

LET’S THINK VAR-BY-VAR: LARGE LANGUAGE MODELS ENABLE *Ad Hoc* PROBABILISTIC REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

A hallmark of intelligence is the ability to flesh out underspecified situations using “common sense.” We propose to extract that common sense from large language models (LLMs), in a form that can feed into probabilistic inference. We focus our investigation on *guesstimation* questions such as “How much are Airbnb listings in Newark, NJ?” Formulating a sensible answer without access to data requires drawing on, and integrating, bits of common knowledge about how Price and Location may relate to other variables, such as Property Type. Our framework answers such a question by synthesizing an *ad hoc* probabilistic model. First we prompt an LLM to propose a set of random variables relevant to the question, followed by moment constraints on their joint distribution. We then optimize the joint distribution p within a log-linear family to maximize the overall constraint satisfaction. Our experiments show that LLMs can successfully be prompted to propose reasonable variables, and while the proposed numerical constraints can be noisy, jointly optimizing for their satisfaction reconciles them. When evaluated on probabilistic questions derived from three real-world tabular datasets, we find that our framework performs comparably to a direct prompting baseline in terms of total variation distance from the dataset distribution, and is similarly robust to noise.

1 INTRODUCTION

Thus, in reasoning we depend very much on *prior information* to help us in evaluating the degree of plausibility in a new problem. This reasoning process goes on unconsciously, almost instantaneously, and we conceal how complicated it really is by calling it common sense. —E. T. Jaynes, *Probability Theory: The Logic of Science* (2003)

Humans constantly reason about novel situations, integrating evidence with prior knowledge. The Jaynes (2003) quote above refers to an everyday example: a policeman sees a masked man with a bag crawling out of the broken window of a jewelry store, and suspects a burglary. How can such conclusions be arrived at—appropriately generating hypotheses and weighing competing evidence?

Like Jaynes, we hope to draw on the very same methods of statistical modeling and inference that allow scientists to reason formally about complex domains like epidemiology, diplomacy, or syntax. For those domains, however, scientists normally invest time in perfecting a durable scientific model that supports many queries. Commonsense reasoning may instead generate a quick-and-dirty *ad hoc* model for each query.

We show that one can construct such ephemeral models automatically by enlisting the existing commonsense knowledge of large language models (LLMs). Of course, today’s LLMs are already smart enough to recognize the above scene as a burglary—either at once, or via a chain of thought that may explicitly generate and evaluate different hypotheses. But there are harder situations that may benefit from systematically eliciting many fragments of relevant knowledge from the LLM, and deriving conclusions from this combined knowledge in a more formal and systematic way.

While one could elicit *logical* propositions and derive conclusions from those (Jung et al., 2022), we consider here the more general case of *probabilistic* knowledge and conclusions. Consider a

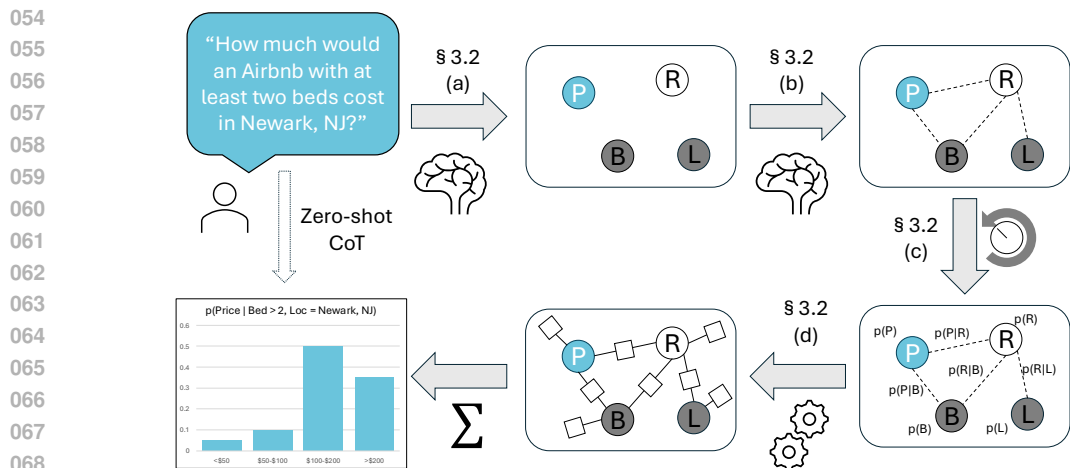


Figure 1: An illustration of our proposed framework applied to answering an example probabilistic question, $Q = \text{“How much would an Airbnb with at least two beds cost in Newark, NJ?”}$. Going **clockwise** from Q , we first prompt an LLM to brainstorm the relevant random variables (§3.2 (a)), producing Price (P), Rating (R), Beds (B), Location (L), where shaded nodes denote variables being conditioned on, blue nodes denote target variables, and white nodes denote latent variables. Then we prompt an LLM to propose interacting pairs $\{v_1, v_2\}$ of proposed variables, and whether to constrain $p(v_1 | v_2)$ or $p(v_2 | v_1)$ (§3.2 (b)). Next we prompt LLMs to propose numeric constraints on the marginal $p(v)$ of every proposed variable, as well as the conditional marginals $p(v_1 | v_2)$ of every proposed pairwise interaction (§3.2 (c)); Finally, we optimize the parameters of a log-linear model with fuzzy maximum entropy objective (2) in order to maximize constraint satisfaction (§3.2 (c)). The final output is an ad hoc probability model that can be used to answer Q . Going **counter-clockwise** from Q is a baseline of asking for an estimate of Q directly using a zero-shot LLM with Chain-of-Thought.

guesstimation question such as “How many people in Nigeria own laptops?”¹ An LLM that has only weak intuitions about this may nonetheless be able to recall various relevant information:

- One route to an answer would estimate Nigeria’s distribution over occupations, and then estimate those occupations’ distributions over computing devices. It is relevant that Nigeria is a developing country and that some developing countries have largely skipped over laptops to mobile phones.
- Another route would estimate Nigeria’s wealth distribution and its ownership rates for other appliances (cars, dishwashers, cellphones), and then guess how a person’s laptop ownership correlates with their wealth and possessions.
- Another route might look at historical data (if known) and try to extrapolate to the present.
- The above bullets estimate Nigeria’s *rate* of laptop ownership, which must be multiplied by Nigeria’s population. If the population is not known, it could be guessed based other facts, such as Nigeria’s physical size and political influence relative to nearby countries, or the relative visibility of Nigerians in global culture.

Integrating all of this information *systematically* may provide a more robust answer than simply asking the LLM to answer directly or to think step-by-step. We do this by constructing an *ad hoc* probability model over situations, with latent variables and their interactions proposed by the LLM.

Though the LLM proposes the model’s structure, we do not expect the LLM to provide its *parameters*. In general, such parameters are not interpretable.² Rather, we ask the LLM to make predictions

¹This may arise in the course of solving another guesstimation question: “If my aging laptop fails during my trip to Lagos, how long will it take to repair?”

²In a Markov random field (MRF), the optimal parameters for one factor are not a property of that factor alone, but depend strongly on what other factors have been added to the model and what their parameters are.

108 about the world—such as marginal probabilities. We set the model parameters so as to align the
 109 model’s predictions with the LLM’s predictions. We can then query our model to answer the
 110 original guesstimation question (via probabilistic inference over the situations described by the model
 111 variables).

112 This paper will focus on specific guesstimation questions where we (as experimenters) are able to
 113 evaluate answer quality. In §5, we evaluate our approach on three real-world datasets, Inside Airbnb
 114 (AIR)³, American Time-Use Survey (ATUS)⁴, and World Values Survey (WVS)⁵, by comparing our
 115 system’s answers to the answers estimated from these datasets. We develop our prompts on subsets of
 116 Inside Airbnb and American Time-Use Survey, and evaluate on held-out subsets of these two datasets,
 117 as well as on World Value Survey, which we held out completely during system development.

118 2 PROBLEM SETUP

119 Let Q denote a question about some *novel* situation to some agent—in the sense that there is not
 120 enough prior experience to answer the question *directly*. Concretely, consider the example question,
 121 “What would the age be for a widow living in California?” Without direct prior knowledge (e.g.
 122 from having met many Widows in California or from looking up census data), formulating sensible
 123 answers to such questions requires drawing on and integrating bits of common knowledge about
 124 how Widowness and Location may relate to *other* variables like Occupation of their spouse, and
 125 whether they have any Children.

126 We can formalize such a question as a probabilistic query for a particular conditional distribution,
 127 $p(y \mid \mathbf{x} \in \mathbb{S})$, where y is the target variable, \mathbf{x} are the conditioning variables, and $\mathbb{S} \subseteq \mathbb{X}$ is the
 128 event being conditioned on. The example question above can be formalized this way as a query for
 129 $p(\text{Age} \mid \text{Location} = \text{California}, \text{Widow} = \text{True})$. Given such a question Q , our task is to generate
 130 an estimate $\hat{p}(y \mid \mathbf{x} \in \mathbb{S})$ without relying on direct data.

131 For simplicity, our investigation will focus on questions where there is a single target variable, where
 132 all variables are discrete (we discretize continuous variables into ranges), but it would not be difficult
 133 to generalize our method to more than one target variable and to handle continuous variables directly
 134 (via approximate inference methods such as Minka (2013)).

135 3 METHOD

136 Our framework extracts knowledge from LLMs and integrates it to build an *ad hoc* probability model
 137 that can be used to answer the kind of questions described in §2. In §3.1, we formalize the notion of
 138 common knowledge relevant to some question Q , and how such knowledge can be integrated in a
 139 principled way to yield an *ad hoc* probability model. In §3.2, we describe how we instantiate the
 140 formalization with a prompted-LLM as the source of common knowledge, as well as specific choices
 141 we made in terms of parameterizing the ad hoc models.

142 3.1 INTEGRATING COMMON KNOWLEDGE VIA FUZZY MOMENT MATCHING

143 **Moment Constraints** Given a question Q , for example, “What would the age be for a widow living
 144 in California?”, what kinds of prior knowledge might be helpful for answering it? Our main insight is
 145 to extract prior knowledge in the form of *moment-matching constraints*, that is, constraints on the
 146 (conditional) marginals over random variables that are relevant to the question Q .

