both subtasks together? To answer these questions, we looked at subtask generation and their utilization by LLMs more closely, which led to two main findings.

The first finding is that in contrast to most mathematical problems tested with LLMs, which require brief solutions (Kim et al. 2023; Frieder et al. 2023), our dataset demands the generation of large outputs. When the added information by subtasks is not exploited effectively and

Model	Method	V_2	V_3	V_4	V_5	V_6	V_7	V_8	V_9	V_{10}	V_{2-5}	V_{6-10}	V_{2-10}
	PAL	66	34	25	17	14	9	6	5	2	35	7	19
	PAL w/NE	85	66	41	27	19	12	5	3	6	54	9	29
GPT3.5	Monte Carlo	79	63	71	65	41	32	33	18	14	69	27	46
	Monte Carlo w/GG	85	82	83	68	42	31	28	18	8	79	25	49
	ProbLog	87	82	88	75	59	52	46	38	35	83	46	62
	PAL	100	84	57	36	31	20	10	14	8	69	17	40
	PAL w/NE	100	95	71	52	46	28	16	16	9	79	23	48
Llama3	Monte Carlo	100	100	96	96	92	85	77	72	64	98	78	87
	ProbLog	90	95	92	87	95	94	87	82	78	91	87	89
	PAL	100	86	70	58	50	27	21	14	7	78	24	48
CPT4	PAL w/NE	99	96	78	64	43	26	14	14	10	84	21	49
0114	Monte Carlo	100	99	98	100	92	94	92	90	88	99	91	94
	Monte Carlo w/GG	100	97	99	98	97	96	88	92	85	99	92	95
	ProbLog	99	98	100	100	96	97	97	98	96	99	97	98

Table 4: GPT3.5, Llama3, and GPT4 accuracy results, presented as percentages, for the PAL, Monte Carlo, and ProbLog methods. w/NE and w/GG denote the inclusion of *Number Extraction* and *Graph Generation*. The columns represent dataset splits V_i , and the average results for smaller BNs V_{2-5} , larger BNs V_{6-10} , and all BNs V_{2-10} .

lengthens the output even more for no reason, it leads to a notable drop in the LLM performance. For example, while the graph structure is intuitively helpful for probabilistic inference, the Python code in PAL does not utilize it directly. This directly addresses Q1 and touches on Q4. The main bottleneck of the ProbLog method was the syntax errors that subtasks could help with, which answers Q2.

The second finding concerns the accuracy of the subtasks, which drops when generated in the same prompt with the main solution (as we prompt the LLM only once). This puts Llama3 in a precarious position regarding the Monte Carlo method with its high accuracy. For Graph Generation to further improve this method, it has to have an accuracy higher than the method to be helpful. However, that is not the case for this configuration for Llama3. For instance, the accuracy of Llama3 in the Monte Carlo method for V_{10} is 64% (Table 4), which is already higher than Graph Generation accuracy for V_{10} that is 56% when generated in the same prompt (See the Supplementary for detailed results). For GPT3.5 and GPT4, Graph Generation accuracy remains high enough, which in the case of GPT3.5 is partially due to its weaker performance in the Monte Carlo method. This finding resolves Q3 and provides further insights into Q4.

Trade-off Between Complexity and Effectiveness Among our 5 methods, the Simplest and the most efficient one is the Basic QA, which generates a few tokens. COT slightly improves the results at the cost of more tokens in the input and output. There is a significant improvement in accuracy, moving to our main methods with external tools. Their LLM code generation time is the same as COT, but they need the additional time to execute the generated program. According to our experiments, the probabilistic inference run-time is negligible compared to the inference run-time of LLM output generation which is a few milliseconds versus seconds taken by LLMs. However, from an algorithmic perspective for probabilistic inference, ProbLog will be more complex compared to Monte Carlo sampling as it needs to deal with logical representations. Given that a probabilistic network question described in a natural language context forms rather small Bayesian Networks, our methods are practical for solving this problem.

Use of External Tools Using external tools with LLMs is an area of research that leverages the LLMs and exploits diverse computational paradigms (Schick et al. 2024). The tools we use are 1) pure Python for PAL and Monte Carlo methods and 2) Python plus the underlying ProbLog engine. All these tools are open-source, and conversion to their Python interface is highly accurate using language models. They are more efficient compared to LLMs, as discussed above, and their operation as black-box executables requires minimal computing resources.

Adaptation of the CLADDER Dataset

We conclude our experiments by testing our methods on an adaptation of the CLADDER dataset (Jin et al. 2023). This dataset is designed to test the causal reasoning capabilities of LLMs. The contexts of this QA dataset describe a probabilistic causal structure with a maximum of 4 variables, designed from natural-sounding templates. Questions in this dataset mostly require a binary yes/no answer and not a probability. Our results are, thus, not comparable to the ones in (Jin et al. 2023). Using the natural-sounding contexts in CLADDER's *hard tests* split, we created challenging queries for the contexts and sample 100 instances. Figure 3 provides an example of this dataset along with one of our generated queries for its context.

