

# Soft Self-Consistency Improves Language Model Agents

Han Wang\*    Archiki Prasad\*    Elias Stengel-Eskin\*    Mohit Bansal

UNC Chapel Hill

{hwang, archiki, esteng, mbansal}@cs.unc.edu

## Abstract

Generations from large language models (LLMs) can be improved by sampling and scoring multiple solutions to select a final answer. Current “sample and select” methods such as self-consistency (SC; Wang et al., 2023) rely on majority voting to score answers. However, when tasks have many distinct and valid answers, selection by voting requires a large number of samples. This makes SC prohibitively expensive for interactive tasks that involve generating multiple actions (answers) sequentially. After establishing that majority voting fails to provide consistent gains on such tasks, we demonstrate how to increase success rates by softening the scoring criterion. We introduce *Soft Self-Consistency* (SOFT-SC), which replaces SC’s discontinuous scoring with a continuous score computed from model likelihoods, allowing for selection even when actions are sparsely distributed. SOFT-SC improves both performance *and* efficiency on long-horizon interactive tasks, requiring half as many samples as SC for comparable or better performance. For a fixed number of samples, SOFT-SC leads to a 1.3% increase over SC in absolute success rate on writing bash programs, a 6.6% increase on online shopping (WebShop), and a 4.7% increase for an interactive household game (ALFWorld). Finally, we show that SOFT-SC can be applied to both open-source and black-box models.<sup>1</sup>

## 1 Introduction

The performance of large language models (LLMs) can be greatly improved by generating multiple samples and scoring their answers before making a final selection. One popular and effective “sample and select” approach is *Self-Consistency* (SC; Wang et al., 2023), which leverages chain-of-thought prompting (Wei et al., 2022) to generate

multiple solutions for each input query and then determines the final answer via a majority vote. While SC has demonstrated consistent benefits on question-answering datasets, we find it provides minimal gains in several interactive settings where LLMs act as agents to generate a sequence of actions. SC’s selection mechanism relies on *exact match* in order to tally votes, i.e., it scores answers based on their frequency. However, in interactive domains, multiple distinct and valid answers – in this case, actions – can be generated at each step. This diminishes the effectiveness of SC over actions because the likelihood of generating identical actions decreases as the number of plausible options grows. For instance, a model tasked with predicting bash commands based on user queries has a very large action space (all bash commands) and could generate semantically equivalent commands that differ in their surface form (e.g., `ls -ltr` vs `ls -trl`).<sup>2</sup> Therefore, deriving a signal from voting in LLM-agent domains would require sampling a large number of actions at each step throughout a lengthy trajectory, reducing efficiency and making SC prohibitively expensive (cf. Fig. 1).

We hypothesize that relaxing the strict scoring criterion from votes tallied by exact match to a continuous score will address the shortcomings of SC in two ways: (i) improving *task performance* in sparse action spaces; and (ii) increasing *sample efficiency*, i.e., higher success rates with fewer samples. We propose *Soft Self-Consistency* (SOFT-SC), a continuous relaxation of exact-match sample and select methods. Unlike match-based voting, SOFT-SC handles cases without a *unique* majority answer. Crucially, for a white-box model, SOFT-SC incurs no additional cost and requires no external tests or metrics, as the probabilities used are already produced. Finally, we show that SOFT-SC can be used

\*Equal Contribution

<sup>1</sup>Our code is publicly available at: [https://github.com/HanNight/soft\\_self\\_consistency](https://github.com/HanNight/soft_self_consistency).

<sup>2</sup>For Bash program prediction with five samples, SC fails to produce a single majority action 86% of the time.

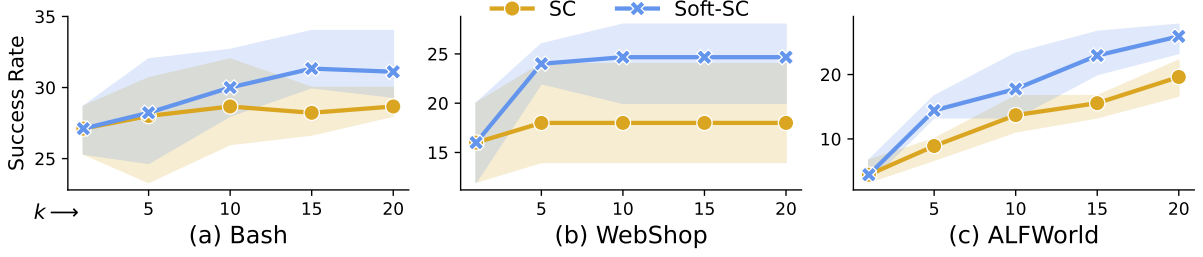


Figure 1: Compared to self-consistency (SC), our method SOFT-SC, exhibits better scaling with respect to the number of samples  $k$ , generally outperforming SC for each  $k$ . We use CodeLlama-34B (Roziere et al., 2023) to compute success rates on the test set of Bash and WebShop. Due to computational cost, for ALFWorld we use Mistral-7B (Jiang et al., 2023) on a 30-task subset of the test set.

to rescore black-box models’ outputs and can be integrated into an efficient variant of SC.