147 Let’s suppose for now that we are supplied with a set of variables that are relevant to the question Q ,⁶
 148 which includes the target variable y , the conditioning variables \mathbf{x} , and some latent variables \mathbf{z} . Our

149 ³<https://insideairbnb.com/>

150 ⁴<https://www.bls.gov/tus/>

151 ⁵<https://www.worldvaluessurvey.org/wvs.jsp>

152 ⁶The judgement of relevance of a random variable x_1 to x_2 is a kind of prior knowledge about their joint
 153 distribution. For example, relevance could be formalized as a threshold on the mutual information $I(x_1, x_2)$,
 154 which can be derived from their joint marginal.

constraints $c_1(p), \dots, c_n(p)$ on the conditional expectations of the joint distribution p take the form

$$c_i(p) : b_i = \mathbb{E}_p [f_i(y, \mathbf{x}, \mathbf{z}) \mid g_i(y, \mathbf{x}, \mathbf{z})] \stackrel{\text{def}}{=} \frac{\mathbb{E}_{(y, \mathbf{x}, \mathbf{z}) \sim p} [f_i(y, \mathbf{x}, \mathbf{z}) \cdot g_i(y, \mathbf{x}, \mathbf{z})]}{\mathbb{E}_{(y, \mathbf{x}, \mathbf{z}) \sim p} [g_i(y, \mathbf{x}, \mathbf{z})]} \quad (1)$$

where g_i is an indicator function and f_i is a real-valued feature function.⁷

Why do we formalize prior knowledge as constraints on the *distribution* p rather than its *parameters*? The optimal parameters of a probability distribution are often interdependent and change with the model structure. Adding new latent variables \mathbf{z} to a model may change the optimal parameters in other parts of the model. However, conditional expectations are stable across different model structures since they are properties of the world, not properties of the model. This makes it possible to elicit them individually from an LLM.

Estimation Objective The constraints will be drawn from an LLM and may not be wholly correct. We optimize p to *approximately* satisfy the constraints via

$$\operatorname{argmin}_p -H(p) + \sum_i w_i (b_i - \mathbb{E}_p [f_i(y, \mathbf{x}, \mathbf{z}) \mid g_i(y, \mathbf{x}, \mathbf{z})])^2 \quad (2)$$

The hyperparameter w_i specifies the importance of each constraint c_i , which controls tradeoffs when it is not possible to satisfy all constraints at once. Rewarding the Shannon entropy $H(p)$ encourages smoother distributions when it is possible to satisfy all constraints (Jaynes, 1957) and even when it is not. The hybrid objective (2) is historically known as the fuzzy maximum-entropy objective (Chen & Rosenfeld, 2000; Dudík et al., 2007) because it does not require the constraints to be satisfied exactly. Other reasonable variants are reviewed by Kazama & Tsujii (2005) and could be used here. Our innovation is to obtain the constraints from an LLM instead of from a data sample as in past work.

3.2 EXTRACTING COMMONSENSE FROM LLMs FOR PROBABILISTIC INFERENCE

We develop a concrete pipeline to build models as in §3.1 with LLMs as the knowledge source. In particular, the pipeline involves three stages of prompting: given a question Q , we identify (a) relevant variables and (b) pairs of interacting variables, allowing us to elicit (c) numerical constraints c . We can then (d) formulate a log-linear family of distributions p and optimize equation (2) over that family.

(a) Brainstorming Relevant Variables Given a question Q expressed in natural language, we prompt an LLM to brainstorm in free-form text, specifying the target variable y , the conditioning variables \mathbf{x} , and any additional variables \mathbf{z} by giving them names as well as a list of possible values $\mathbb{Y}, \mathbb{X}, \mathbb{Z}$ that they can take on.

Specifically, we prompt with the system message in Appendix A.1.1 followed by the single (1-shot) example in Appendix A.1.2. The example’s input is not from any of the domains we evaluate on; we obtained the example’s output by lightly editing the 0-shot output from a strong LLM (namely GPT-4o).

We then prompt the LLM to translate this free-form answer into a machine-readable JSON object, including variable definitions.

For evaluation purposes, we also supply in user prompt name of the target variable y , all its possible values defined in the dataset, and encourage the LLM to include it in its variables. However, to ensure that the target variable is always used exactly, we do not extract it during translation and instead add it into the list of variables.

⁷In our experiments, f_i will always be an indicator function as well, so the conditional expectations are simply conditional probabilities of our discrete random variables. However, allowing real-valued f_i would let us constrain the means, variances, and covariances of random variables. In the future, we might further broaden the constraint language. For example, one might ask the LLM about the differences or ratios of conditional expectations—“cats weigh less than dogs on average”—or the conditional entropy or mutual information of random variables. The LLM could also be asked for prediction intervals rather than point estimates, resulting in interval constraints.

(b) Choosing Quantities to Constrain We prompt the LLM to brainstorm interacting pairs of variables from stage (a), choose the best few pairs, and finally decide for each chosen pair $\{v_1, v_2\}$ whether to constrain $p(v_1 | v_2)$ or $p(v_2 | v_1)$. This prompt includes the brainstorming message from stage (a).

As before, we then prompt the LLM to translate this free-form list of conditional distributions into a JSON object. We then drop any z and x variables from the model that are not connected (directly or indirectly) to the target variable y , and thus drop conditional distributions mentioning those variables.

(c) Eliciting the Numerical Targets Now, for each surviving conditional distribution $p(v_1 | v_2)$, we ask the LLM to supply the numerical conditional probabilities. Specifically, for each $v_2 \in \mathbb{V}_2$, we prompt the LLM to generate a natural language query Q' for the distribution $p(v_1 | v_2 = v_2)$, and then prompt the LLM separately to return that distribution as a vector of dimension $|\mathbb{V}_1|$.

(In principle, we could constrain the distribution of v_1 for only certain proposed values $v_2 = v_2$. We leave this possibility to future work, along the possibility of eliciting conditional or joint probabilities involving more than 2 variables.)

Using the same method of generating natural language questions Q' , we prompt for the unary marginal distribution $p(v)$ for each variable v . We similarly prompt for the distribution $p(y | x = \mathbf{x})$, which corresponds to the original question Q (or a backed-off version of it, if some of the variables in \mathbf{x} were dropped).

(d) Optimizing a Log-linear Model We now choose a distribution p that approximately has the elicited conditional and marginal probabilities, by optimizing equation (2). Specifically, we define a log-linear family of models p_θ and optimize θ by batch gradient descent. The features of the log-linear model are all and only the indicator functions f_i and g_i that are necessary to express the list of unary and pairwise constraints (but not necessarily Q). The factor graph of this joint model contains only pairwise and unary potential functions that correspond to the proposed constraints.

We use brute force summation to exactly compute the conditional probabilities in equation (2).⁸ As for the weights w_i in equation (2), we use $w_i = c$ for some constant c to balance between constraint satisfaction and entropy smoothing.⁹ We empirically choose c on development data.

4 RELATED WORK

Large language models perform remarkably well on a diverse and challenging set of benchmarks (Ouyang et al., 2022; Anthropic Team, 2024; Gemini Team, 2024). Their effectiveness (Bubeck et al., 2023) is perhaps unsurprising, as they absorb vast amounts of world knowledge from their pretraining data (Petroni et al., 2019; Alkhamissi et al., 2022). On the other hand, their reasoning is brittle and is often based on shortcuts rather than sound inference rules (Saparov & He, 2023; ?; Dziri et al., 2023). Some studies suggest that learning sound reasoning from samples may be too challenging due to statistical shortcuts (Geirhos et al., 2020), even if a deep architecture like Transformer (?) can in principle implement it (Zhang et al., 2022). Many methods have thus been developed to extract better reasoning from LLMs in hopes of making better predictions with them. Within this direction, two ideas are immediately relevant to our work.

The first idea is using LLMs to brainstorm various pieces of relevant common knowledge about a question and then aggregating them to arrive at a prediction. Wang et al. (2023); Yao et al. (2023); Besta et al. (2024); Jung et al. (2022) all do so by aggregating over multiple reasoning paths. Viewed through the lens of brainstorming relevant knowledge and aggregation, our work introduces a new unit of common knowledge—that of a moment constraint on a probability distribution. We also propose a corresponding aggregation procedure of optimizing a shared underlying probabilistic model to agree with all the constraints.

⁸In our experiments, we instruct the LLM to propose at most 4 variables, and to select no more edges than variables, which makes this feasible. Scaling up to larger models will require approximate inference algorithms which may introduce additional sources of error.

⁹A more sophisticated option would be to place more weight on constraints where the LLM is more confident in the target value b_i . Another possibility would be to downweight constraints on variables and pairs of variables with many values, so that the objective function is not dominated by the many constraints that they yield.

Another related idea is to augment LLMs with formal reasoning components such as external symbolic reasoning engines and soft verifiers (Lyu et al., 2023; Xu et al., 2024; Pan et al., 2023; Bostrom et al., 2022; Ling et al., 2023). Our method can be viewed as augmenting LLMs with a formal reasoning engine that includes both fuzzy moment matching to infer the parameters of a graphical model and probabilistic inference to make predictions from the graphical model. While the cited works focus on improving the *logical* reasoning of LLMs, we study how to improve the *probabilistic* reasoning of LLMs.

Particularly worth mentioning is the maieutic prompting method of Jung et al. (2022), which takes inspiration from both lines of ideas—they brainstorm latent propositions by abductive reasoning, and then solve a joint constraint satisfaction problem to guess which propositions are true (and in particular, whether the original query Q is true). Their method can be viewed as performing *MAP* inference under a factor graph consisting of binary random variables corresponding to propositions, and with unary factors and binary factors whose parameters are extracted from LLMs and pretrained NLI models. They use a recursive algorithm to create an initial tree of propositions, and later add edges between all pairs of propositions. On the other hand, our method performs *marginal* inference over a factor graph of categorical variables corresponding to properties of situations in the world; our graph structure is directly proposed by an LLM and is usually sparser. The parameters of our graphical model are found by optimizing a set of LLM-proposed constraints on its various marginal distributions.

Probabilistic reasoning using LLMs has been relatively under-explored as a research problem. In a position paper, Dohan et al. (2022) propose to view prompted LLMs as conditional distributions over strings and the orchestration of LLM calls as a probabilistic program over strings (van de Meent et al., 2021). More recently, Nafar et al. (2024) use LLMs to generate probabilistic programs that get executed to produce distributions that answer probabilistic questions. However, crucially, their focus is more on abstract reasoning problems and requires as input the definition of a probabilistic model. Our work focuses on building that probabilistic model with the help of a LLM.