The results of our tests on this adaptation of CLADDER are shown in Table 5, which follow the same trend seen in BLInD for smaller BNs. For example, adding *Graph Generation* to the Monte Carlo method does not improve the model here as it was most helpful when the number of variables was large. This consistency with BLInD evaluations further solidifies our claims. When testing our methods on the CLADDER, we used the same in-context examples of BLInD without tailoring them to the more natural contexts of CLADDER as we found it unnecessary. The performance



Figure 3: An example from the hard test subset of CLAD-DER dataset and a corresponding generated probabilistic query. The top section displays the context with events in bold font (white box), a query (yellow box), and the binary (Yes/No) answer (purple box). The bottom section presents an example of a probabilistic query derived from the same context, which requires a probability-based response.

Method	GPT3.5	Llama3	GPT4
Basic QA ZS	0	0	20
Basic QA FS	0	0	0
COT ZS	9	47	65
COT FS	3	38	64
PAL	26	91	96
PAL w/NE	39	96	96
Monte Carlo	75	96	98
Monte Carlo w/GG	75	95	97
ProbLog	71	84	97

Table 5: Accuracy results of the CLADDER dataset as percentages. w/NE, w/GG, ZS and FS denote use of *Number Extraction*, *Graph Generation*, zero-shot and few-shot.

remained very high, as seen in Table 5. Based on the results from BLInD and CLADDER, our experiments suggest that the difficulty of probabilistic reasoning over text is not directly correlated with the naturalness and sophistication of the language. Instead, it depends on the depth of reasoning required and the number of variables involved.

Conclusion and Future Work

In this work, we introduced BLInD, a new dataset for dealing with uncertain text and evaluating the capabilities of LLMs on probabilistic reasoning over text with explicitly quantified uncertainty. We proposed several prompt engineering techniques, mapping the problem to different formal representations, including Python low-level arithmetic computations, approximate inference algorithms, and probabilistic logical programming. Our evaluations demonstrated that our main methods significantly improve the performance of LLMs on BLInD and on an adapted version of another dataset, CLADDER, with natural-sounding contexts.

Our methods solve probabilistic questions without any fine-tuning or modification to the architecture of LLMs. As an interesting direction to continue this work, future research could explore alterations to open-source LLMs' architectures and training objectives specifically designed for probabilistic inference.

Acknowledgments

This project is supported by the Office of Naval Research (ONR) grant N00014-23-1-2417. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of Office of Naval Research.

References

AI@Meta. 2024. Llama 3 Model Card. https://github.com/ meta-llama/llama3/blob/main/MODEL_CARD.md. Accessed: 2024-07-01.

Ankan, A.; and Panda, A. 2015. pgmpy: Probabilistic graphical models using python. In *Proceedings of the 14th Python in Science Conference (SCIPY 2015)*. Citeseer.

Bratko, I. 2000. *Prolog Programming for Artificial Intelligence*. Harlow, England: Pearson Addison-Wesley, 3 edition. ISBN 978-0-201-40375-6.

Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877– 1901.

Bubeck, S.; Chandrasekaran, V.; Eldan, R.; Gehrke, J.; Horvitz, E.; Kamar, E.; Lee, P.; Lee, Y. T.; Li, Y.; Lundberg, S.; et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.

Cobbe, K.; Kosaraju, V.; Bavarian, M.; Chen, M.; Jun, H.; Kaiser, L.; Plappert, M.; Tworek, J.; Hilton, J.; Nakano, R.; et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

De Raedt, L.; Kimmig, A.; and Toivonen, H. 2007. ProbLog: A Probabilistic Prolog and Its Application in Link Discovery. In *Proceedings of the 20th International Joint Conference on Artifical Intelligence*, IJCAI'07, 2468–2473. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Burstein, J.; Doran, C.; and Solorio, T., eds., *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)*, 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics.

Dries, A.; Kimmig, A.; Davis, J.; Belle, V.; and de Raedt, L. 2017. Solving Probability Problems in Natural Language. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, 3981–3987.

Frieder, S.; Pinchetti, L.; Chevalier, A.; Griffiths, R.-R.; Salvatori, T.; Lukasiewicz, T.; Petersen, P. C.; and Berner, J. 2023. Mathematical Capabilities of ChatGPT. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*

Gao, L.; Madaan, A.; Zhou, S.; Alon, U.; Liu, P.; Yang, Y.; Callan, J.; and Neubig, G. 2023. PAL: program-aided language models. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Google. 2023. Google Gemini AI. https://blog.google/ technology/ai/google-gemini-ai/#availability. Accessed: 2024-07-01.

He-Yueya, J.; Poesia, G.; Wang, R. E.; and Goodman, N. D. 2023. Solving math word problems by combining language models with symbolic solvers. *arXiv preprint arXiv:2304.09102*.