We test SOFT-SC on three diverse interactive domains: Bash (Yang et al., 2023), WebShop (Yao et al., 2022), and ALFWorld (Shridhar et al., 2021).

#### Summary of Key Findings:

1. We demonstrate that SOFT-SC *outperforms* SC with the same number of samples, e.g., by up to 6.6% on WebShop using CodeLlama-34B.
2. SOFT-SC exhibits better sample efficiency i.e., produces better performance than SC with fewer samples (cf. Fig. 1).
3. SOFT-SC scales better with model size than SC, increasing performance by 8.8% on Bash as model size increases from 7B to 70B, as opposed to only 5.8% improvement by SC.
4. SOFT-SC can be combined with smaller LMs to *score generations from black-box models*. We observe that SOFT-SC outperforms SC on closed-source models such as GPT-4 (OpenAI, 2023) by up to 4% on WebShop.

## 2 Methodology

### 2.1 Soft Self-Consistency (SOFT-SC)

Following Wang et al. (2023), for a given input  $\mathbf{x}$  containing the task description, we generate  $k$  solutions using temperature-based sampling (Ackley et al., 1985; Ficler and Goldberg, 2017). To perform selection, we score the action  $y_i$  resulting from each solution using the aggregated probability of the action’s tokens. For an action  $\mathbf{y}$  composed of tokens  $y_1, \dots, y_n$ , we define  $\text{score}(\mathbf{y}) = f(\{P_{\text{LM}}(y_i|y_{<i}, \mathbf{x}) \mid \forall i \in [1, n]\})$  where  $f \in \{\min, \text{mean}, \text{product}\}$ . We choose the aggregation method based on dev set performance. We use mean probability for Bash and ALFWorld and min probability for Webshop. We then choose an action  $\hat{\mathbf{y}}$  with the highest score, i.e.,  $\hat{\mathbf{y}} = \arg \max_{j=1}^k \text{score}(\mathbf{y}_j)$ . Further details and

results for  $f$  options are provided in Appendix A.6.

### 2.2 Adaptive Soft Self-Consistency

To improve efficiency, Aggarwal et al. (2023) introduce adaptive-consistency, which reduces the number of samples ( $k$ ) by approximating the final vote tally per example via sampling discrete vote distributions from a prior and stopping when the samples converge. Instead of sampling from discrete distributions, we choose  $k$  by aggregating likelihood scores until a score threshold  $\tau$  is reached. Following Stengel-Eskin and Van Durme (2023b), we use the minimum probability for comparing with the threshold. We sample one action at a time, stopping when  $\sum_{j=1}^k \min_{i=1}^{|y_j|} P_{\text{LM}}(y_i|y_{<i}, \mathbf{x}) \geq \tau$ , where  $\tau$  is chosen on the dev set (cf. Appendix A.8).

### 2.3 Datasets

We test on three representative English LLM agent datasets; further details can be found in Appendices A.3 to A.5.

**Bash.** We use Yang et al. (2023)’s bash data, which consists of 200 user queries or instructions that can be completed via bash actions. We split these into 50 dev and 150 test. The agent’s performance is measured via success rate. Bash represents a domain with a large action space, as the space of possible bash commands is very large, and many of the queries involve stringing multiple functionalities together into a complex command.

**WebShop.** WebShop (Yao et al., 2022) is a simulated online shopping website environment. Success is measured both by a score  $\in [0, 1]$  reflecting how well the purchased product matches the user’s criteria; the success rate is the rate of perfect scores. Following Zhou et al. (2023), we report performance on a subset of 50 user queries. WebShop also has a large action space, as there are 1.18 million real-world products to select from.