Researchers in Psychology and Cognitive Science have long explored the probability judgments in humans. Our work is also motivated by theories suggesting that a coherent probability judgment should be an accurate one. Osherson et al. (1994; 2001) proposed to extract from human intuitions a coherent distribution that reconciles a person’s different instances of probability judgments. More recently, Zhu & Griffiths (2024) showed that LLMs exhibit similar statistical properties in their probability judgments. However, despite the theoretical soundness, empirical results in this area have been mixed (Zhu et al., 2022), and there often is a lack of correlation between a coherent judgement and an accurate judgment.

5 EXPERIMENTS

We perform two experiments. §5.1 studies whether our model-building pipeline helps end-to-end performance in answering questions of the form introduced in §2. §5.2 tests the effectiveness of our two prompting stages (§3.2), by measuring the effect of intervening on their results in various ways.¹⁰ All of our experiments use the following setup.

Task As described in §2, the task is to provide an estimate \hat{p} (a normalized vector of size $|\mathbb{Y}|$) to a probability distribution $p(y \mid \mathbf{x} \in \mathbb{S})$ described by a natural language question Q .

Metric To evaluate the quality of an estimate \hat{p} , we compute its Total Variation Distance from a reference distribution p ,

$$\text{TVD}(p, \hat{p}) = \frac{1}{2} \sum_{y \in \mathbb{Y}} |\hat{p}(y) - p(y)| \quad (3)$$

Datasets To evaluate our system, we need questions Q paired with reference distributions p . To do so, we derived questions from three publicly available tabular datasets spanning domains including short term rentals (Inside Airbnb), daily activities rental (American Time-Use Survey), and personal

¹⁰This may be reminiscent of interventional studies on internal activations of neural networks (mechanistic interpretability).

attitudes (World Values Survey). We first describe the datasets briefly, then how we generate a set of questions given the contents of the dataset.

The Inside Airbnb¹¹ dataset (**AIR**) is a publicly available dataset of property rental listings across cities in the United States during 2023. Data for a city is collected by Inside Airbnb if its part of a list of major cities, or upon community request. Among the available cities, we randomly sample six cities to use in our evaluation, plus one more for tuning prompts and hyper-parameters.

The American Time-Use Survey¹² (**ATUS**) is a publicly available census dataset that collects meta-data about how people in the United States spend their time over the course of the week. The data is published yearly, and we choose data from years 2018, 2020, 2022 for evaluation, while using 2023 data for development.

The World Values Survey¹³ (**WVS**) is a survey dataset that collects demographic data about individuals in various countries and their responses to questions that probe their values. We randomly sample six countries for evaluation, and hold out this domain entirely for evaluation.

More details on the three datasets and their pre-processing is discussed in Appendix B.

Question Generation We randomly sample formal probability queries with n conditions based on the schema of the datasets, and translate them to natural language with the help of a LLM (we generate natural language questions given a formal query, and manually fix any errors). Specifically, for each dataset, and each $n \in \{0, 1, 2\}$, we first generate the set of all possible queries of the form $p(y \mid x_1 = x_1, \dots, x_n = x_n)$, and then filter it down by requiring that at least one of the conditions changes the distribution over the target variable y by ≥ 0.05 in terms of total variation distance. Then we sample 6 questions uniformly from this set. For comparability, the questions for a given dataset and n are reused across all values of the split variable (city for AIR, year for ATUS, or country for WVS), with the question being additionally conditioned on this value. We refer to these as the Main questions.

For AIR and ATUS, we also generate a Focus set of questions by repeating the same sampling process described above, except with an additional filter that the target variable y must be Price or Activity, respectively. This provides a set of questions that is more focused. We chose Activity and Price because they potentially interact with many other random variables from their respective domains.

LLM Calls Unless otherwise noted, we use GPT-4o-mini as the LLM in our experiments. All LLM calls are made at temperature 0.2, with a max token of 4096 (the default in LangChain OpenAI).

5.1 END-TO-END EVALUATION

We evaluate our pipeline end-to-end on the World Values Survey (WVS), which was not used to develop the pipeline. For completeness, we also evaluate on the held-out subsets of Inside Airbnb (AIR) and American Time-Use Survey (ATUS).

Direct Prompting We compare against the obvious baseline of simply asking the LLM to answer Q , using a chain-of-thought prompt (“zero-shot CoT”) at temperature 0.2. To ensure that the baseline enjoys a comparable amount of computation time, we actually call the LLM many times and average the resulting distributions \hat{p} . The number of calls is chosen to match the average number of calls made for extracting moment constraints in stage (b) of our pipeline.

Restricted Variables We also report the performance of our pipeline when we prompt it to use only variables in the dataset’s schema (see §5.2 below for details).

The results are given in Table 1. Figure 2 breaks them down by the number of conditions $|x|$ specified in the question. Figure 4 in the appendices compares TVD of our method to the baseline on each question separately, using a scatterplot.

¹¹<https://insideairbnb.com/>

¹²<https://www.bls.gov/tus/>

¹³<https://www.worldvaluessurvey.org/wvs.jsp>

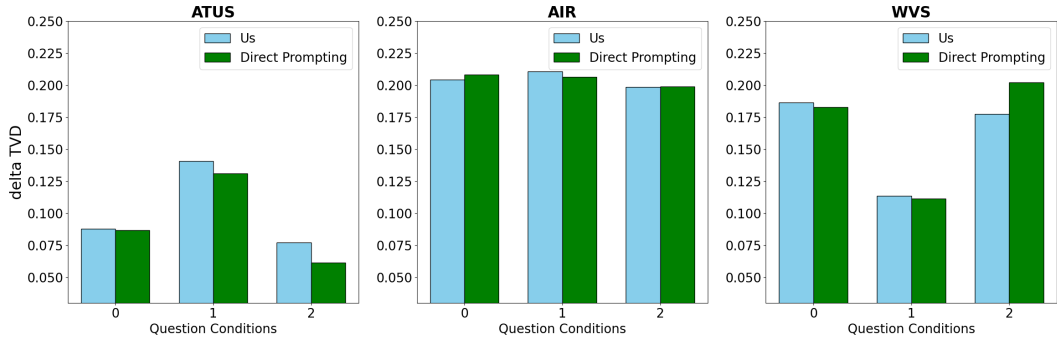


Figure 2: Breakdown of the end-to-end evaluation (§5.1) by number of conditions in the question.

Table 1: Average total variation distance against dataset distribution over questions as well as splits for subsets of questions Main and Focus respectively. We boldface the best result in each column along with all results that are not significantly worse (paired permutation test, $p < 0.05$).

	Main			Focus		
	ATUS	AIR	WVS	ATUS	AIR	WVS
Direct (§5.1)	0.094	0.204	0.166	0.123	0.175	—
Ours (§3)	0.099	0.203	0.163	0.123	0.180	—
Ours, restricted vars (§5.1)	0.105	0.190	—	0.182	0.164	—

Discussion of Results Unfortunately, constructing and querying an *ad hoc* model was not more accurate than simply asking the LLM. The target questions Q that we derived from these datasets were arguably too easy for our rather powerful LLM, GPT-4o-mini. The baseline system was already able to answer them with rather low TVD.

As a consolation, at least our method did not hurt. There are many ways that it could have gone wrong: after all, we were using natural language to obtain many imperfect numeric constraints and feeding them into a joint optimization problem. We had feared that the compounded noise in this process might swamp the signal. However, in practice the elicited constraints on both Q and other conditional probabilities tended to be rather accurate in this domain.¹⁴ Respecting these additional constraints simply did not change the answer much, either for better or for worse (see Figure 4).

Thus, an optimistic interpretation of the results is that our approach is viable, but that we would need to construct more difficult guesstimation problems or commonsense reasoning problems to show its value. Our approach will only help on problems where the LLM does not know how to answer the target question Q , but does know how to identify and answer other questions whose answers jointly imply an answer to Q .

We also discuss possible improvements to our method in §6, which might help on such a domain or on the current domain.

5.2 INTERVENTION EXPERIMENTS

We wish to study whether our method finds useful latent variables,¹⁵ whether stage (b)’s proposed directions are helpful, and whether the elicited numeric constraints are accurate. This leads to the following set of interventions:

1. Randomly replacing a latent variable z with a different one after stages (a) and (b). This affects the natural-language questions that we ask at stage (c).

¹⁴We assessed them during pilot experiments on AIR to have an average TVD of 0.11. However, those results used the stronger GPT-4o model; we will add a formal evaluation using GPT-4o-mini.

¹⁵This requires stage (a) to propose the variables and also requires stage (b) not to discard them (see §3.2).

2. Randomly reverting the direction of the query $v_1 \mid v_2$ to be $v_2 \mid v_1$ after stage (b). Again, this affects the questions at stage (c).
3. Interpolating each elicited numeric constraint after stage (c) with the oracle value computed from the dataset.

For all intervention experiments, we omit the constraint on $p(y \mid \mathbf{x} = \mathbf{x})$, which corresponds to the original question Q . This constraint often has so much influence on the final result that it would mask the effect of the intervention.

Intervention 3 is possible only when the proposed variables appear as fields in the dataset so that we can get oracle values. Therefore, in that experiment—for both intervention 3 and its control condition—we modify the prompt of stage (a) to include the dataset schema (variable names along with their possible values) and to instruct the LLM to confine its brainstorming to these options.

We also use this modified prompt for intervention 1 and its control condition. This ensures a controlled comparison: it asks whether the LLM chooses wisely from among the schema variables, compared to the random choice of schema variable made by intervention 1. With the original prompt, the difference in performance might only reflect whether schema variables are more or less useful than non-schema variables.

1 and 2 are ablations that we expect to hurt performance. For 1, we randomly choose $i \in \{0, 1, 2\}$ number of variables that is not the target or the condition, and substitute uniformly from variables from the schema that’s not already included. For 2, we randomly chose $j \in \{0, 1, 2, 3\}$ pairwise constraints to flip the direction. For both 1 and 2, since not all graphs have enough variables / edges that can be intervened on, we restrict our analysis to the subset of questions where the proposed model supports interventions of 3 node substitutions and 3 pairwise constraint reversals. (See Figure 3, columns 1 and 2.)

3 is an oracle intervention that we expect to help performance. We mix proposed distributions in stage (b) with the oracle distribution computed from the dataset at weight $w \in [0, 0.2, 0.4, 0.6, 0.8, 1.0]$, where $w = 0$ corresponds to no intervention, and $w = 1.0$ corresponds to using oracle numeric constraints. We also tried *hurting* performance by substituting a random distribution for the oracle distribution, drawing it uniformly from the simplex of probability distributions. (See Figure 3, columns 3 and 4.)

Discussion of Results Columns 1 and 2 of Figure 3 suggest that perturbing the selection of variables or the direction of the conditional probabilities did not significantly affect the average gap between our method and the baseline. In other words, the LLM may not have made the best choices at these steps, despite our prompts.

