

# Are language models aware of the road not taken? Token-level uncertainty and hidden state dynamics

Amir Zur<sup>1,2</sup> Atticus Geiger<sup>3,2</sup> Ekdeep Singh Lubana<sup>4,5</sup> Eric Bigelow<sup>4,6</sup>

## Abstract

When a language model generates text, the selection of individual tokens might lead it down very different reasoning paths, making uncertainty difficult to quantify. In this work, we consider whether reasoning language models represent the alternate paths that they could take during generation. To test this hypothesis, we use hidden activations to control and predict a language model’s uncertainty during chain-of-thought reasoning. In our experiments, we find a clear correlation between how uncertain a model is at different tokens, and how easily the model can be steered by controlling its activations. This suggests that activation interventions are most effective when there are alternate paths available to the model—in other words, when it has not yet committed to a particular final answer. We also find that hidden activations can predict a model’s future outcome distribution, demonstrating that models implicitly represent the space of possible paths.

## 1. Introduction

In recent years, Large Language Models (LLMs) have shown impressive capabilities which emerge during next-word prediction (Brown et al., 2020). Despite this high performance on many intelligence benchmarks, LLMs often fail in unexpected and sometimes dramatic ways, for example confidently hallucinating wrong answers or outputting harmful text. Under the surface, when an LLM samples a sequence of text to output, at each token there is a possibility that the LLM might seemingly “change its mind” and say something very different. The goal of this

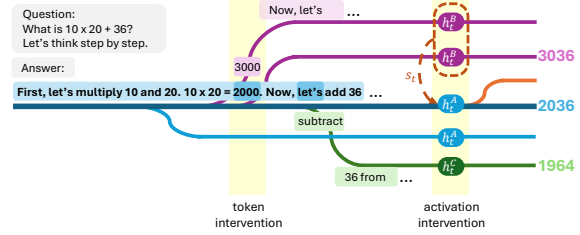


Figure 1. Our experimental set-up. By intervening on the generated tokens, we create branching paths to estimate the model’s outcome distribution. By intervening on the model’s activations, we steer the base generation towards a desired outcome.

work is to study the latent neural representations underlying these token-level uncertainty dynamics, inspired by work evaluating the abilities of LLMs and interpreting their inner workings (Templeton et al., 2024; Panickssery et al., 2023; Ferrando et al., 2024).

Uncertainty estimation is a foundational problem in evaluating and interpreting LLMs (Geng et al., 2024; Xiong et al., 2024). For example, if an LLM “hallucinates” a wrong answer to a question, but the probability of this answer is only 1%, this LLM should be evaluated differently than another LLM that outputs the same wrong answer with 100% confidence. However, the challenge of estimating uncertainty in LLMs is considerably more difficult in settings with long-form text responses, such as with reasoning models where a chain of reasoning text precedes the final answer (Kojima et al., 2022; Jaech et al., 2024; Guo et al., 2025). For each token that an LLM generates, sampling a different token might lead the LLM to generate a different answer. In particular, Bigelow et al. (2025) develop an approach, called Forking Path Analysis, to demonstrate the important role that individual tokens play in determining model certainty.

There is much recent work in AI interpretability which seeks to understand the latent representations in LLMs, and to develop methods for steering or controlling LLMs by intervening on hidden representations. One common approach to interpretability is to study a particular phenomenon in isolation, such as induction heads (Olsson et al., 2022), with small autoregressive language models trained on toy prob-

<sup>1</sup>Department of Linguistics, Stanford University <sup>2</sup>Pr(Ai)<sup>2</sup>R Group <sup>3</sup>Goodfire <sup>4</sup>Physics of Intelligence Group, NTT Research <sup>5</sup>Center for Brain Science, Harvard University <sup>6</sup>Department of Psychology, Harvard University. Correspondence to: Amir Zur <amirzur@stanford.edu>, Eric Bigelow <ebigelow@g.harvard.edu>.