Method	# Samples ( $k$ )	Bash	WebShop	ALFWorld	
		SR	Score	SR	SR
Greedy decoding	1	27.1 $\pm$ 1.7	33.1 $\pm$ 2.8	16.0 $\pm$ 4.0	18.7 $\pm$ 2.1
Self-Consistency (Wang et al., 2023)	10	28.7 $\pm$ 3.1	36.4 $\pm$ 3.3	18.0 $\pm$ 5.3	20.5 $\pm$ 2.9
Adaptive-Consistency (Aggarwal et al., 2023)	[5.0, 7.3] <sup>†</sup>	27.3 $\pm$ 2.4	38.8 $\pm$ 2.4	19.3 $\pm$ 4.2	20.8 $\pm$ 3.2
SOFT-SC	5	28.2 $\pm$ 3.7	44.2 $\pm$ 3.8	24.0 $\pm$ 2.0	22.7 $\pm$ 2.5
SOFT-SC	10	30.0 $\pm$ 2.4	46.0 $\pm$ 6.0	24.6 $\pm$ 4.2	25.2 $\pm$ 3.2
Adaptive SOFT-SC	[5.0, 5.9] <sup>†</sup>	30.0 $\pm$ 2.7	44.5 $\pm$ 4.1	23.3 $\pm$ 2.3	23.9 $\pm$ 2.9

Table 1: Success rates and scores from CodeLlama-34B, averaged across three seeds ( $\pm$  standard deviation). With a fixed  $k = 10$ , SOFT-SC outperforms self-consistency by an average of 4.2%, across datasets. Adaptive sampling uses fewer samples on average than adaptive-consistency while also increasing performance.

<sup>†</sup>Adaptive methods result in differing average  $k$  for each dataset, range reported here.

**ALFWorld.** ALFWorld (Shridhar et al., 2021) is a text-game adaption (Côté et al., 2019) of the embodied ALFRED benchmark (Shridhar et al., 2020) in which an agent performs household chores (e.g., cleaning a mug) via a series of low-level actions. We evaluate on 134 unseen tasks and report the overall success rate. ALFWorld requires agents to generate long action sequences, involving thousands of valid actions at each step for some tasks.

**Metrics.** All these interactive tasks provide a goal and associated environments to execute the LLM-generated actions to accomplish said goal. After executing each action, the environment returns the observation and reward. The observation is a natural language description of the state of the system after executing the action, and the reward indicates if the goal was successfully achieved. The reward can be used to obtain a *success rate*, the percentage of examples with the maximum reward possible. Further details on the rewards for each domain can be found in Appendices A.3 to A.5.

## 2.4 Baselines

We compare SOFT-SC against the following:

**Greedy Decoding.** We sample a single solution with greedy decoding on all datasets; all prompts are given in Appendix C. This is equivalent to both SC and SOFT-SC when  $k = 1$ , as no selection is needed for a single sample.

**Self-Consistency (SC).** We use self-consistency as described by Wang et al. (2023), with majority voting as the selection criterion. We tally votes towards each response using exact match.

**Adaptive-Consistency (AC).** As described in Sec. 2.2, Aggarwal et al. (2023) introduce an adaptive version of SC that improves efficiency by adap-

tively reducing the number of samples. We implement their Beta estimator for all of our settings. Further details can be found in Appendix A.8.

## 3 Results and Discussion

Unless mentioned otherwise, we report average performance on 3 random seeds for each test set.

**For the same number of samples  $k$ , SOFT-SC outperforms SC.** In Table 1, we compare SOFT-SC against the baselines on all datasets using CodeLlama-34B on the test sets. While both SC and SOFT-SC boost performance over the greedy decoding baseline, we find SOFT-SC results in a larger margin of improvement, 8.6% on WebShop (SC only yields 2%). For the same number of samples ( $k = 10$ ), SOFT-SC outperforms SC by 1.3%, 6.6%, and 4.7% (success rate) on Bash, WebShop, and ALFWorld respectively. Comparing the adaptive version of SOFT-SC with Aggarwal et al. (2023), our likelihood-based scores not only improve efficiency by generally using fewer samples, but *also* outperforms AC, e.g., by 4% on WebShop and 3.1% on ALFWorld.