Column 3 of Figure 3 provides a sanity check that as our constraints move towards the oracle, the error moves to 0. Unfortunately, this plot alone does not tease apart the contributions of moving the brainstorming queries produced by stages (a) and (b) towards oracle and moving the query corresponding to Q towards the oracle. Even though we don’t explicitly add it in the intervention experiments, stages (a) and (b) often propose a query corresponding to the question Q by themselves. This suggests additional studies to separate the effect of a good answer Q during brainstorming, and the effect of good answers to *other* related queries. Fortunately, column 4 of Figure 3 shows that artificial IID noise does not hurt our method by more than it hurts the direct-prompt baseline.

6 FUTURE WORK

Further prompt engineering might potentially help our system find crucial combinations of constraints that would improve on the baseline system. We cannot rule out the possibility that such constraints existed in our experiments and we simply failed to find them; we could use brute force exploration to check if they exist.

Stage (d) of our pipeline (§3.2) adds constraints to our model, but at the same time it expands the model family by creating additional parameters to help satisfy those constraints. As this may lead to overfitting, it might be wise to regularize our model objective (2) beyond the entropy term $H(p)$.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

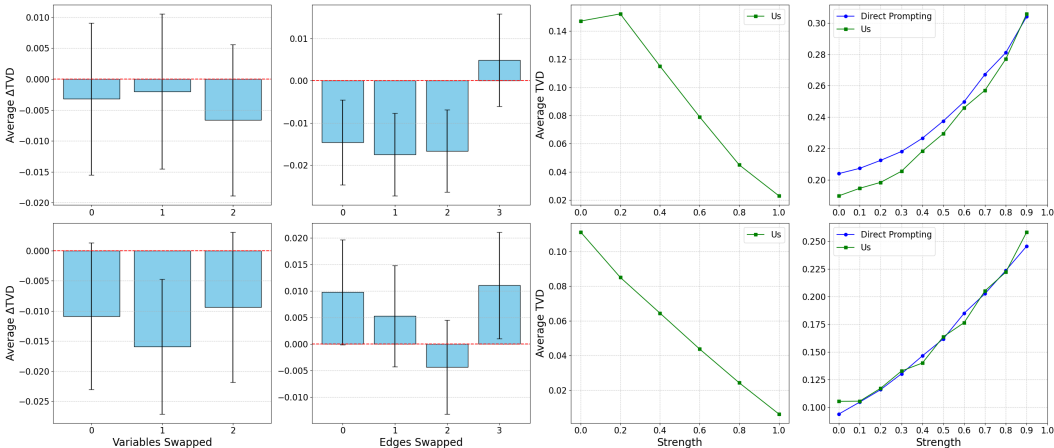


Figure 3: Results of intervention experiments (§5.1). “Us” in this figure refers to our approach. Top row corresponds to results on the Main set of questions on AIR domain, bottom row corresponds to the Main set of questions on ATUS domain. **Columns 1 and 2** visualize results of interventions 1 and 2, which randomly replaces zero to two latent variables with a different one after stages (a) and (b) of §3.2, and randomly reverses the direction of zero to three queries $v_1 \mid v_2 \mid v_2 \mid v_1$ after stage (b), respectively. Their x -axes denote the number of intervened nodes/queries, and their y -axes denote the average $TVD(p, \hat{p}_{us}) - TVD(p, \hat{p}_{direct\ prompt})$. The error bars denote one standard deviation of the average. Columns 3 and 4 correspond to intervention 3-oracle and intervention 3-noise. Their x -axes are the interpolation coefficient, and their y -axes are $TVD(p, \hat{p}_{us})$.

The LLM could also provide more precise information about how to penalize deviations from each constraint c_i , for example by providing a weight w_i , an interval on the target b_i , or a full loss function. The objective (2) could also be extended by asking the model p to satisfy other kinds of constraints extracted from the LLM, such as relative probabilities (see footnote 7).

For simplicity, our implementation focused on models with a small number of categorical variables and only unary and binary factors. Future work should extend this to continuous variables as well as larger models, which may require approximate inference algorithms such as belief propagation and expectation propagation.

Our method builds an *ad hoc* model p_θ that can answer the original question Q , but p_θ can be interrogated further with additional probabilistic queries about its variables. Answers to those questions may be useful for interpreting the answer to the original question Q , and they may be compared against reference distributions computed from datasets to further assess the model.

Furthermore, p_θ can identify likely situations and marginally likely values for y and z . In principle, those could be fed back into a second round of brainstorming to further refine the model in high-probability regions of the outcome space—for example by introducing new latent variables or adjusting the granularity of existing variables.

We primarily used GPT-4o-mini for our experiments due to limited budgets. However, most LLM calls are spent on eliciting numerical targets in stage (c), we can use more powerful LLMs for stage (a) and (b), which can potentially improve the design of the *ad hoc* model.

Finally, future work should investigate when to trust the LLM. Confidence estimation could be used to upweight more accurate constraints in the optimization objective. In some cases, the LLM estimates might be improved (calibrated) with a small amount of supervised training data. For example, we might discover that the LLM tends to overestimate certain kinds of probabilities, and attempt to automatically correct these.

540 ACKNOWLEDGMENTS

541
542 We thank X for the quote from Jaynes (2003), and for their helpful discussions.
543

544 REFERENCES

- 545 Anthropic Team. The claude 3 model family: Opus, sonnet, haiku | papers with code, 2024. URL
546 <https://paperswithcode.com/paper/the-claude-3-model-family-opus-sonnet-haiku>.
547
548
- 549 Gemini Team. Gemini: A family of highly capable multimodal models, 2024.
550
- 551 Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. A review
552 on language models as knowledge bases, 2022.
553
- 554 Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi,
555 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoeffler. Graph
556 of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI*
557 *Conference on Artificial Intelligence*, 38(16):17682–17690, Mar. 2024. doi: 10.1609/aaai.v38i16.
558 29720. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29720>.
- 559 Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, and Greg Durrett. Natural language deduction
560 through search over statement compositions. In Yoav Goldberg, Zornitsa Kozareva, and Yue
561 Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp.
562 4871–4883, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
563 Linguistics. doi: 10.18653/v1/2022.findings-emnlp.358. URL [https://aclanthology.org/](https://aclanthology.org/2022.findings-emnlp.358)
564 [2022.findings-emnlp.358](https://aclanthology.org/2022.findings-emnlp.358).
- 565 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar,
566 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio
567 Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4,
568 2023.
569
- 570 Stanley F. Chen and Ronald Rosenfeld. A survey of smoothing techniques for ME models. *IEEE*
571 *Trans. on Speech and Audio Processing*, 8(1):37–50, January 2000.
572
- 573 David Dohan, Winnie Xu, Aitor Lewkowycz, Jacob Austin, David Bieber, Raphael Gontijo Lopes,
574 Yuhuai Wu, Henryk Michalewski, Rif A. Saurous, Jascha Sohl-dickstein, Kevin Murphy, and
575 Charles Sutton. Language model cascades, 2022.
- 576 Miroslav Dudík, Steven J. Phillips, and Robert E. Schapire. Maximum entropy density estimation
577 with generalized regularization and an application to species distribution modeling. *Journal*
578 *of Machine Learning Research*, 8(44):1217–1260, 2007. URL [http://jmlr.org/papers/v8/](http://jmlr.org/papers/v8/dudik07a.html)
579 [dudik07a.html](http://jmlr.org/papers/v8/dudik07a.html).
580
- 581 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean
582 Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, Soumya Sanyal, Xiang
583 Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. Faith and fate: Limits of transformers on
584 compositionality. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
585 URL <https://openreview.net/forum?id=Fkckkr3ya8>.
- 586 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
587 Bethge, and Felix A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine*
588 *Intelligence*, 2(11):665–673, November 2020. ISSN 2522-5839. doi: 10.1038/s42256-020-00257-z.
589 URL <http://dx.doi.org/10.1038/s42256-020-00257-z>.
590
- 591 E. T. Jaynes. Information theory and statistical mechanics. *Phys. Rev.*, 106:620–630, May 1957. doi:
592 10.1103/PhysRev.106.620. URL <https://link.aps.org/doi/10.1103/PhysRev.106.620>.
593
- E. T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, 2003.

- 594 Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brahman, Chandra Bhagavatula, Ronan Le Bras, and
595 Yejin Choi. Maieutic prompting: Logically consistent reasoning with recursive explanations. In
596 Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on*
597 *Empirical Methods in Natural Language Processing*, pp. 1266–1279, Abu Dhabi, United Arab
598 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
599 emnlp-main.82. URL <https://aclanthology.org/2022.emnlp-main.82>.
- 600 Jun’ichi Kazama and Jun’ichi Tsujii. Maximum entropy models with inequality constraints: A case
601 study on text categorization. *Machine Learning*, 60:159–194, 2005.
- 602 Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su.
603 Deductive verification of chain-of-thought reasoning. In *Thirty-seventh Conference on Neural*
604 *Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=I5rsM4CY2z>.
- 605 Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki,
606 and Chris Callison-Burch. Faithful chain-of-thought reasoning, 2023.
- 607 Thomas P. Minka. Expectation propagation for approximate bayesian inference, 2013. URL
608 <https://arxiv.org/abs/1301.2294>.
- 609 Aliakbar Nafar, Kristen Brent Venable, and Parisa Kordjamshidi. Probabilistic reasoning in generative
610 large language models, 2024.
- 611 Daniel Osherson, Eldar Shafir, and Edward E. Smith. Extracting the coherent core of human
612 probability judgement: a research program for cognitive psychology. *Cognition*, 50(1):299–
613 313, 1994. ISSN 0010-0277. doi: [https://doi.org/10.1016/0010-0277\(94\)90033-7](https://doi.org/10.1016/0010-0277(94)90033-7). URL <https://www.sciencedirect.com/science/article/pii/0010027794900337>.
- 614 Daniel Osherson, David Lane, Peter Hartley, and Richard R. Batsell. Coherent probability from
615 incoherent judgment, 2001. URL <https://pubmed.ncbi.nlm.nih.gov/11577618/>.
- 616 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
617 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,
618 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and
619 Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- 620 Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. Logic-lm: Empowering large
621 language models with symbolic solvers for faithful logical reasoning, 2023.
- 622 Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller,
623 and Sebastian Riedel. Language models as knowledge bases?, 2019.
- 624 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis
625 of chain-of-thought. In *The Eleventh International Conference on Learning Representations*, 2023.
626 URL <https://openreview.net/forum?id=qFVVBzXxR2V>.
- 627 Jan-Willem van de Meent, Brooks Paige, Hongseok Yang, and Frank Wood. An introduction to
628 probabilistic programming, 2021.
- 629 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha
630 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
631 models. In *The Eleventh International Conference on Learning Representations*, 2023. URL
632 <https://openreview.net/forum?id=1PL1NIMMrw>.
- 633 Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical
634 reasoning via symbolic chain-of-thought, 2024.
- 635 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R
636 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In
637 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=5Xc1ecx01h>.
- 638 Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang, and Guy Van den Broeck. On the
639 paradox of learning to reason from data, 2022.