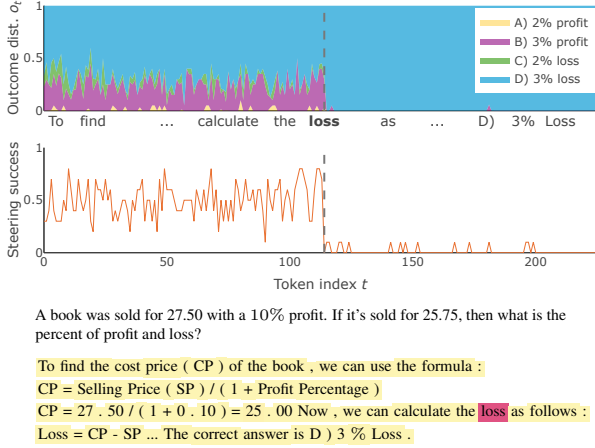


Figure 2. Comparison of the model outcome distribution  $o_t$  (top) and steering success (bottom) across tokens. The outcome distribution and steering success have similar dynamics, with the same change points detected by the CPD algorithm (highlighted text).

lems such as sequences of characters or functional input-output pairs (Akyürek et al., 2022; Xie et al., 2021). An alternative approach to studying LLMs is to instead study in-context learning dynamics and representation learning in pre-trained LLMs, at the level of individual tokens (Park et al., 2025; Bigelow et al., 2024). While these approaches emphasize the role of in-context learning and autoregressive text generation, these aspects of text generation are typically not considered in research on steering. Steering methods strive to construct interventions on hidden representations, such as a vector representation of a concept which can be added to the residual stream of a transformer LLM to steer its behavior (Marks & Tegmark, 2023; Li et al., 2023; Turner et al., 2023; Rimsky et al., 2024). For example, Ardit et al. (2024) construct a vector that makes an LLM more or less likely to refuse answering a user’s question. Fei et al. (2024) develop a method for steering an LLM at the token level, using a smaller LLM as a guide to intervening by replacing specific tokens with alternate strings that lead the LLM down a different reasoning path.

The goal of this work is to study how the hidden representations in LLMs change over the course of in-context learning and next-word prediction, and to consider how model steering relates to token-level uncertainty. More specifically, we apply Forking Paths Analysis (Bigelow et al., 2025) to estimate an LLM’s certainty at each individual token in text generation. We then steer LLMs at different tokens and consider how uncertainty dynamics can predict which points an LLM can or cannot be successfully steered. Finally, we investigate whether token-level uncertainty dynamics can be directly predicted from an LLM’s hidden states.

## 2. Estimating Outcome Distribution with Forking Paths Analysis

We build on the work of Bigelow et al. (2025), who present a method for estimating token-level uncertainty dynamics for black-box neural text generation. The goal of this approach is to understand how, for a given text completion sequence (a *base path*), re-sampling at different tokens or token indices causes changes in the model’s uncertainty.

Their method, Forking Paths Analysis, involves a few main steps: first, given an arbitrary prompt for long-form text generation, such as a multi-hop reasoning question and a chain-of-thought prompt, sample a single *base path* completion  $x^*$ , along with token logit probabilities and top-N alternate token probabilities. In the second step, for each token index  $t$  in the base path, concatenate all tokens in the base path up to  $t$  (i.e.  $x_{<t}^*$ ) along with each top-N alternate token  $w$  (i.e.  $x_t = w$ ), and re-sample  $S$  text completions from the LLM  $x_{>t}^{(s)}$  conditioned on this prompt. Using an answer extraction method, such as concatenating a string “Therefore, the answer is: \_\_\_” and re-prompting the LLM, collect a final answer – or *outcome* – as a one-hot vector  $R(x)$  for each text completion. For the next step, collect these responses  $R(x)$  and token probabilities  $p(x_t = w | x_{<t}^*)$ , and  $p(x_{>t}^{(s)} | x_{<t}^*, x_t = w)$  into a weighted *outcome distribution*  $o_t$  for each token index  $t$ :

$$o_t = \mathbb{E}_{w,s} [R(x_{<t}^*, x_t = w, x_{>t}^{(s)})]$$

The *outcome distribution*  $o_t$  thus represents the LLM’s uncertainty over final answers at each token index  $t$ . In the final step of Forking Paths Analysis, statistical models are used to analyze  $o_t$  and look for specific token indexes  $t$  where the outcome distribution changes suddenly. Bayesian Change Point Detection (CPD) models are used to estimate the probability  $p(\tau = t | o_t)$  that a change point occurs at each token index  $t$ . Bigelow et al. (2025) analyze uncertainty dynamics of GPT-3.5 on a variety of common LLM benchmarks and find many striking cases of “Forking Tokens” where the outcome distribution suddenly shifts from one distribution of certainty to another. A major drawback of this approach, however, is the computational cost – on the order of millions of tokens to analyze a single base path completion. In this paper, we investigate whether hidden states provide meaningful information about the outcome distribution during generation of reasoning text, without needing to generate new tokens.

