

# Paradoxical parsimony: How latent complexity favors explanatory simplicity

Anonymous CogSci submission

## Abstract

Investigating how people evaluate more or less complex explanations has been a focal point of research. However, previous studies have either focused on choice between a limited set of explanations or do not systematically quantify the explanations' complexity. We provide a new approach for modeling explanation selection that foregrounds the balance between observed and latent structure in the mechanism being explained. We combine a Bayesian framework with program induction, enabling coverage of unbounded partially observable model space through sampling, and reflecting how a simplicity bias emerges naturally in this setting. Through simulation, we identify two novel principles: (1) simpler explanations should be favored as latent uncertainty (the number of hidden variables) increases; (2) latent structure is attributed a larger role when the observable patterns become less compressible. We found that these principles were reflected in human judgments, indicating that people are sensitive to latent uncertainty when selecting between explanations.

**Keywords:** explanation; mechanism; program induction; simplicity; hidden variables; inductive reasoning

## Introduction

We acquire our causal theories by interacting with our environment in limited ways and observing what can, and propagate them socially by providing one another with causal explanations. Inevitably these explanations fall short of replicating their explananda perfectly, leaving unexplained variance, limited predictive power and uncertain generalizability. This raises deep questions: What is the right level of complexity for an explanation? And, what does this balance depend on?

Explanations can strike as both inappropriately complex or overly simplistic, given the context and the evidence they are based on. But it can be hard to pin down how much of this comes down to personal preference or if there is sometimes a normative answer to how complex a good explanation should be. For example, for the average person, conspiracy theories strike a poor balance, typically presenting as implausible even as they tie together diverse evidence (Wojtowicz & DeDeo, 2020) while superstitious beliefs arguably overdetermine effects by imputing both physical and cosmic causes (Jin, Jensen, Gottlieb, & Ferrera, 2022). Within scientific research, a critical consideration is the bias-variance trade-off when selecting a model — a foundational question that dictates the balance between overfitting and underfitting (Doroudi & Rastegar, 2023; Lucas, Griffiths, Williams, & Kalish, 2015; Brighton & Gigerenzer, 2015). However, this is typically treated as a purely statistical question rather than one that can depend on the explanatory context and the explanation's ontological implications.

Investigating how people evaluate more or less complex explanations of the same evidence has been a focal research

point (Lombrozo, 2007; Zemla, Sloman, Bechlivanidis, & Lagnado, 2023; Johnson, Valenti, & Keil, 2019). However, whether people evaluate an explanation as good or bad may not directly reflect whether they hold the corresponding inner beliefs or not: at times, explanation serves purposes of communication somewhat distinct from communicators' representation or understanding (Lombrozo, 2010). This is particularly evident when individuals have to use natural language to express their beliefs (Zemla et al., 2023; Sulik, van Parijs, & Lupyan, 2023). Furthermore, while many studies concentrate on the diagnostic or *token* level of explanation (i.e. a particular explanation for how or why a specific event happened), fewer have investigated how people make predictions or represent mechanisms at the *type* level of causation (i.e. a general theory to explain how a type of events comes into being). Additionally, studies have predominantly focused on a few specific causal structures, such as the common-effect (e.g. where a symptom might stem from one disease or several; Pacer & Lombrozo, 2017; Johnson et al., 2019) and chain structures (e.g. where intermediary variables are incorporated into an explanatory narrative or sequence of events; Johnson & Ahn, 2015; Johnson et al., 2019). In reality, causal mechanisms could be arbitrarily more mixed, diverse or awkward to describe.

Recent studies have also prompted people to express their beliefs by drawing causal graphs, enabling the collection of inner representations without reliance on natural language (Bramley, Lagnado, & Speekenbrink, 2015; Bramley, Dayan, Griffiths, & Lagnado, 2017; Gong, Gerstenberg, Mayrhofer, & Bramley, 2023; Gong & Bramley, 2023; Kushnir, Gopnik, Lucas, & Schulz, 2010). However, again these have predominantly focused on a limited set of relationships (mainly generative) or functional forms (often noisy-OR disjunctive combinations of causes). This is largely a consequence of computing convenience since modeling inferences over a larger hypothesis space that covers various types of functions as well as connectivity patterns rapidly gets unwieldy, making it difficult or intractable to compute the posterior distributions necessitated by a traditional Bayesian analysis. Meanwhile, these graphical approaches have paid less attention to the complexity of beliefs and rarely investigated the factors that can influence what level of complexity is appropriate.