**SOFT-SC exhibits better scaling with  $k$ .** In Table 1, even with  $k = 5$ , SOFT-SC can outperform SC with  $k = 10$ , e.g., with 2.2% improvement on ALFWorld. In Fig. 1, we compare this trend across more values of  $k$ , showing the scaling of SOFT-SC and SC with an increasing  $k$ . We observe that SC provides minimal gains even as  $k$  increases, e.g., on Bash increasing  $k$  from 5 to 20 only yields 1% point improvement in success rate. On the other hand, SOFT-SC consistently improves success rates with  $\sim 3\%$  points improvement as  $k$  goes from 5 to 20. While SC does improve the success rate of Mistral-7B on ALFWorld with increasing  $k$ , SOFT-SC yields greater performance gains us-

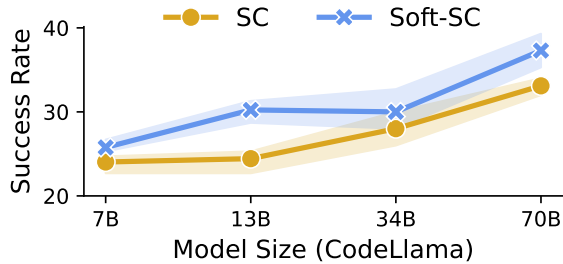


Figure 2: Scaling with model size on Bash (test). SOFT-SC improves over SC for all model sizes.

ing fewer samples, e.g., SOFT-SC with  $k = 5$  is comparable to SC with  $k = 10$ .

**SOFT-SC effectively scales with model size.** As we scale up the size of the LM, we find that SOFT-SC continues to provide improvements over SC. Fig. 2 shows the scaling trends for CodeLlama models ranging from 7B to 70B parameters on Bash and WebShop with a fixed  $k = 10$ . For each LM, SOFT-SC always outperforms SC. Furthermore, SOFT-SC often allows smaller LMs to outperform larger members of the same model class, e.g., CodeLlama-13B with SOFT-SC outperforms CodeLlama-34B with SC. This points to additional efficiency gains from SOFT-SC, as it can allow smaller models to replace larger ones.

**SOFT-SC improves black-box models more than SC.** SOFT-SC requires access to token probabilities to score actions. However, the most performant LLMs are typically black-box models, often with limited or no access to logits (OpenAI, 2023; Pichai, 2023; Anthropic, 2023). In Fig. 3, we study whether (smaller) open-source LMs can be used to score generations from GPT-3.5 and GPT-4. Here, we observe that SOFT-SC offers improvements over SC for a given black-box model, e.g., 4% for GPT-4 on WebShop and 1.8% on Bash when SOFT-SC uses the *same* number of generations from the black-box models as SC. Furthermore, even though Soft-SC requires 2 model calls (one to the black-box model and one to a smaller open-source model), SOFT-SC with  $k = 5$  (total 10 calls) outperforms SC with  $k = 15$  (total 15 calls to the black-box LLM), which shows that our method is significantly more efficient and effective since it can achieve better performance with fewer calls. Note that half of the calls for SOFT-SC are to a 7B model, likely making them much less expensive than calls to the black-box model.

**Calibration is not required for strong SOFT-SC performance.** Given that SOFT-SC selects op-

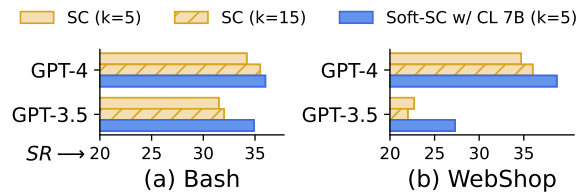


Figure 3: SOFT-SC can be used to score outputs from black-box models on Bash and Webshop (test), improving success rate (SR) over self-consistency.

k	SC	SOFT-SC (logit)	SOFT-SC (verb.)
5	28.0 $\pm$ 4.1	<b>28.2<math>\pm</math> 3.7</b>	27.8 $\pm$ 2.2
10	28.7 $\pm$ 3.1	<b>30.0<math>\pm</math> 2.4</b>	27.6 $\pm$ 2.0

Table 2: Success rates for CodeLlama-34B on Bash with logit-based confidence vs. verbalized (verb.) confidence, averaged across three seeds ( $\pm$  std. dev.).

tions using scores based on token probabilities, we investigate whether a model has to be well-calibrated for SOFT-SC to work. We compute the correlations between two standard calibration metrics – ECE (Naeini et al., 2015) and AUROC – and absolute SOFT-SC performance for CodeLlama-34B across seeds and values of  $k$  on WebShop and Bash test sets. The full plot is shown in Appendix B. We find a moderate negative correlation with AUROC ( $r = -0.55$ ) on Bash and no significant correlation on WebShop); there is no significant correlation for ECE. In other words, having a well-calibrated model is *not* a prerequisite for SOFT-SC. This may be because calibration metrics do not measure *ranking* performance, which is central to our approach.