## 3. Steering Outcome Distribution with Hidden State Interventions

Forking Paths Analysis reveals interesting uncertainty dynamics, where an LLM’s uncertainty can dramatically change upon sampling a single token (see top of Figure 2).

In this work, we analyze how these dynamics are reflected in the hidden representations of an LLM. For instance, can we steer an LLM away from its path once it has decided on an answer? By intervening on linear subspaces in a model’s hidden activations, we investigate whether steering success depends on the particular token where steering occurs, and how this relates to uncertainty dynamics in the model’s output.

**Methods** We apply difference-in-means interventions as in Marks & Tegmark (2023) to steer a model during generation. To steer our model towards a desired outcome  $A$ , we sample  $n = 500$  generations whose final outcome is  $A$  and  $n$  generations whose final outcome is some other answer, denoted as  $\bar{A}$ . Let  $\mathbf{h}_t^{(A)}$  be the list of hidden activations over all the generations leading to answer  $A$  at token  $t$ , and  $\mathbf{h}_t^{(\bar{A})}$  be the list of activations over all generations leading to an alternate outcome. We set up a linear mean-mass probe (Marks & Tegmark, 2023) to create a steering vector  $s_t$ .

$$s_t^{(A)} = \frac{1}{n} \sum \mathbf{h}_t^{(A)} - \frac{1}{n} \sum \mathbf{h}_t^{(\bar{A})}$$

We select the token position  $t$  that best separate between generations leading to  $A$  and ones leading to  $\bar{A}$ . Specifically, we treat  $s_t^{(A)}$  as a linear separator between the sets of activations and construct the classifier  $S_t^{(A)}(x) = \text{sigmoid}(\hat{s}_t^{(A)} \cdot x)$ . We choose the token  $t$  whose corresponding classifier  $S_t^{(A)}$  has the highest classification accuracy on hidden activations from a held-out set of generations.

Difference-in-means steering vectors over later tokens and middle layers ( $t = 200$ ,  $l = 12$  for the example in Figure 2) achieve the highest classification accuracy ( $\sim 80\%$ ). It is worth noting that these token positions do not correspond to change points in the model generation (as described in Section 2) but instead to the end of the generation, when the final answer likely already appears in the generated text.

Let  $s^{(A)}$  be steering vector with the highest classification accuracy for answer  $A$ . At every token position  $t$  in the base path, we apply difference-in-means steering by adding the pre-computed steering vector  $s^{(A)}$  to the activation  $h_t$  at token  $t$  in the base path. We then sample  $k = 10$  continuations from the intervened model, repeatedly adding  $s^{(A)}$  to the activation at every generated token. We measure steering success at each token  $t$  as the number of times out of  $k$  that the steered outcome is answer  $A$ , subtracted by the model’s original outcome distribution  $o_t(A)$  for  $A$  at token  $t$ , i.e., the difference between steering success rate and base success rate without steering. High steering success corresponds to token positions at which we can control generation by intervening on linear subspaces in the model’s activations.

**Data and Models** We analyze the Llama-3.2 3B Instruct LLM prompted with zero-shot chain-of-

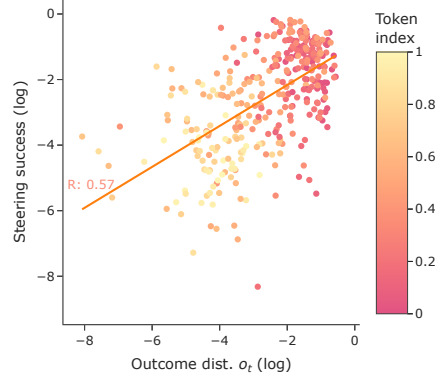


Figure 3. Correlation between steering success ( $y$ -axis) and base outcome probability ( $x$ -axis) across token positions.

thought reasoning (i.e., prefixing “Let’s think step by step” before the model’s completion). Given the computational complexity of Forking Paths Analysis, we consider four examples from three reasoning datasets on which the LLM is uncertain. See Appendix A for more details.