In this study, we take an approach that supports studying induction in an unbounded hypothesis space and unpacks why beliefs are inevitably bounded in their complexity. One way to make sense of inference in an unbounded hypothesis space is adopt a program induction approach — assuming learners sample explanatory hypotheses by composing them

stochastically from a sufficiently expressive grammar. We here consider probabilistic context-free grammars (PCFGs), and concretely assume a grammar that can be used to produce any causal rule expressible in propositional logic applied to the question of how a set of putative causes combine to determine an effect (Buchanan, Tenenbaum, & Sobel, 2010). Similar formalisms have been applied to a range of concept learning settings (Bramley & Xu, 2023; Goodman, Tenenbaum, Feldman, & Griffiths, 2008; Zhao, Lucas, & Bramley, 2024; Fränken, Theodoropoulos, & Bramley, 2022), and recent work has argued for the framework’s general applicability to explaining our capacity for induction (Piantadosi, 2021; Bramley, Zhao, Quillien, & Lucas, 2023). A core feature of a PCFGs is a built-in inductive bias favoring shorter and simpler explanations – longer strings typically involve more rule applications, which typically imply geometrically decreasing probabilities.

Using this representational approach, this paper will explicate a specific theoretical question: How does knowledge about the existence of hidden components the environment influence the selection between more complex vs. simple causal explanations? A familiar domain, such as a physical or social setting, could contain different levels of complexity due to the number of hidden (i.e. known to exist yet practically impossible to observe) variables (Johnson et al., 2019; Lucas, Holstein, & Kemp, 2014; Valentin, Bramley, & Lucas, 2022). Even though we do not know the states of such hidden variables, knowing they exist and could feature in the mechanism of interest can influence how we explain the roles that the explicit variables play. In particular, as we will show, this affects how much variance or noise we should tolerate in our explanations.

Suppose you know you are in a fully observed and deterministic setting and have observed outcomes that cannot be explained by any small set of variables. Determinism demands that outcomes be perfectly explainable without recourse to randomness or noise, so a rational observer will invoke as many variables as necessary. However, if you have not observed all the variables in the setting, it could still be that a small set of the variables you have observed are causative of the outcome. The residual variance could be explained by the action of the hidden variables. To illustrate, suppose a learner has a dozen encounters with a system involving four binary inputs  $X_1$  to  $X_4$ . Suppose every outcome but one can be explained by the presence or absence of  $X_1$ ; but the one remaining outlier can then only be explained via positing a complicated conjunction of all four observable variables. With no hidden variables, the learner has no better option than to posit this complex explanation. However, if they learn that even one hidden variable exists  $H_1$ , the learner could more simply attribute the outlier to the action of that hidden variable, resulting in a simple but imperfect explanation in terms of the observable features. That is, their explanation may evokes only  $X_1$  explicitly and implicitly marginalizes over a hidden complicating possibility, the

$X_1$	$X_2$	$E$	<b>AND:</b> 0001, 1000
0	0	0	<b>OR:</b> 0111, 1110
0	1	0	<b>XOR:</b> 0110, 1001
1	0	0	<b>Singular:</b> 0011, 1100, 0101, 1010
1	1	1	<b>One Exception:</b> 0010, 1011, 0100, 1101

Figure 1: Stimuli tested in the paper. Each trial involved four observations that show the outcomes of four combinations of  $X_1$  and  $X_2$ . Fourteen stimuli were classed into five categories given their logical class. “Singular” stimuli can be explained by a single observed variable; the “one exception” group includes stimuli that cannot be explained in terms of the observable variables using any of the standard two place Boolean combinators. The brown color indicates stimuli that must invoke negation to explain outcomes under the “AND”, “OR”, and “Singular” categories.

unobserved states and involvement of  $H_1$ . As the number of hidden variables increases it becomes increasingly plausible that the outcome could be entirely due to some such hidden mechanism, increasing the credence that a normative learner should give a “null hypothesis” type explanation where none of the measured variables matter at all for the outcome: at least implicitly the randomness of the outcome is driven by the increasing chance of its being controlled entirely by unobserved mechanisms.

We will first introduce our model, and use simulations to demonstrate the idea above. We will then conduct an empirical experiment to see whether people exhibit the same sorts of sensitivity to latent complexity as the model.

## A program induction approach to express causal rules

In this section, we demonstrate how explanations in terms of observable variables should shift from complex to simple as the number of hidden variables in a situation increases. We adopt a Bayesian approach where the hypotheses and prior are defined via program induction. We work with a deterministic likelihood such that a complete mechanistic explanation can only be correct if it fully explains all observations.