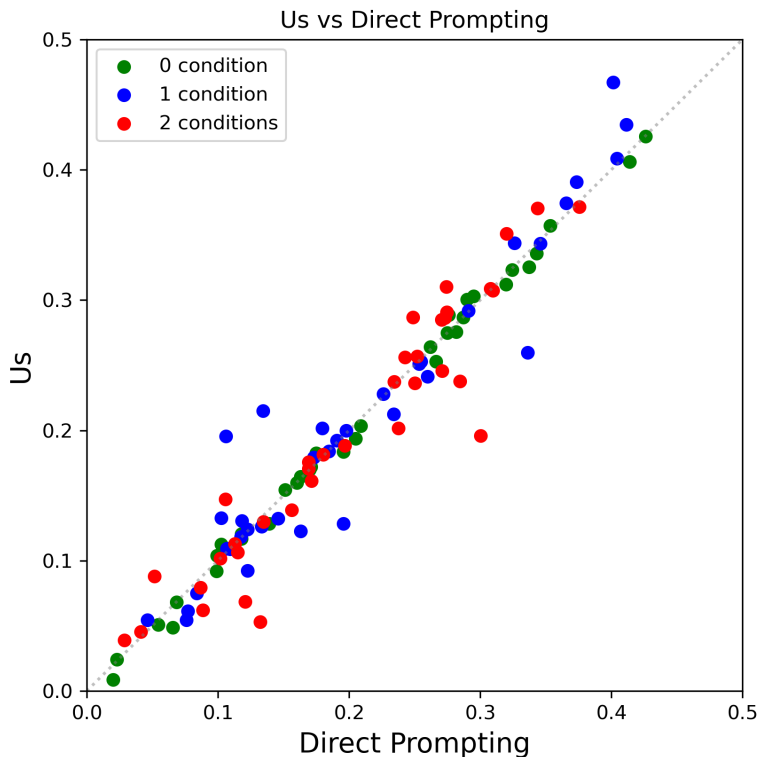


Figure 4: Scatterplot of the total variation distance against reference, Us versus Direct Prompt, on the Main set of questions for Inside Airbnb. Each point in the plot corresponds to a question from Main on a particular evaluation split (one of Ashville, Austin, Chicago, New Orleans, Pacific Grove, and Rhode Island), averaged over three random executions at temperature 0.2. The color of a point denote the number of conditions in the question. The other domains (ATUS and WVS) and the other set of questions (Focus) show a similar pattern in their scatterplots (not shown here).

Jian-Qiao Zhu and Thomas L. Griffiths. Incoherent probability judgments in large language models, 2024. URL <https://arxiv.org/abs/2401.16646>.

Jian-Qiao Zhu, Philip W.S. Newall, Joakim Sundh, Nick Chater, and Adam N. Sanborn. Clarifying the relationship between coherence and accuracy in probability judgments. *Cognition*, 223: 105022, 2022. ISSN 0010-0277. doi: <https://doi.org/10.1016/j.cognition.2022.105022>. URL <https://www.sciencedirect.com/science/article/pii/S0010027722000105>.

A APPENDIX

A.1 PROMPTS

A.1.1 VARIABLE PROPOSAL

SYSTEM

You are a data scientist.

You must design a graphical model to estimate conditional probabilities in a certain domain. The domain and the requested probabilities will be specified informally in REQUEST, so you must formalize the REQUEST into an Outcomes space, and defining categorical random Variables with mutually exclusive Values on the Outcomes space. The Outcomes space is a set of tuples that

702 represents all the possible Values combinations. DOMAIN, on the other hand, is one concise sentence
 703 that summarizes the entire Outcomes space. DOMAIN is thus a succinct summary of the population
 704 of the model, to provide to your colleagues. For simplicity, DOMAIN should omit information that
 705 the colleagues would assume by default. In other words, DOMAIN includes simple background
 706 information, and only some Variables whose sets of Values on the Outcomes are non-trivial.

707 The Variables must be sufficient to answer the REQUEST and contain all possible Values on the
 708 Outcomes. That is, you must design your model so that the REQUEST can be formalized as $P(X_0 |$
 709 $X_1 \in x_1, \dots, X_k \in x_k)$, where each X_i is one of your model Variables and x_i is some subset of
 710 its defined Values.

711 Although the REQUEST may specify Values of the Variables, the size of the model's Outcomes space
 712 is up to you to decide, i.e. you may define Values not specified in the REQUEST. This may allow us
 713 to use the model for other similar REQUESTs. In other words, the model may have an Outcomes
 714 space larger than implied in the REQUEST, on which the Variables can take Values that are not
 715 mentioned by REQUEST. If you choose a larger Outcomes space, make sure that the REQUEST can
 716 still be formalized exactly in the model.

717 You may also include additional Variables not mentioned by the REQUEST. These Variables can be
 718 very useful in mediating the relationships between X_0 and X_1, \dots, X_k . However, they should be
 719 concrete and unambiguous.

721 There can be at most 4 Variables in total, each with at most 5 Values. However, fewer Values
 722 enables more robust estimations. You can reduce the number of Values by grouping them into coarser
 723 categories, while still having enough granularity to express the REQUEST.

724 For simplicity, do not model missing data or edge cases, i.e. do not include Values like 'Unspecified'
 725 or 'Other'.

726 ****Keep the Values categories concrete and precise. DO NOT USE Values like 'Other'.****

727 For Variables that are numeric in nature, use precise ranges.

729 Variables with one Value will not count toward the limit of 4, but still must be mentioned in the
 730 DOMAIN.

731 Some Variables may be supplied to you. Include them if appropriate.

733 Discuss the design details of your model. First identify the potential Variables in the REQUEST, and
 734 whether it will help to include other Values. Again, you are free to define Values not in the REQUEST
 735 so we can reuse the model, or fit the model with a larger dataset.

736 From this, decide the appropriate Outcomes size, i.e. whether a larger Outcomes will be helpful.
 737 Afterwards, define the Variables (including any additional Variables not in REQUEST) accordingly,
 738 with all possible Values on the Outcomes. Then, derive the DOMAIN from discussing each Variable's
 739 Values, and whether the Values set needs to be mentioned.

740 In the end, state the Variables, the DOMAIN in one sentence, and the formalized REQUEST.

743 A.1.2 ONE-SHOT EXAMPLE FOR VARIABLE PROPOSAL

744
 745 USER

746 REQUEST: For a horror movie released in United States between 1980 and 2000 rated as 'R', what
 747 is its probability distribution of total profit?

748 Supplied Variables:

749 Profit Range: 'Less than \$10M'; '\$10M - \$50M'; '\$50M - \$100M'; 'Over \$100M'

750 ASSISTANT

751 ### Design Details

752 ### Potential Variables in the REQUEST:

753 1. ****Genre****: The genre is specified as horror.

756 2. **Country of Release**: This is specified as the United States.
757
758 3. **Release Period**: This is specified as between 1980 and 2000.
759 4. **Rating**: This is specified as 'R'.
760 5. **Profit Range**: Supplied with specific categories.
761
762 **### Additional Variables:**
763 1. **Production Budget**: This could be an important variable to include, as it can significantly
764 affect the total profit.
765 2. **Critical Acclaim**: This is subjective and hard to obtain data for, so it will not be included.
766
767 **### Deciding on the Outcomes Space**
768 To make the model more reusable and provide more robust estimations, I will include other categories
769 for the rating. This will allow us to use the model for other similar requests in the future. We also add
770 production budget as an additional Variable.
771
772 **### Variables and Their Values**
773 1. **Genre**:
774 - Values: 'Horror' (only horror movies are of interest, fixed and does not count toward the limit of 4,
775 mention in DOMAIN).
776 2. **Country of Release**:
777 - Values: 'United States' (only US releases are of interest, fixed and does not count toward the limit
778 of 4, mention in DOMAIN).
779 3. **Release Period**:
780 - Values: '1980-1989'; '1990-2000' (to allow for some granularity in the release years, can be
781 summarized as 1980 - 2000, mention in DOMAIN).
782 4. **Rating**:
783 - Values: 'G'; 'PG'; 'PG-13'; 'R'; 'NC-17' (enlarging the rating categories to make the model more
784 reusable, summarized as all ratings, can be left out of DOMAIN).
785 5. **Profit Range**:
786 - Values: 'Less than \$10M'; '\$10M - \$50M'; '\$50M - \$100M'; 'Over \$100M' (as supplied, all profit
787 ranges, can be left out of DOMAIN).
788 6. **Production Budget**:
789 - Values: 'Less than \$10M'; '\$10M - \$50M'; 'Over \$50M' (to provide granularity in the budget,
790 summarized as all production ranges, can be left out of DOMAIN).
791
792 **### DOMAIN**
793 We need to mention the Horror genre, release in United States, released between 1980 - 2000. On
794 the other hand, Rating (expanded), Profit Range, and Production Budget have covered all possible
795 Values, and are rather trivial to mention. We summarize the model population as:
796
797 'Horror movies released in the United States between 1980 and 2000'
798
799 **### Formalized REQUEST**
800 The formalized REQUEST is: $P(\text{Profit Range} \mid \text{Genre} = \text{'Horror'}, \text{Country of Release} = \text{'United States'}, \text{Release Period} \in \{ \text{'1980-1989'}, \text{'1990-2000'} \}, \text{Rating} = \text{'R'})$.
801
802 **### Summary**
803 ****Variables:****
804 1. **Genre**:
805
806
807
808
809