**Results** Our main steering analysis is presented in Figure 2. In this example, sampling token  $t = 114$  (“loss”) drastically shifts the outcome distribution away from “B) 3% profit”, which is the correct answer, and towards “D) 3% loss”, which is the model’s generated answer.

The success rate of steering in this case has very similar dynamics to the outcome distribution. Steering success is relatively high until  $t = 114$ , at which point it abruptly drops down to nearly zero. Given the illustrated example, we hypothesize that controllability (measured by steering success) correlates with uncertainty (measured by the outcome distribution  $o_t$ ). That is, *models are most steerable when they are least certain* about the final outcome.

In Figure 3, we plot the steering success for the answer “B) 3% profit” over the model’s original outcome probability for answer B) across different time steps on a log-log scale. We find a moderate correlation between steering success and base probability ( $R = 0.57$ ), as shown in Figure 3. Broadly speaking, our results suggest that steering is most effective when the model is unsure about its final answer. Hence, steering success is a promising estimate of uncertainty during generation.

**Discussion** The success of difference-in-means steering depends on the model uncertainty during generation. We find a *moderate correlation between controllability and uncertainty*, as measured by steering success and base outcome probabilities, respectively. Hence, uncertainty estimates may be indicative of times during generation at which to intervene on a model’s activations to control its output.

Estimating the uncertainty of a model during generation with Forking Paths Analysis is computationally expensive,

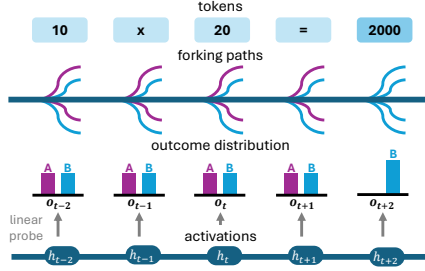


Figure 4. Our experimental set-up for Section 4. At every token position  $t$ , we train a linear probe to predict the distribution of outcomes  $o_t$  from re-sampled paths starting at  $t$ , given the hidden representation  $h_t$  over that token.

since we must simulate the model’s outcomes for branching continuations. Given that hidden state dynamics correlate with a model’s uncertainty dynamics, can we efficiently predict the model’s outcome distribution directly from its hidden states?

#### 4. Predicting Outcome Distribution from Hidden States

Forking paths analysis estimates the distribution of outcomes  $o_t$  by sampling alternate paths that diverge at token  $t$ . This is a computationally expensive process, because it must simulate the LLM’s generation across numerous samples. However, research in faithfulness suggests that the reasoning chain leading to token  $t$  is predictive of the final outcome distribution  $o_t$  (Lanham et al., 2023). In this section, we try to estimate an LLM’s outcome distribution  $o_t$  from the activations over its reasoning process leading up to token  $t$ .

Is the semantic information in chain-of-thought text enough to predict a model’s outcome distribution, or do the model’s hidden activations contain information about its underlying decision-making that doesn’t surface at the token level?

We hypothesize that for a given token  $t$ , the embeddings  $h_t$  from Llama-3.2 3B Instruct and  $h'_t$  from Gemma-2 2B Instruct contain the same *semantic information* about the chain-of-thought text up to token  $t$ , since they are both capable language models. However, only  $h_t$  provides *information about the underlying uncertainty dynamics* of the Llama-3.2 3B Instruct model.

We test our hypothesis by reporting KL losses for linear probes trained to predict the original model’s outcome distribution from either the original model’s or a separate model’s hidden embeddings. Comparable KL losses indicate that model-specific hidden state dynamics aren’t especially useful in predicting its outcome distribution (i.e., they are interchangeable with embeddings from a different LLM). Meanwhile, a lower KL loss for probes trained on the original

model indicate that a model’s hidden states are more predictive of its outcome distribution than the chain-of-thought text alone.

**Methods** Figure 4 illustrates our experimental set-up. For each token index  $t$  in the base path, consider the residual stream activation  $h_t$  over that token. We train a linear probe to predict  $o_t$  from  $h_t$  over a set of token indices. Since  $o_t$  is a full distribution, we train our linear probe with KL divergence loss and report the KL loss on a held-out validation set of token indices.

We also train a linear probe to predict the original model’s outcome distribution  $o_t$  from the text embedding  $h'_t$  of a separate model, Gemma-2 2B Instruct, at token  $t$ . We report the probe’s KL loss on the same validation set of token indices.