### Probabilistic context-free grammars

We use AND ( $\wedge$ ), OR ( $\vee$ ), and NOT ( $\neg$ ) as basic primitives to express causal rules. For example,  $(X_1 \wedge X_2) \vee (\neg X_3)$  means that for the effect to appear, you need either the combined presence of  $X_1$  and  $X_2$ , or the absence of  $X_3$ . To sample explanations, we use a simple disjunctive normal form grammar (DNF) as outlined in Goodman et al. (2008) and Buchanan et al. (2010):

$$S \rightarrow \forall x, l(x) \Leftrightarrow (D) \quad (1)$$

$$D \rightarrow (C) \vee D \quad (2)$$

$$D \rightarrow \text{False} \quad (3)$$

$$C \rightarrow N \wedge C \quad (4)$$

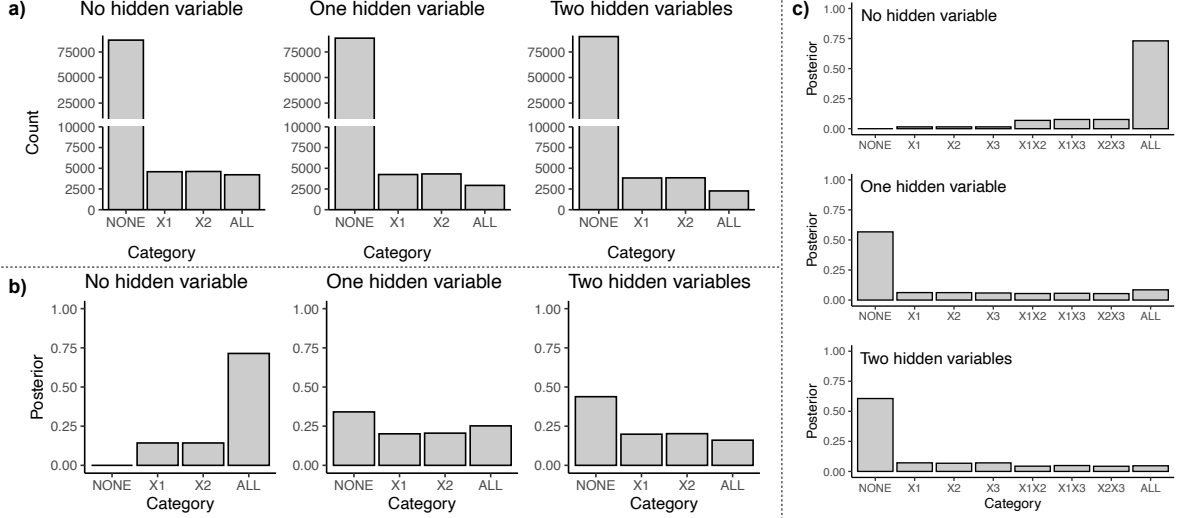


Figure 2: Model simulation results. a) The distribution of simulated hypotheses in each category under a sample of 100,000. b) Posterior distributions averaged from stimuli in Figure 1. c) Posterior distributions averaged from all possible stimuli for three explicit variables.

$$C \rightarrow \text{True} \quad (5)$$

$$N \rightarrow \neg N \quad (6)$$

$$N \rightarrow P \quad (7)$$

$$P \rightarrow X_1 \dots X_n \quad (8)$$

For each sampled hypothesis, the generative mechanism a disjunctive form placeholder ( $D$ ) and iteratively replaces the terms according to the production rules above until all terms in the expression reach their terminal status (*True*, *False*, or one of the variables from  $X_1$  to  $X_n$ ). We apply a principle of indifference, assuming that during each replacement process, each valid production rule has equal probability of application. Production rules (6) and (7) result in a higher probability for a generative relationship than a preventative relationship (i.e. direct assertions are the default and an additional negation step is required to invert them making negative rules less likely to be generated than positive). Of note, we chose to do this rather than treating positive and negative assertions as symmetric because this both reflects the intuition that negation increases an assertion’s complexity and aligns with empirical findings that generative relationships are easier to discover than preventative ones especially when the effect is assumed to be absent by default (Gong & Bramley, 2023; Cheng, 1997). Future work will investigate the extent to which this affects predictions.

### Likelihood and the hidden variables

The approach outlined above allows us to generate a prior sample of hypotheses that are naturally biased toward simplicity, and in the limit of infinite sampling would cover all expressions in propositional logic relating the variables to the outcome. Since we assume a deterministic setting we can simply filter this sample, ruling out all models that cannot explain the observation to arrive at a posterior sample and marginalize over this to calculate posterior probabilities for involvement of different variables.

When there are hidden variables, those variables whose presence or absence is unknown are included in  $X_1 \dots X_n$  as well, and the likelihood is calculated by marginalizing over all possible combinations of states of the hidden variables (assuming uniform priors on whether they are present or absent on each trial). For example, if there are two hidden variables, the likelihood would be marginalized (averaged) over the four states of presence and absence of the two hidden variables ( $\{0,0\}$ ,  $\{0,1\}$ ,  $\{1,0\}$ ,  $\{1,1\}$ ).