810 - Values: 'Horror'
811
812 2. **Country of Release**:
813 - Values: 'United States'
814
815 3. **Release Period**:
816 - Values: '1980-1989'; '1990-2000'
817
818 4. **Rating**:
819 - Values: 'G'; 'PG'; 'PG-13'; 'R'; 'NC-17'
820
821 5. **Profit Range**:
822 - Values: 'Less than \$10M'; '\$10M - \$50M'; '\$50M - \$100M'; 'Over \$100M'
823
824 6. **Production Budget**:
825 - Values: 'Less than \$10M'; '\$10M - \$50M'; 'Over \$50M'
826 **DOMAIN**: 'Horror movies released in the United States between 1980 and 2000.'
827 **Formalized REQUEST**: $P(\text{Profit Range} \mid \text{Genre} = \text{'Horror'}, \text{Country of Release} = \text{'United States'},$
828 $\text{Release Period} \in \{ \text{'1980-1989'}, \text{'1990-2000'} \}, \text{Rating} = \text{'R'})$
829
830 A.1.3 VARIABLE PROPOSAL EXAMPLE DIALOGUE
831
832 USER
833
834 REQUEST: For an Airbnb listing in Austin, TX in 2023 with a price between \$201 to \$500 and is
835 listed to accommodate a maximum of 4 or more people, what is the probability distribution of room
836 type?
837 Supplied Variables:
838 Room Type: 'entire home or apartment'; 'private room'; 'shared or hotel room'
839
840 ASSISTANT
841 **Design Details**
842 **Potential Variables in the REQUEST**:
843
844 1. **Location**: This is specified as Austin, TX.
845
846 2. **Year**: This is specified as 2023.
847
848 3. **Price Range**: This is specified as between \$201 and \$500.
849
850 4. **Maximum Occupancy**: This is specified as accommodating 4 or more people.
851
852 5. **Room Type**: This is supplied with specific categories.
853
854 **Additional Variables**:
855
856 1. **Property Type**: This could be useful to include as it may influence the room type (e.g., house,
857 apartment, condo).
858
859 2. **Amenities**: This could also be relevant, but it is too broad and subjective for this model.
860
861 **Deciding on the Outcomes Space**
862
863 To make the model more reusable, I will include additional categories for the property type. This
will allow us to use the model for other similar requests in the future. The price range and maximum
occupancy will also be defined clearly.
Variables and Their Values
1. **Location**:
- Values: 'Austin, TX' (fix to Austin, does not count toward the limit of 4, mention in DOMAIN).

- 864 2. **Year**:
- 865 - Values: '2023' (fix to 2023, does not count toward the limit of 4, mention in DOMAIN).
- 866
- 867 3. **Price Range**:
- 868 - Values: '\$201 - \$300'; '\$301 - \$400'; '\$401 - \$500' (providing granularity within the specified
- 869 range).
- 870
- 871 4. **Maximum Occupancy**:
- 872 - Values: '1-3 people'; '4 or more people' (to capture the specified condition).
- 873
- 874 5. **Room Type**:
- 875 - Values: 'Entire home or apartment'; 'Private room'; 'Shared or hotel room' (as supplied).
- 876
- 877 6. **Property Type**:
- 878 - Values: 'House'; 'Apartment'; 'Condo'; 'Townhouse'; 'Other' (to provide additional context for the
- 879 listing).

880 ### DOMAIN

881 We need to mention the location (Austin, TX), the year (2023), and the price range (between \$201

882 and \$500). The maximum occupancy condition is also relevant but can be summarized in the context

883 of the room type. The room type and property type have covered all possible values and are rather

884 trivial to mention. We summarize the model population as:

885

886 "Airbnb listings in Austin, TX in 2023 with a price between \$201 and \$500."

887 ### Formalized REQUEST

888 The formalized REQUEST is: P(Room Type | Location = 'Austin, TX', Year = '2023', Price Range

889 \in {'\$201 - \$300', '\$301 - \$400', '\$401 - \$500'}, Maximum Occupancy = '4 or more people').

890

891 ### Summary

892 **Variables**:

893 1. **Location**:

894 - Values: 'Austin, TX'

895

896 2. **Year**:

897 - Values: '2023'

898

899 3. **Price Range**:

900 - Values: '\$201 - \$300'; '\$301 - \$400'; '\$401 - \$500'

901

902 4. **Maximum Occupancy**:

903 - Values: '1-3 people'; '4 or more people'

904

905 5. **Room Type**:

906 - Values: 'Entire home or apartment'; 'Private room'; 'Shared or hotel room'

907

908 6. **Property Type**:

909 - Values: 'House'; 'Apartment'; 'Condo'; 'Townhouse'; 'Other'

910

911 **DOMAIN**: "Airbnb listings in Austin, TX in 2023 with a price between \$201 and \$500."

912 **Formalized REQUEST**: P(Room Type | Location = 'Austin, TX', Year = '2023', Price Range \in

913 {'\$201 - \$300', '\$301 - \$400', '\$401 - \$500'}, Maximum Occupancy = '4 or more people').

914

915 A.1.4 VARIABLE PROPOSAL EXAMPLE TRANSLATION

916 SYSTEM You will be given the transcript of a statistician's model designing process. Transcribe the

917 final Variables and DOMAIN in the provided message into JSON using a provided schema. You can

918 find the Variables and DOMAIN toward the end of the message in a summary. Do not extract Room
919 Type.
920
921 USER
922 ### Design Details
923 ##### Potential Variables in the REQUEST:
924
925 1. **Location**: This is specified as Austin, TX.
926 2. **Year**: This is specified as 2023.
927 3. **Price Range**: This is specified as between \$201 and \$500.
928 4. **Maximum Occupancy**: This is specified as accommodating 4 or more people.
929 5. **Room Type**: This is supplied with specific categories.
930
931 ##### Additional Variables:
932
933 1. **Property Type**: This could be useful to include as it may influence the room type (e.g., house,
934 apartment, condo).
935 2. **Amenities**: This could also be relevant, but it is too broad and subjective for this model.
936
937 ### Deciding on the Outcomes Space
938
939 To make the model more reusable, I will include additional categories for the property type. This
940 will allow us to use the model for other similar requests in the future. The price range and maximum
941 occupancy will also be defined clearly.
942
943 ### Variables and Their Values
944
945 1. **Location**:
946 - Values: 'Austin, TX' (fix to Austin, does not count toward the limit of 4, mention in DOMAIN).
947 2. **Year**:
948 - Values: '2023' (fix to 2023, does not count toward the limit of 4, mention in DOMAIN).
949 3. **Price Range**:
950 - Values: '\$201 - \$300'; '\$301 - \$400'; '\$401 - \$500' (providing granularity within the specified
951 range).
952 4. **Maximum Occupancy**:
953 - Values: '1-3 people'; '4 or more people' (to capture the specified condition).
954 5. **Room Type**:
955 - Values: 'Entire home or apartment'; 'Private room'; 'Shared or hotel room' (as supplied).
956 6. **Property Type**:
957 - Values: 'House'; 'Apartment'; 'Condo'; 'Townhouse'; 'Other' (to provide additional context for the
958 listing).
959
960 ##### DOMAIN
961
962 We need to mention the location (Austin, TX), the year (2023), and the price range (between \$201
963 and \$500). The maximum occupancy condition is also relevant but can be summarized in the context
964 of the room type. The room type and property type have covered all possible values and are rather
965 trivial to mention. We summarize the model population as:
966
967 "Airbnb listings in Austin, TX in 2023 with a price between \$201 and \$500."
968
969 ### Formalized REQUEST
970
971 The formalized REQUEST is: $P(\text{Room Type} \mid \text{Location} = \text{'Austin, TX'}, \text{Year} = \text{'2023'}, \text{Price Range} \in \{\text{'$201 - $300'}, \text{'$301 - $400'}, \text{'$401 - $500'}\}, \text{Maximum Occupancy} = \text{'4 or more people'})$.

972 **### Summary**
973 ****Variables:****
974
975 1. ****Location****:
976 - Values: 'Austin, TX'
977
978 2. ****Year****:
979 - Values: '2023'
980
981 3. ****Price Range****:
982 - Values: '\$201 - \$300'; '\$301 - \$400'; '\$401 - \$500'
983
984 4. ****Maximum Occupancy****:
985 - Values: '1-3 people'; '4 or more people'
986
987 5. ****Room Type****:
988 - Values: 'Entire home or apartment'; 'Private room'; 'Shared or hotel room'
989
990 6. ****Property Type****:
991 - Values: 'House'; 'Apartment'; 'Condo'; 'Townhouse'; 'Other'
992 ****DOMAIN****: "Airbnb listings in Austin, TX in 2023 with a price between \$201 and \$500."
993 ****Formalized REQUEST****: $P(\text{Room Type} \mid \text{Location} = \text{'Austin, TX'}, \text{Year} = \text{'2023'}, \text{Price Range} \in \{\text{'$201 - $300'}, \text{'$301 - $400'}, \text{'$401 - $500'}\}, \text{Maximum Occupancy} = \text{'4 or more people'})$.
994
995 ASSISTANT
996
997 {"Variables": [{"Name": "Location", "Value": ["Austin, TX"]}, {"Name": "Year", "Value": ["2023"]},
998 {"Name": "Price Range", "Value": ["\$201 - \$300", "\$301 - \$400", "\$401 - \$500"]}, {"Name":
999 "Maximum Occupancy", "Value": ["1-3 people", "4 or more people"]}, {"Name": "Room Type",
1000 "Value": ["Entire home or apartment", "Private room", "Shared or hotel room"]}, {"Name": "Property
1001 Type", "Value": ["House", "Apartment", "Condo", "Townhouse", "Other"]}, {"Name": "DOMAIN": "Airbnb
1002 listings in Austin, TX in 2023 with a price between \$201 and \$500."}]
1003
1004
1005 A.1.5 INTERACTION PROPOSAL EXAMPLE DIALOGUE
1006
1007 SYSTEM
1008
1009 You are designing an undirected graphical model for DOMAIN The model will consist of only binary
1010 factors. You will be given a list of categorical variable definition with Target variable and potentially
1011 some Observed variables and Latent variables. Propose pairs of variables from this list that should be
1012 coupled by a factor in your model.
1013
1014 Your model should be good at predicting Target VARIABLE when observing Observed variables.
1015 Latent variables may help reasoning about different situation and contexts.
1016
1017 You are free to omit Observed or Latent variables that are not helpful for this prediction. Some Latent
1018 variables may exist only to aid prediction of other Observed or Latent variables, and not directly
1019 correlate with the Target variable.
1020
1021 For a potential factor between two variables, consider their relationship IN STRICT ISOLATION, i.e.
1022 WITHOUT referring to any other variable. To do this, when considering (A, B), ALWAYS explicitly
1023 state a question: does A vary significantly with different values of B? Then answer the question. If
1024 you are not very confident, it means they should not be connected.
1025
1026 Avoid focusing on only the Target variable. Aim for few but strong factors. Add factors sparingly.
1027 Specifically, with N variables, NEVER propose more than N factors. For the variables you do include
1028 in the model, the variables should form a connected graph.
1029
1030 Start by repeating the maximum number of factors allowed, then discuss the model design.