**Data and Models** We analyze Llama-3.2 3B Instruct as in Section 3. We report results on 10 randomly sampled examples from the AQUA dataset (Ling et al., 2017), including the example illustrated in Figure 2.

**Results** Figure 5 shows the KL loss for linear probes trained on  $h_t$  and  $h'_t$ , averaged across our 10 data points. Both  $h_t$  and  $h'_t$  are more predictive than a random baseline which always predicts the uniform distribution (1.53 loss, not pictured) and a majority baseline which always predicts a one-hot distribution over the majority class (0.85 loss). Both probes are also most predictive around the middle layers of their respective model, around layers 6-10. This suggests that the reasoning chains have meaningful semantic information that’s predictive of the model’s outcome distribution.

The validation loss for the linear probe trained over the original activations is lower than the loss for the probe trained over activations from a different model (0.11 vs. 0.19 at layer 8). The loss continues to be low for probes trained on the original model’s later layers, while the loss rises for probes trained on the separate model’s later layers. These re-

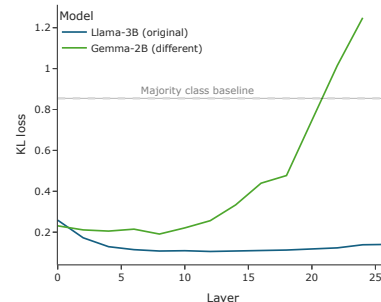


Figure 5. KL loss (lower is better) for linear probes predicting the outcome distribution of Llama from the hidden representations of Llama (blue) and Gemma (green) at the same token mid-generation. Low loss suggests that hidden states over chain-of-thought text are predictive of Llama’s outcome distribution.

sults suggest that the original model’s hidden activations  $h_t$  may carry information that is used to determine the model’s future actions, beyond the information carried by the output tokens alone.

**Discussion** The reasoning chain at token  $t$  and its corresponding hidden activation,  $h_t$ , are both predictive of the model’s outcome distribution  $o_t$ . However, hidden activations capture underlying decision-making information that is more predictive of the final outcome than the embedding of the same text by a different model. Probing a model’s hidden activations is a promising direction for efficiently estimating its outcome distribution during generation.

## References

- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Arditi, A., Obeso, O., Syed, A., Paleka, D., Panickssery, N., Gurnee, W., and Nanda, N. Refusal in language models is mediated by a single direction. *arXiv preprint arXiv:2406.11717*, 2024.
- Bigelow, E., Holtzman, A., Tanaka, H., and Ullman, T. Forking paths in neural text generation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Bigelow, E. J., Lubana, E. S., Dick, R. P., Tanaka, H., and Ullman, T. D. In-context learning dynamics with random binary sequences. *International Conference on Learning Representations (ICLR)*, 2024.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Fei, Y., Razeghi, Y., and Singh, S. Nudging: Inference-time alignment via model collaboration. *arXiv preprint arXiv:2410.09300*, 2024.
- Ferrando, J., Sarti, G., Bisazza, A., and Costa-Jussà, M. R. A primer on the inner workings of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.
- Geng, J., Cai, F., Wang, Y., Koepl, H., Nakov, P., and Gurevych, I. A survey of confidence estimation and calibration in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6577–6595, 2024.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.



- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., et al. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*, 2023.
- Li, K., Patel, O., Viégas, F., Pfister, H., and Wattenberg, M. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36:41451–41530, 2023.
- Ling, W., Yogatama, D., Dyer, C., and Blunsom, P. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 158–167, 2017.
- Marks, S. and Tegmark, M. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *arXiv preprint arXiv:2310.06824*, 2023.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- Panickssery, N., Gabrieli, N., Schulz, J., Tong, M., Hubinger, E., and Turner, A. M. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
- Park, C. F., Lee, A., Lubana, E. S., Yang, Y., Okawa, M., Nishi, K., Wattenberg, M., and Tanaka, H. Iclr: In-context learning of representations. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- Rimsky, N., Gabrieli, N., Schulz, J., Tong, M., Hubinger, E., and Turner, A. Steering llama 2 via contrastive activation addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15504–15522, 2024.
- Templeton, A., Conerly, T., Marcus, J., Lindsey, J., Bricken, T., Chen, B., Pearce, A., Citro, C., Ameisen, E., Jones, A., Cunningham, H., Turner, N. L., McDougall, C., MacDiarmid, M., Freeman, C. D., Summers, T. R., Rees, E., Batson, J., Jermyn, A., Carter, S., Olah, C., and Henighan, T. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- Turner, A. M., Thiergart, L., Leech, G., Udell, D., Vazquez, J. J., Mini, U., and MacDiarmid, M. Steering language models with activation engineering. *arXiv preprint arXiv:2308.10248*, 2023.
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- Xiong, M., Hu, Z., Lu, X., Li, Y., Fu, J., He, J., and Hooi, B. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *International Conference on Learning Representations (ICLR)*, 2024.