### Clustering the hypotheses

In the sampling process, we will find that many hypotheses are syntactically different but semantically identical (e.g., both  $X_1 \wedge X_2$  and  $X_2 \wedge X_1$  express a conjunction between  $X_1$  and  $X_2$  at the semantic level). We also note that rules can be harder or easier to articulate, especially when the rule is complicated or involves hidden variables. Therefore, instead of focusing solely on the syntax or semantics, we concentrate on a higher level by clustering hypotheses into different categories depending on which observed variables they involve.

When there are two explicit variables,  $X_1$  and  $X_2$ , we identify four categories: Rules in which (1) None of  $X_1$  and  $X_2$  are relevant; (2) only  $X_1$  is relevant; (3) only  $X_2$  is relevant; (4) both  $X_1$  and  $X_2$  are relevant. These categories are referred to as “NONE”, “X1”, “X2”, and “ALL” throughout the rest of the paper. Each hypothesis belongs to one of these categories based on whether the outcome predictions can depend on the state of  $X_1$  or  $X_2$ .

### Simulation

Figure 2a illustrates the distribution of hypotheses across different categories when sampling 100,000 hypotheses. The majority of hypotheses fall into the “NONE” category, with the fewest falling into the “ALL” category. This distribution aligns with the intuition that explanation complexity in-

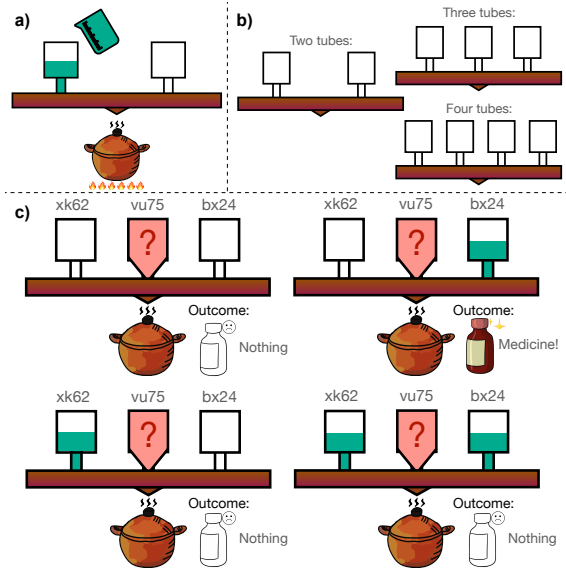


Figure 3: The cover story and stimulus example. a) The pot and equipment. b) Different number of tubes. c) One experimental trial that represents “0100” situation in Figure 1.

creases in the order:  $NONE < X_1 = X_2 < ALL$ . The “ALL” category involves longer expressions that encompass relationships invoking roles for both  $X_1$  and  $X_2$ .

A dominance in prior does not imply that “NONE” is always the best answer. With no hidden variables, the “NONE” category can only explain observations where the effect always appears or never occurs, irrespective of  $X_1$  and  $X_2$ . We tested the model’s predictions on 14 types of observations where the outcome does not remain constant (see Figure 1), a set that will be empirically validated later. As shown in Figure 2b, when no hidden variables are present, explanations often require both  $X_1$  and  $X_2$  to account for the observations. However, this dynamic shifts as the number of hidden variables increases. The posterior probability of the “NONE” category increases, along with those of the “ $X_1$ ” and “ $X_2$ ” categories. It implies that complex situations can more plausibly be attributed fully or partially to the influence of hidden variables. Consequently, our best explanations, regarding the role of the observed variables, revert to simpler ones.

Meanwhile, the role of hidden variables could be significant or minimal depending on the types of stimuli. If the ground truth of a stimulus is already simple (i.e., the prior is already high), the pressure to leverage hidden variables is low. Figure 5 shows the priors of five different types of stimuli (in the strip) and the model predictions (marked as X). The differences among no, one, and two hidden conditions increases when the ground truth prior decreases (from top to the bottom). This means the model only makes the switch when the explanation would otherwise be highly complex.

Figure 2c shows how the model predictions extend to scenarios involving three explicit variables. Although in this simulation section we applied a limited number of primitives and variables, future work will test it on more complex situations that involve different primitives and more variables.

## Experiment

### Methods

**Participants** 90 participants (31 female, 58 male, 1 preferred not to say, aged  $41 \pm 11$ ) were recruited via Prolific Academic and were paid £2. The task took around 18 minutes. The anonymous data and analysis code, as well as experiment procedure are available (<https://bit.ly/47X1Zv6>).

**Design** Participants were asked to imagine themselves as “medical alchemists” who needed to determine the roles of different ingredients in producing medicine (Figure 3a). For each trial, they observed the brewing process of four potions, which were identical before any ingredients were added. Ingredients were added through equipment that could contain two, three, or four tubes (Figure 3b), and then the outcome showed whether the potion successfully became medicinal or not (Figure 3c). Participants were further instructed that some tubes might be covered with a red cloth (Figure 3c), indicating that participants would not know whether corresponding ingredients were added to each potion or not.