1026 After potential options, select a final set of few but strong factors. For each factor, consider whether
1027 you will be better at guessing the values of $p(A|B)$ or $p(B|A)$. If $P(A|B)$ is easier to guess, list them in
1028 an ordered pair (A, B) to represent (Target, Condition). Otherwise, if $P(B|A)$ is easier to guess, give
1029 (B, A).

1030 In the end, state them a list, each of which is an ordered pair of 2 strings.

1031 USER

1032 **Target Variable:**

1033 Room Type: 'entire home or apartment'; 'private room'; 'shared or hotel room'

1034 **Observed Variable:**

1035 Maximum Occupancy: '1-3 people'; '4 or more people'

1036 **Latent Variable:**

1037 Price Range: '\$201 - \$300'; '\$301 - \$400'; '\$401 - \$500'

1038 Property Type: 'House'; 'Apartment'; 'Condo'; 'Townhouse'

1039 ASSISTANT

1040 The maximum number of factors allowed is equal to the number of variables in the model. In this
1041 case, we have 5 variables: 1 Target Variable, 1 Observed Variable, and 3 Latent Variables. Therefore,
1042 we can propose a maximum of 5 factors.

1043 ### Model Design Discussion

1044 1. **Room Type (Target Variable)** and **Maximum Occupancy (Observed Variable)**:

1045 - Question: Does Room Type vary significantly with different values of Maximum Occupancy?

1046 - Answer: Yes, different room types can accommodate different maximum occupancies. For example,
1047 an entire home or apartment typically accommodates more people than a shared or hotel room.

1048 2. **Room Type (Target Variable)** and **Price Range (Latent Variable)**:

1049 - Question: Does Room Type vary significantly with different values of Price Range?

1050 - Answer: Yes, different room types are often associated with different price ranges. For instance,
1051 entire homes or apartments may be priced higher than shared or hotel rooms.

1052 3. **Room Type (Target Variable)** and **Property Type (Latent Variable)**:

1053 - Question: Does Room Type vary significantly with different values of Property Type?

1054 - Answer: Yes, the type of property can influence the room type. For example, a house is more likely
1055 to be listed as an entire home or apartment compared to a condo.

1056 4. **Maximum Occupancy (Observed Variable)** and **Price Range (Latent Variable)**:

1057 - Question: Does Maximum Occupancy vary significantly with different values of Price Range?

1058 - Answer: Yes, properties that accommodate more people may also be priced differently, often higher.

1059 5. **Maximum Occupancy (Observed Variable)** and **Property Type (Latent Variable)**:

1060 - Question: Does Maximum Occupancy vary significantly with different values of Property Type?

1061 - Answer: Yes, different property types can have different maximum occupancies. For example, a
1062 townhouse may accommodate more people than a condo.

1063 ### Final Set of Factors

1064 After considering the relationships, I will select the following pairs of variables to form the factors in
1065 the model:

1066 1. (Room Type, Maximum Occupancy) - It is easier to guess the Room Type given the Maximum
1067 Occupancy.

- 1080 2. (Room Type, Price Range) - It is easier to guess the Room Type given the Price Range.
1081
1082 3. (Room Type, Property Type) - It is easier to guess the Room Type given the Property Type.
1083
1084 4. (Maximum Occupancy, Price Range) - It is easier to guess the Maximum Occupancy given the
1085 Price Range.
1086
1087 5. (Maximum Occupancy, Property Type) - It is easier to guess the Maximum Occupancy given the
1088 Property Type.

1087 ### List of Factors

- 1089 1. (Maximum Occupancy, Room Type)
1090
1091 2. (Price Range, Room Type)
1092
1093 3. (Property Type, Room Type)
1094
1095 4. (Price Range, Maximum Occupancy)
1096
1097 5. (Property Type, Maximum Occupancy)

1096 This selection of factors creates a connected graph that effectively captures the relationships between
1097 the variables while focusing on strong predictive capabilities for the target variable.
1098

1099
1100 A.1.6 INTERACTION PROPOSAL EXAMPLE TRANSLATION

1102 SYSTEM You will be provided a message transcript, where at the end there is a list of factor pairs in the
1103 form of [Target, Condition]. Translate the proposed factor pairs with Target and Condition into JSON.
1104 You must use the names exactly as provided at the very end of the message, under 'VARIABLE
1105 NAMES'

1106 USER

1108 The maximum number of factors allowed is equal to the number of variables in the model. In this
1109 case, we have 5 variables: 1 Target Variable, 1 Observed Variable, and 3 Latent Variables. Therefore,
1110 we can propose a maximum of 5 factors.

1111 ### Model Design Discussion

- 1112 1. **Room Type (Target Variable)** and **Maximum Occupancy (Observed Variable)**:
1113
1114 - Question: Does Room Type vary significantly with different values of Maximum Occupancy?
1115
1116 - Answer: Yes, different room types can accommodate different maximum occupancies. For example,
1117 an entire home or apartment typically accommodates more people than a shared or hotel room.
1118
1119 2. **Room Type (Target Variable)** and **Price Range (Latent Variable)**:
1120
1121 - Question: Does Room Type vary significantly with different values of Price Range?
1122
1123 - Answer: Yes, different room types are often associated with different price ranges. For instance,
1124 entire homes or apartments may be priced higher than shared or hotel rooms.
1125
1126 3. **Room Type (Target Variable)** and **Property Type (Latent Variable)**:
1127
1128 - Question: Does Room Type vary significantly with different values of Property Type?
1129
1130 - Answer: Yes, the type of property can influence the room type. For example, a house is more likely
1131 to be listed as an entire home or apartment compared to a condo.
1132
1133 4. **Maximum Occupancy (Observed Variable)** and **Price Range (Latent Variable)**:
1134
1135 - Question: Does Maximum Occupancy vary significantly with different values of Price Range?
1136
1137 - Answer: Yes, properties that accommodate more people may also be priced differently, often higher.
1138
1139 5. **Maximum Occupancy (Observed Variable)** and **Property Type (Latent Variable)**:
1140
1141 - Question: Does Maximum Occupancy vary significantly with different values of Property Type?

1134 - Answer: Yes, different property types can have different maximum occupancies. For example, a
 1135 townhouse may accommodate more people than a condo.

1136 ### Final Set of Factors

1137 After considering the relationships, I will select the following pairs of variables to form the factors in
 1138 the model:

1139 1. (Room Type, Maximum Occupancy) - It is easier to guess the Room Type given the Maximum
 1140 Occupancy.

1141 2. (Room Type, Price Range) - It is easier to guess the Room Type given the Price Range.

1142 3. (Room Type, Property Type) - It is easier to guess the Room Type given the Property Type.

1143 4. (Maximum Occupancy, Price Range) - It is easier to guess the Maximum Occupancy given the
 1144 Price Range.

1145 5. (Maximum Occupancy, Property Type) - It is easier to guess the Maximum Occupancy given the
 1146 Property Type.

1147 ### List of Factors

1148 1. (Maximum Occupancy, Room Type)

1149 2. (Price Range, Room Type)

1150 3. (Property Type, Room Type)

1151 4. (Price Range, Maximum Occupancy)

1152 5. (Property Type, Maximum Occupancy)

1153 This selection of factors creates a connected graph that effectively captures the relationships between
 1154 the variables while focusing on strong predictive capabilities for the target variable.

1155 **VARIABLE NAMES:** (Use Name exactly as defined below)

1156 'Price Range'; 'Maximum Occupancy'; 'Property Type'; 'Room Type';

1157 ASSISTANT

1158 {"Features": [{"Target": {"Name": "Room Type"}, "Condition": {"Name": "Maximum Occu-
 1159 pancy"}}, {"Target": {"Name": "Room Type"}, "Condition": {"Name": "Price Range"}}, {"Target":
 1160 {"Name": "Room Type"}, "Condition": {"Name": "Property Type"}}, {"Target": {"Name": "Max-
 1161 imum Occupancy"}, "Condition": {"Name": "Price Range"}}, {"Target": {"Name": "Maximum
 1162 Occupancy"}, "Condition": {"Name": "Property Type"}}]}

1163

1164

1165 A.1.7 QUESTION GENERATION EXAMPLE

1166 SYSTEM You will translate a formalized request for a (conditional) distribution into a concise question
 1167 expressed in natural language. The formalized request contains three entries: DOMAIN, TARGET,
 1168 CONDITION. formalized request is for $P(\text{TARGET} \mid X_0 \in x_0, \dots, X_k \in x_k)$ under the DOMAIN.

1169 DOMAIN provides background information for the request.

1170 TARGET is the VARIABLE whose distribution is requested.

1171 CONDITION provide a (possibly empty) list of conditions, each in the form ' $X_i \in x_i$ '. The request
 1172 is conditioned on VARIABLE X_i taking any value in the set x_i .

1173 USER

1174 DOMAIN: Airbnb listings in Austin, TX in 2023 with a price between \$201 and \$500.

1175 TARGET: Maximum Occupancy: ['1-3 people'; '4 or more people']

1176 CONDITION: Price Range \in 2208 ['\$401 - \$500']

1177

1178 ASSISTANT

1188 For an Airbnb listing in Austin, TX in 2023 with a price between \$201 to \$500 and is listed to
1189 accommodate a maximum of 4 or more people, what is the probability distribution of room type?

1190 USER

1192 DOMAIN: Airbnb listings in Austin, TX in 2023 with a price between \$201 and \$500.

1193 TARGET: Maximum Occupancy: ['1-3 people'; '4 or more people']

1194 CONDITION: Price Range \in ['\$401 - \$500']

1196 ASSISTANT

1197 For an Airbnb listing in Austin, TX in 2023 priced between \$401 and \$500, what is the probability
1198 distribution of maximum occupancy?
1199

1200

1201 A.1.8 NUMERIC CONSTRAINT PROPOSAL EXAMPLE DIALOGUE

1202

1203 SYSTEM

1204 You will estimate a probability distribution given the definition of the Target Variable. For each
1205 probability, give the best exact estimate instead of intervals.

1206 USER

1208 For an Airbnb listing in Austin, TX in 2023 priced between \$401 and \$500, what is the probability
1209 distribution of maximum occupancy?

1210 '1-3 people'; '4 or more people'

1211

1212 ASSISTANT

1213 To estimate the probability distribution of maximum occupancy for an Airbnb listing in Austin, TX,
1214 priced between \$401 and \$500 in 2023, we can consider general trends in the market.