## A. Data Selection and Additional Examples

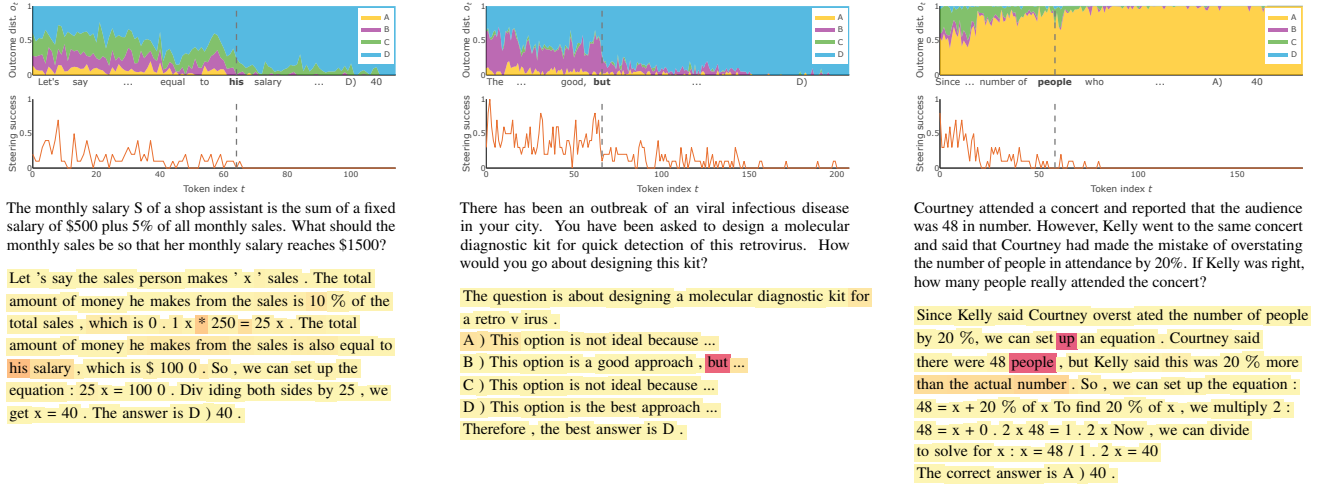


Figure 6. Three additional examples of steering analysis. Each column corresponds to a single example, with the following structure: (Top) outcome distribution  $o_t$  across token positions, estimated by re-sampling completions at alternate tokens for each token index (see Section 2). (Middle) Steering success across token positions, estimated by the number of times a steering vector successfully changes the model’s final answer (see Section 3). (Bottom) answer generated with greedy sampling, with highlighted change points in the outcome distribution  $o_t$  (see Section 2). Across all examples, the outcome distribution and steering success display similar dynamics, with a sharp shift at the highlighted change points.

Due to the computational complexity of Forking Paths Analysis, we consider a handful of examples from the multiple-choice reasoning datasets: **GSM8k**, a collection of grade-school math questions (Cobbe et al., 2021); **AQuA**, a collection of algebraic word problems (Ling et al., 2017); and **GPQA**, a collection of graduate-level science questions (Rein et al., 2024).

We select our data points by sampling from each dataset and keeping questions for which the LLM we analyze (Llama-3.2 3B Instruct) is uncertain. In particular, we sample  $k = 10$  generations from the LLM, and extract the answer from each generation by prompting the model at the end with “What is your final answer?”. We choose examples where the frequency of the most common answer across the 10 examples is between 4 and 6. That is, the LLM is only 40–60% likely to generate that answer. Figure 6 shows the uncertainty and steering dynamics for our selected examples. The steering correlations reported in Section 3 pertain to the example in Figure 2 in the main text. The correlation coefficient for steering success and outcome distribution (on a log-log scale as in Figure 3) averaged across our four examples is  $R = 0.64$ .