As a first foray, we only tested the setting with two explicit variables and 14 stimuli (Figure 1). We examined three conditions where the number of hidden variables varied from 0 to 2. This implies that equipment with two, three, or four tubes would always have zero, one, or two tubes covered by red cloth, respectively. The conditions and stimuli were arranged using the Latin Square design so that each participant experienced all observed patterns, but only in one of the three hidden complexity conditions. This resulted in three stimulus lists to which participants were randomly assigned.

After viewing the four observations for each trial, participants were asked a forced-choice question: which one of the following statements best reflects the truth about the two focal causes: None of  $X_1$  and  $X_2$  is relevant; only  $X_1$  is relevant; only  $X_2$  is relevant; both  $X_1$  and  $X_2$  are relevant. Here,  $X_1$  and  $X_2$  correspond to ingredient names that varied from trial to trial. Participants were also asked to provide a confidence rating for their response on a 0-100 scale.

**Procedure** Before beginning the task, participants were instructed about the cover story, the meaning, and examples of relevance (generation, prevention, or rules that combine it with other ingredients). The deterministic setting was emphasized by explaining that if the same ingredients are added through the equipment, the result will always be the same. Participants had to pass comprehension check questions before the task. The order of the four observations in each trial and the positions of tubes covered by red cloths were randomized among trials. The order of trials and the order of four forced choices were randomized among participants.

### Results

Figure 4 shows participants’ answer for each category under three conditions with different numbers of hidden variables. Three answer distribution differed ( $\chi^2$  test of independence:  $\chi^2(6) = 19.52$ ,  $p = .003$ ). The distribution of answers sig-

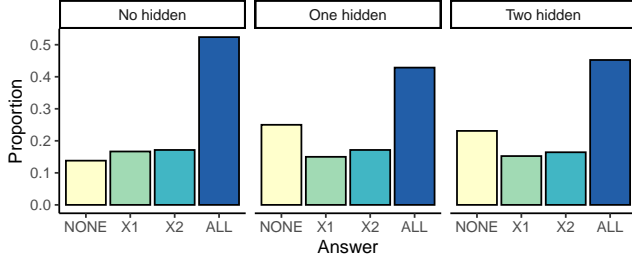


Figure 4: The proportion of participants' answer for each option.

nificantly differed between “No Hidden” vs. “One Hidden” ( $\chi^2(3) = 17.92, p < .001$ ), or “No Hidden” vs. “Two Hidden” ( $\chi^2(3) = 12.34, p = .006$ ), but not “One Hidden” vs. “Two Hidden” conditions ( $\chi^2(3) = 0.66, p = .88$ ). When hidden variables were present, the proportion choosing “NONE” increased, while the proportion choosing “ALL” decreased, which reflects the previous simulation results (Figure 2b).

**Proportions by stimulus types** Figure 5 shows participants' answers broken down by type of stimuli. For each type of stimuli, we focus on whether participants chose none, all, or one of the variables as related (i.e. “X1” and “X2” are merged as the category “ONE” here). The stimulus types are ordered according to how likely the corresponding deterministic observable variable explanation (Figure 1) would be sampled from the PCFGs. This means that any prior sampled hypothesis would be more likely to fall under singular ground truth than AND ground truth, and so on. For each type, there was a more or less tendency to favor “NONE” category and disfavor “ALL” category when the number of hidden variables increased from zero to one or two. However, the answer distributions only significantly differed between conditions for “OR” ( $\chi^2(4) = 13.37, p = 0.01$ , Fisher's exact test was used due to small numbers in some cells:  $p = 0.006$ ) and “One exception” ( $\chi^2(4) = 13.08, p = 0.01$ ) stimuli. These two types also have lower priors than all other types except for the “XOR”. This is aligned with our model prediction that when the situation is more complex, people will turn to simpler explanations whenever it is possible. We will later discuss why we think “XOR” was an exception.

**Asymmetry between generation and prevention** We divided “Singular”, “AND”, and “OR” stimuli into generation vs. prevention types, according to whether they need to contain a negation (prevention) of  $X_1$  or  $X_2$  in the rule in order to explain the outcomes (see Figure 1). The prevention rules are harder to sampled in the hypothesis space and hence receive low prior (Figure 6). Accordingly, participants more often referred to the “NONE” option, and less often referred to the “ALL” option, when facing a prevention ground truth ( $ps < 0.001$ , Figure 6). The tendency here did not further differ between the hidden variable conditions.