1215 Given the price range, it's likely that listings in this category are more spacious and may cater to
1216 larger groups. Therefore, the probability distribution can be estimated as follows:

1217

1218 - Probability of '1-3 people': 30%

1219

1220 - Probability of '4 or more people': 70%

1221

1222 This distribution reflects the expectation that higher-priced listings are more likely to accommodate
1223 larger groups.

1223

1224 A.1.9 NUMERIC CONSTRAINT PROPOSAL EXAMPLE TRANSLATION

1225

1226 SYSTEM You will extract a statistician's probability distribution estimate from a given message into
1227 JSON. The estimate is for the Variable:

1228

1229 Target: Maximum Occupancy: '1-3 people'; '4 or more people'

1230

1231 Use exactly this naming and ordering of the Variable Value for the extraction. Extract a distribution
1232 of 2 numbers corresponding to each of its Value.

1232

1233 USER

1234 To estimate the probability distribution of maximum occupancy for an Airbnb listing in Austin, TX,
1235 priced between \$401 and \$500 in 2023, we can consider general trends in the market.

1236

1237 Given the price range, it's likely that listings in this category are more spacious and may cater to
1238 larger groups. Therefore, the probability distribution can be estimated as follows:

1239

1240 - Probability of '1-3 people': 30%

1241

1242 - Probability of '4 or more people': 70%

1243

1244 This distribution reflects the expectation that higher-priced listings are more likely to accommodate
1245 larger groups.

1242 ASSISTANT
1243
1244 {'Target': {'Name': 'Maximum Occupancy', 'Value': ['1-3 people', '4 or more people']}, 'Probabil-
1245 ity': [30, 70]}

1246 PROCESSED
1247 {"Target": {"Name": "Maximum Occupancy", "Value": ["1-3 people", "4 or more people"]}, "Condi-
1248 tion": {"Name": "Price Range", "Value": ["401-500"]}, "Probability": [0.3, 0.7]}

1249
1250 A.1.10 ZERO-SHOT CHAIN-OF-THOUGHT EXAMPLE DIALOGUE
1251
1252 SYSTEM
1253
1254 You will estimate a probability distribution given the definition of the Target Variable. For each
1255 probability, give the best exact estimate instead of intervals.

1256 USER
1257
1258 For an Airbnb listing in Austin, TX in 2023 with a price between \$201 to \$500 and is listed to
1259 accommodate a maximum of 4 or more people, what is the probability distribution of room type?
1260 'entire home or apartment'; 'private room'; 'shared or hotel room'

1261 ASSISTANT
1262
1263 To estimate the probability distribution of room types for an Airbnb listing in Austin, TX in 2023,
1264 given the specified criteria (price between \$201 to \$500 and accommodating 4 or more people), we
1265 can provide the following estimates based on trends observed in the market:
1266 - **Entire home or apartment**: 70%
1267 - **Private room**: 25%
1268 - **Shared or hotel room**: 5%
1269
1270 These estimates reflect the likelihood of each room type being listed under the given conditions.
1271

1272 A.1.11 ZERO-SHOT CHAIN-OF-THOUGHT EXAMPLE TRANSLATION
1273
1274 SYSTEM
1275
1276 You will extract a statistician's probability distribution estimate from a given message into JSON.
1277 The estimate is for the Variable: Room Type: 'entire home or apartment'; 'private room'; 'shared or
1278 hotel room'
1279
1280 Use exactly this naming and ordering of the Variable Value for the extraction. Extract a distribution
1281 of 3 numbers corresponding to each of its Value.

1282 USER
1283
1284 To estimate the probability distribution of maximum occupancy for an Airbnb listing in Austin, TX,
1285 priced between \$401 and \$500 in 2023, we can consider general trends in the market.
1286
1287 Given the price range, it's likely that listings in this category are more spacious and may cater to
1288 larger groups. Therefore, the probability distribution can be estimated as follows:
1289 - Probability of '1-3 people': 30%
1290 - Probability of '4 or more people': 70%
1291
1292 This distribution reflects the expectation that higher-priced listings are more likely to accommodate
1293 larger groups.

1294 ASSISTANT
1295 {"Target": {"Name": "Room Type", "Value": ["entire home or apartment", "private room", "shared or
1296 hotel room"]}, "Probability": [70, 25, 5]}

1297 PROCESSED

1296 {"Target": {"Name": "Room Type", "Value": ["entire home or apartment", "private room", "shared or
1297 hotel room"]}, "Probability": [0.7, 0.25, 0.05]}

1299 B ADDITIONAL DETAILS ON DATASET AND PREPROCESSING

1301 B.1 DATASET SPLITS

- 1303 1. On Inside Airbnb, we use Ashville, Austin, Chicago, New Orleans, Pacific Grove, and
1304 Rhode Island for evaluation, and Twin Cities for development.
- 1305 2. On American Time-Use survey, we use 2018, 2020, 2022 as evaluation, and 2023 as
1306 development.
- 1307 3. On World Values Survey, we use Malaysia, New Zealand, Rwanda, Sweden, United States,
1308 and Uruguay for evaluation, and no development.

1310 B.2 PREPROCESSING

1312 For each dataset, we use a subset of all available columns. We also discretize any continuous data
1313 into ranges, and coarsen any discrete variables with too many values. All such choices were made
1314 before any significant tuning of the prompts and hyper-parameters of our pipeline or the prompt for
1315 zero-shot Chain-of-Thought baseline.

1317 B.2.1 INSIDE AIRBNB

1318 Many columns of the Inside Airbnb dataset have missing values for a significant proportion of rows.
1319 We thus ignored any column with too high a proportion of missing values, and then manually picked
1320 a subset of 8 columns that we judged to be interesting. The processed variables and their possible
1321 values are included in Table 2.

1323 Table 2: Schema for our processed Inside Airbnb dataset.

Column Name	Possible Values
Number of Bedrooms	studio or 1 bedroom, 2 bedrooms, 3 bedrooms, 4 or more bedrooms
Number of Bathrooms	shared or single bathroom, 2 bathrooms, 3 or more bathrooms
Superhost Status	Superhost, Not Superhost
Room Type	entire home or apartment, private room, shared or hotel room
Total Beds	1 bed, 2 beds, 3 beds, 4 or more beds
Review Score	less than 4.4, 4.5 to 4.8, at least 4.9
Max Accommodates	1, 2, 3, 4 or more
Price	under \$50, \$51 to \$100, \$101 to \$200, \$201 to \$500, at least \$501

1335 B.2.2 AMERICAN TIME-USE SURVEY

1336 We use most of the frequently used subset¹⁶ of ATUS. The processed variables and their possible
1337 values are included in Table 3.

1339 B.2.3 WORLD VALUES SURVEY

1341 Again, we manually picked most of the objective demographics variables as well as columns that are
1342 not too granular. The processed variables and their possible values are included in Table 4.

1344 C EXAMPLE QUESTIONS

1346 C.1 AIRBNB

1348 **Split:** Chicago, IL

1349 ¹⁶<https://www.bls.gov/tus/other-documentation/freqvariables.pdf>

1350

1351

Table 3: Schema for our processed American Time-Use Survey dataset.

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

Table 4: Schema for our processed World Values Survey dataset.

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Column Name	Possible Values
Sex	Male, Female
Age	15-29, 30-44, 45-64, 65-85
Region	Northeast, Midwest, South, West
Marital Status	Married, Widowed, Divorced, Separated, Never Married
Metropolitan Residency Status	Metropolitan, Non-metropolitan
Labor Force Status	Employed, Unemployed, Not in Labor Force
Household Composition	Children Under 18 Present in Household, No Children Under 18 in Household
Day of Week	Weekday, Weekend
High School/College Enrollment	Currently Enrolled, Not Currently Enrolled
Activity	Personal Care, Sleep, and Sustenance, Leisure, Sports, and Social, Traveling and Commuting, Work and Education, Household and Other

Column Name	Possible Values
Importance of family in life	Not at all important, Not very important, Rather important, Very important
Importance of friends in life	Not at all important, Not very important, Rather important, Very important
Importance of leisure time in life	Not at all important, Not very important, Rather important, Very important
Importance of politics in life	Not at all important, Not very important, Rather important, Very important
Importance of work in life	Not at all important, Not very important, Rather important, Very important
Importance of religion in life	Not at all important, Not very important, Rather important, Very important
Member of religious organization	Member, Not member
Member of sport or recreational organization	Member, Not member
Member of art, music or educational organization	Member, Not member
Member of labour union	Member, Not member
Member of political party	Member, Not member
Member of environmental organization	Member, Not member
Member of humanitarian or charitable organization	Member, Not member
Marital Status	Married, Divorced, Separated, Widowed, Single
Age	18-29, 30-44, 45-64, 65+,
Sex	Male, Female
Labor Force Status	Employed, Unemployed, Not in Labor Force

Target:

Number of Bathrooms: shared or single bathroom; 2 bathrooms; 3 or more bathrooms

Conditions:

Number of Bedrooms = 3 Bedrooms

Natural Language Question:

For an Airbnb listing with 3 bedrooms in Chicago, IL in 2023, what is the probability distribution of its number of bathrooms?

1404 **Answer:**
1405 shared or single bathroom: 0.435
1406 2 bathrooms: 0.476
1407 3 or more bathrooms: 0.089
1408
1409 C.2 AMERICAN TIME-USE SURVEY
1410
1411 **Split:** 2020
1412 **Target:**
1413 Labor Force Status : Employed; Unemployed; Not in Labor Force
1414
1415 **Conditions:**
1416 Age = 30 - 44
1417 **Natural Language Question:**
1418 For a person aged 30-44 in the United States population in 2020, what is the probability distribution
1419 of their Labor Force Status?
1420 **Answer:**
1421 Employed: 0.797
1422 Unemployed: 0.051
1423 Not in Labor Force: 0.153
1424
1425 C.3 WORLD VALUE SURVEY
1426
1427 **Split:** Sweden
1428 **Target:**
1429 Importance of politics in life : Not at all important; Not very important; Rather important; Very
1430 important
1431 **Conditions:**
1432 Member of humanitarian or charitable organization = Member
1433
1434 **Natural Language Question:**
1435 For a person in Sweden aged 18 or older in 2010-2014 who is not a member of a humanitarian
1436 or charitable organization, what is the probability distribution of their views on the importance of
1437 politics in their life?
1438 **Answer:**
1439 Not at all important: 0.11
1440 Not very important: 0.299
1441 Rather important: 0.441
1442 Very important: 0.15
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457