Heritage, J. 2013. Action formation and its epistemic (and other) backgrounds. *Discourse Studies*, 15(5): 551–578.

Jin, Z.; Chen, Y.; Leeb, F.; Gresele, L.; Kamal, O.; LYU, Z.; Blin, K.; Adauto, F. G.; Kleiman-Weiner, M.; Sachan, M.; and Schölkopf, B. 2023. CLadder: A Benchmark to Assess Causal Reasoning Capabilities of Language Models. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Kim, J.; Kim, Y.; Baek, I.; Bak, J.; and Lee, J. 2023. It Ain't Over: A Multi-aspect Diverse Math Word Problem Dataset. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Proceedings* of the 2023 Conference on Empirical Methods in Natural Language Processing, 14984–15011. Singapore: Association for Computational Linguistics.

Koller, D.; and Friedman, N. 2009. *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press, hard-cover edition. ISBN 9780262013192.

Landmark, A. M. D.; Gulbrandsen, P.; and Svennevig, J. 2015. Whose decision? Negotiating epistemic and deontic rights in medical treatment decisions. *Journal of Pragmatics*, 78: 54–69. Epistemics and Deontics in Conversational Directives.

Mishra, S.; Mitra, A.; Varshney, N.; Sachdeva, B.; Clark, P.; Baral, C.; and Kalyan, A. 2022. NumGLUE: A Suite of Fundamental yet Challenging Mathematical Reasoning Tasks. In Muresan, S.; Nakov, P.; and Villavicencio, A., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 3505– 3523. Dublin, Ireland: Association for Computational Linguistics.

Nafar, A.; Venable, K. B.; and Kordjamshidi, P. 2024. Teaching Probabilistic Logical Reasoning to Transformers. In Graham, Y.; and Purver, M., eds., *Findings of the Association for Computational Linguistics: EACL 2024*, 1615– 1632. St. Julian's, Malta: Association for Computational Linguistics.

Oaksford, M.; and Chater, N. 2007. *Bayesian Rationality: The probabilistic approach to human reasoning*. Oxford University Press. ISBN 9780198524496.

Oaksford, M.; and Chater, N. 2009. Précis of Bayesian Rationality: The Probabilistic Approach to Human Reasoning. *Behavioral and Brain Sciences*, 32(1): 69–84.

OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774.

Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.

Ozturkler, B.; Malkin, N.; Wang, Z.; and Jojic, N. 2023. ThinkSum: Probabilistic reasoning over sets using large language models. In Rogers, A.; Boyd-Graber, J.; and Okazaki, N., eds., *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1216–1239. Toronto, Canada: Association for Computational Linguistics.

Poesia, G.; Gandhi, K.; Zelikman, E.; and Goodman, N. D. 2023. Certified Deductive Reasoning with Language Models. arXiv:2306.04031.

Poggi, I.; D'Errico, F.; Vincze, L.; et al. 2019. Uncertain words, uncertain texts. perception and effects of uncertainty in biomedical communication. *Acta Polytechnica Hungarica*, 16(2): 13–34.

Rajaby Faghihi, H.; Guo, Q.; Uszok, A.; Nafar, A.; and Kordjamshidi, P. 2021. DomiKnowS: A Library for Integration of Symbolic Domain Knowledge in Deep Learning. In Adel, H.; and Shi, S., eds., *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 231–241. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.

Rajaby Faghihi, H.; Nafar, A.; Zheng, C.; Mirzaee, R.; Zhang, Y.; Uszok, A.; Wan, A.; Premsri, T.; Roth, D.; and Kordjamshidi, P. 2023. GLUECons: A Generic Benchmark for Learning under Constraints. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(8): 9552–9561.

Saeed, M.; Ahmadi, N.; Nakov, P.; and Papotti, P. 2021. RuleBERT: Teaching Soft Rules to Pre-Trained Language Models. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 1460–1476. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.

Schick, T.; Dwivedi-Yu, J.; Dessì, R.; Raileanu, R.; Lomeli, M.; Hambro, E.; Zettlemoyer, L.; Cancedda, N.; and Scialom, T. 2024. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36.

Shi, Z.; Zhang, Q.; and Lipani, A. 2022. Stepgame: A new benchmark for robust multi-hop spatial reasoning in texts. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, 11321–11329.

Suster, S.; Fivez, P.; Totis, P.; Kimmig, A.; Davis, J.; de Raedt, L.; and Daelemans, W. 2021. Mapping probability word problems to executable representations. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Proceedings of the 2021 Conference on Empirical Methods in Nat-* *ural Language Processing*, 3627–3640. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.

Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; ichter, b.; Xia, F.; Chi, E.; Le, Q. V.; and Zhou, D. 2022a. Chainof-Thought Prompting Elicits Reasoning in Large Language Models. In Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; and Oh, A., eds., *Advances in Neural Information Processing Systems*, volume 35, 24824–24837. Curran Associates, Inc.

Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.