**The one-exception cases** We finally explore the patterns for the one exception stimuli. They have either only one presence ( $\{0,0,1,0\}, \{0,1,0,0\}$ ) or only one absence in the outcome ( $\{1,1,0,1\}, \{1,0,1,1\}$ ). Preferences seemed to differ be-

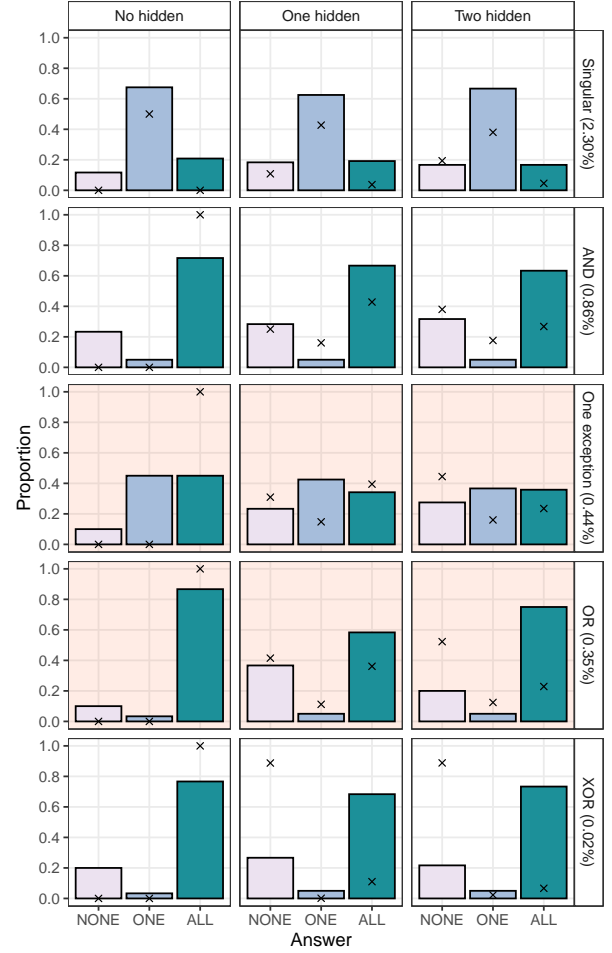


Figure 5: Participants' answer for under each stimulus type. Model predictions are marked as “X”. Numbers in the strips indicate the stimulus probabilities in the prior. Coral color indicates distributions significantly differ between conditions.

tween this two types. When facing one-present stimuli, participants, especially in the one-hidden condition, tend to refer to one variable rather than select “ALL” or “NONE”. This reflects that in the situation when a singular variable can explain all but one data point (e.g.  $X_1$  can fully explain  $\{0,0,1,0\}$  as long as we turn the fourth outcome from 0 to 1), participants may leverage the hidden variable in the system to help maintain a singular and generative answer, an answer that is more informative than “NONE” and less complex than “ALL”. This tendency was not statistically significant, partially because of the sample size. This needs to be further tested in future work probably with situations that involve more explicit variables.

**Confidence** Participants confidence ratings differed between conditions ( $F(2,89) = 12.04, p < .001$ ). This was significant between the no-hidden condition ( $75 \pm 18$ ) and one-hidden condition ( $68 \pm 19, t(89) = 4.12, p < .001$ ), or the no-hidden condition and two-hidden condition ( $66 \pm 21, t(89) = 4.84, p < .001$ ), but not between two hidden variable conditions ( $t(89) = 1.42, p = .48$ ). This is consistent with previous results where the answer proportions did not significantly differ between one hidden variable and two hidden



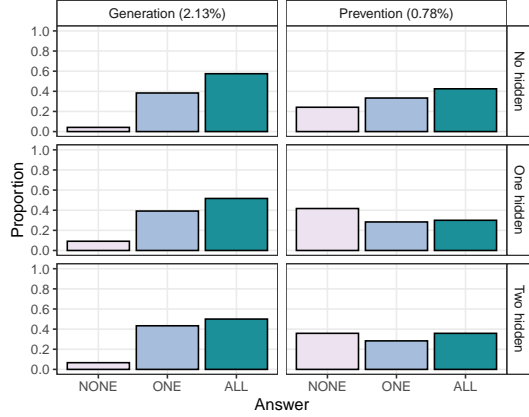


Figure 6: Participants’ answers in generation vs. prevention types of stimuli.

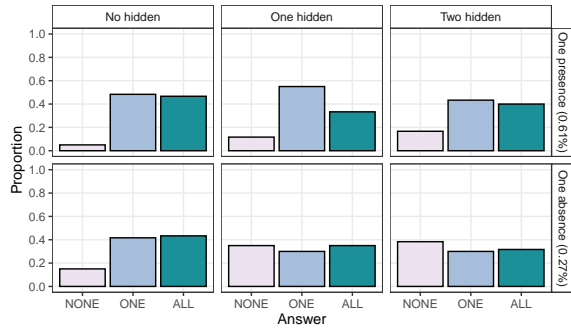


Figure 7: Participants’ answers in two subsets of stimuli under the “one exception” category.

variables (Figure 4). Participants’ confidence ranked from stimuli AND ( $75 \pm 24$ ), Singular ( $70 \pm 23$ ), OR ( $70 \pm 25$ ), XOR ( $68 \pm 25$ ), to one exception ( $67 \pm 26$ ,  $F(4, 89) = 3.90$ ,  $p = .006$ ). Only the difference between AND and one exception was statistically significant ( $t(89) = 3.80$ ,  $p = .003$ ).

## Discussion

In this paper, we provide a model that demonstrates how explanations should differ their complexity in environments involving varying numbers of hidden variables. By leveraging a program induction approach, our model allowed for consideration of a wide variety of belief hypotheses, reflecting the reality that the causal explanations people could form in principle is practically unbounded. With this approach, we demonstrate that complex beliefs are demanded when the environment is fully observed yet lacks a simple explanatory mechanism, while simpler explanations persist when the environment is more uncertain.

Our model provides a new insight into a long-standing philosophical question: How to choose among theories. It also speaks to the modern statistical question of how scientists should choose between models (Doroudi & Rastegar, 2023). A good theory or model should not only be simple but also informative in explaining phenomena (Jefferys & Berger, 1992; Pacer & Lombrozo, 2017). We demonstrate how this tendency towards simplicity and informativeness is reflected in a Bayesian framework, via the prior and likelihood respec-

tively, and how Bayes’ rule can help guide the combination of the two to provide a balanced answer.

With the model simulations, we identified two principles that could be tested empirically: (1) When there are hidden variables, a learner should shift from complex explanations toward simpler explanations; (2) this shift should be more pronounced when the only fully observable explanation for the pattern is more complex (i.e., has a low prior). We showed both these principles are reflected in human preferences. Participants were more inclined to choose simpler explanations when there were hidden variables, and their response distributions were more varied across conditions for the “One exception” and “OR” stimuli in line with their lower priors.

We also observe deviations between human performance and our account. For example, the second principle mentioned above did not manifest for “XOR” stimuli, which received the lowest prior under our PCFGs. The low prior is a consequence of the grammar not containing XOR as a primitive (having to construct it by combining AND, OR and NOT). Although XOR has historically been treated as complex and difficult case due to its non-linearity, particularly in the early connectionist literature, recent research suggests that XOR, representing “either A or B” is a salient possibility for that human reasoners will readily entertain (Jiang & Lucas, 2024; Bramley & Xu, 2023; Gerstenberg & Icard, 2020). Therefore, future work will attempt to adapt the primitives and architecture of the model to explore what better reflects human inductive biases.

Other deviations are suggestive of how human interpretation of the task differs from the assumptions based into the normative model. In our experiment, when no hidden variable exists, a rational learner should never choose “NONE”, given that the outcomes in our stimuli always change according to the changing variable states. However, the “NONE” category still received 14% of answers (Figure 4). Most of these answers appeared in the prevention trials (Figure 6). We suspect a subset of participants had a strong asymmetrical preference toward generative over prevention causation and may have misinterpreted the instruction about judging relevance as pertaining to generative rather than preventative influence (Szollosi, Grigoras, Quillien, Lucas, & Bramley, 2023). Although this asymmetry is reflected in our grammars, future work may need to further test how strong the asymmetry would be by fitting the probability of grammars in the model. Meanwhile, participants were generally more often to choose “ALL” than the model (Figure 4 vs. Figure 2b), and they were relatively insensitive to the difference between one and two hidden variables. This shows that in contrast to our model, human cognition may have a different way of dealing hidden variables; instead of exhaustively incorporating all possible states of hidden variables in the inference process, they may use a more heuristic approach, ignoring the larger, unknown space of things unless they are necessary to explain the observation (Gershman, 2019). Future work needs to examine how cognitive resource limitations factor into this process.

## References

- Bramley, N. R., Dayan, P., Griffiths, T. L., & Lagnado, D. A. (2017). Formalizing neurath's ship: Approximate algorithms for online causal learning. *Psychological Review*, 124(3), 301–338.
- Bramley, N. R., Lagnado, D. A., & Speekenbrink, M. (2015). Conservative forgetful scholars: How people learn causal structure through sequences of interventions. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 41(3), 708–731.
- Bramley, N. R., & Xu, F. (2023). Active inductive inference in children and adults: A constructivist perspective. *Cognition*, 238, 105471.
- Bramley, N. R., Zhao, B., Quillien, T., & Lucas, C. G. (2023). Local search and the evolution of world models. *Topics in Cognitive Science*.
- Brighton, H., & Gigerenzer, G. (2015). The bias bias. *Journal of Business Research*, 68(8), 1772–1784.
- Buchanan, D., Tenenbaum, J., & Sobel, D. (2010). Edge replacement and nonindependence in causation. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 32).
- Cheng, P. W. (1997). From covariation to causation: A causal power theory. *Psychological Review*, 104(2), 367–405.
- Doroudi, S., & Rastegar, S. A. (2023). The bias–variance tradeoff in cognitive science. *Cognitive Science*, 47(1), e13241.
- Fränken, J.-P., Theodoropoulos, N. C., & Bramley, N. R. (2022). Algorithms of adaptation in inductive inference. *Cognitive Psychology*, 137, 101506.
- Gershman, S. J. (2019). How to never be wrong. *Psychonomic Bulletin & Review*, 26, 13–28.
- Gerstenberg, T., & Icard, T. (2020). Expectations affect physical causation judgments. *Journal of Experimental Psychology: General*, 149(3), 599–607.
- Gong, T., & Bramley, N. R. (2023). Continuous time causal structure induction with prevention and generation. *Cognition*, 240, 105530.
- Gong, T., Gerstenberg, T., Mayrhofer, R., & Bramley, N. R. (2023). Active causal structure learning in continuous time. *Cognitive Psychology*, 140, 101542.
- Goodman, N. D., Tenenbaum, J. B., Feldman, J., & Griffiths, T. L. (2008). A rational analysis of rule-based concept learning. *Cognitive Science*, 32(1), 108–154.
- Jefferys, W. H., & Berger, J. O. (1992). Ockham's razor and bayesian analysis. *American Scientist*, 80(1), 64–72.
- Jiang, C., & Lucas, C. G. (2024). Actively learning to learn causal relationships. *Computational Brain & Behavior*, 1–26.
- Jin, Y., Jensen, G., Gottlieb, J., & Ferrera, V. (2022). Superstitious learning of abstract order from random reinforcement. *Proceedings of the National Academy of Sciences*, 119(35), e2202789119.
- Johnson, S. G., & Ahn, W.-k. (2015). Causal networks or causal islands? the representation of mechanisms and the transitivity of causal judgment. *Cognitive Science*, 39(7), 1468–1503.
- Johnson, S. G., Valenti, J., & Keil, F. C. (2019). Simplicity and complexity preferences in causal explanation: An opponent heuristic account. *Cognitive Psychology*, 113, 101222.
- Kushnir, T., Gopnik, A., Lucas, C. G., & Schulz, L. (2010). Inferring hidden causal structure. *Cognitive Science*, 34(1), 148–160.
- Lombrozo, T. (2007). Simplicity and probability in causal explanation. *Cognitive Psychology*, 55(3), 232–257.
- Lombrozo, T. (2010). Causal–explanatory pluralism: How intentions, functions, and mechanisms influence causal ascriptions. *Cognitive Psychology*, 61(4), 303–332.
- Lucas, C. G., Griffiths, T. L., Williams, J. J., & Kalish, M. L. (2015). A rational model of function learning. *Psychonomic Bulletin & Review*, 22(5), 1193–1215.
- Lucas, C. G., Holstein, K., & Kemp, C. (2014). Discovering hidden causes using statistical evidence. In P. Bello, M. Guarini, M. McShane, & B. Scassellati (Eds.), *Proceedings of the 36th annual conference of the cognitive science society* (pp. 892–897).
- Pacer, M., & Lombrozo, T. (2017). Ockham's razor cuts to the root: Simplicity in causal explanation. *Journal of Experimental Psychology: General*, 146(12), 1761–1780.
- Piantadosi, S. T. (2021). The computational origin of representation. *Minds and machines*, 31, 1–58.
- Sulik, J., van Paridon, J., & Lupyan, G. (2023). Explanations in the wild. *Cognition*, 237, 105464.
- Szollosi, A., Grigoras, V., Quillien, T., Lucas, C., & Bramley, N. R. (2023). How do instructions, examples, and testing shape task representations? In *Proceedings of the annual meeting of the cognitive science society* (Vol. 45).
- Valentin, S., Bramley, N. R., & Lucas, C. G. (2022). Discovering common hidden causes in sequences of events. *Computational Brain & Behavior*, 1–23.
- Wojtowicz, Z., & DeDeo, S. (2020). From probability to consilience: How explanatory values implement bayesian reasoning. *Trends in Cognitive Sciences*, 24(12), 981–993.
- Zemla, J. C., Sloman, S. A., Bechlivanidis, C., & Lagnado, D. A. (2023). Not so simple! causal mechanisms increase preference for complex explanations. *Cognition*, 239, 105551.
- Zhao, B., Lucas, C. G., & Bramley, N. R. (2024). A model of conceptual bootstrapping in human cognition. *Nature Human Behaviour*, 8, 125–136.