**Logit-based score outperforms verbalized confidence score.** Recent work has explored prompting language models to express uncertainty or confidence score in human language (Lin et al., 2022; Tian et al., 2023; Xiong et al., 2024). We study whether verbalized confidence scores can be used for selection instead of logit-based scores. We follow Lin et al. (2022) in prompting models to generate verbalized scores, which we then use for selection. As shown in Table 2, verbalized scores perform poorly when used in place of logit-based scores on Bash: Soft-SC with logits outperforms the verbalized method by 2.4% with  $k = 10$ .

## 4 Related Work

**Sample and Select Methods for LLMs.** Ensembling via voting over or aggregating outputs (Breiman, 1996; Freund and Schapire, 1997) can



improve a classifier’s performance. Wang et al. (2023) apply this paradigm to improve LLMs on reasoning tasks, introducing self-consistency (SC). We find that the majority voting used in SC is not suited for LLM-agent domains because the LLM’s generations may not exactly match when the action space is large. Chen et al. (2023b) generalize SC by prompting the LLM to determine consistency. However, LLMs still struggle to determine consistency in interactive domains where the task is partially observable (Ruan et al., 2023). In contrast to SOFT-SC, past work examining re-ranking strategies in code generation (Chen et al., 2022; Li and Xie, 2024) or reasoning (Golovneva et al., 2023; Prasad et al., 2023b) rely on external test cases or model-based metrics to score responses.

**LLM-Agents.** LLMs have proven to be effective agents across a diverse array of multi-step tasks, e.g., mathematical reasoning (Wei et al., 2022), tool-usage (Schick et al., 2023; Qin et al., 2023), robotic navigation (Ahn et al., 2022; Singh et al., 2023), and code-generation (Yang et al., 2023). Standard LLM-agent solutions employ chain of thought prompting (Wei et al., 2022) interleaved with permissible actions within an environment (Yao et al., 2023b). Several follow-up works improve upon this pipeline by building feedback over multiple trials (Shinn et al., 2023), decomposing tasks (Prasad et al., 2023a), or searching over trajectories (Yao et al., 2023a). SOFT-SC is complementary to these approaches, which can be seen as improvements to CoT for a single generation. Note that our work focuses on a single LLM agent (Andreas, 2022) interacting with an external environment to accomplish tasks; this single agent is compatible with other lines of work on discussion among multiple LLM agents (Du et al., 2023; Chen et al., 2023a).

## 5 Conclusion

After establishing the shortcomings of standard voting-based SC in interactive tasks, we introduced SOFT-SC, which relaxes the exact-match scoring function used by SC to a continuous score. On three commonly used interactive benchmarks, we showed that SOFT-SC results in improved performance and increased efficiency. We also show that SOFT-SC is compatible with both white-box and black-box models and that it can be integrated into a more efficient adaptive variant of self-consistency. Finally, we find that a well-calibrated model is not

required for SOFT-SC to work well, and that logits outperform verbalized confidence scores.

## 6 Limitations and Broader Impacts

**Limitations.** In Sec. 1, we pointed out that excessive diversity can lead to failures for SC, as no majority will emerge. However, both SC and SOFT-SC rely on some amount of output diversity: if the model generates  $k$  identical samples, then the output will be no better than generating one. One major motivation for SOFT-SC is efficiency; SOFT-SC substantially improves performance and is able to do so with fewer samples than SC, but it still requires multiple samples from an LLM. Thus, like all sample and select methods, SOFT-SC has a greater cost than greedy decoding. In Sec. 3, we demonstrate that SOFT-SC can be used to rerank outputs from other models that do not consistently provide logits. While SOFT-SC shows major improvements in reranking the outputs of black-box models, it could be applied directly without a smaller scoring model if the generation model’s underlying logits (which exist by design) were made accessible to users.

**Broader Impacts.** Large language models have the potential for negative applications and malicious use (Weidinger et al., 2021; Bommasani et al., 2021). Our work improves LLM performance, meaning it could also be negatively applied. As our work is applied to LLMs operating as agents, it shares the inherent risk of all LLM agent work, namely that the LLM agent could potentially make mistakes and that its actions could lead to negative outcomes for the user. Overall, we believe this risk is mitigated by our use of simulated benchmarks (i.e., no agent we evaluate or develop can affect the world) and by the fact that our work improves agent accuracy, making adverse outcomes less likely.

## Acknowledgements

We thank Justin Chen and Swarnadeep Saha for their valuable help and feedback on the paper. This work was supported by NSF-AI Engage Institute DRL-2112635, DARPA Machine Commonsense (MCS) Grant N66001-19-2-4031, and the Accelerate Foundation Models Research program. The views contained in this article are those of the authors and not of the funding agencies.

## References

- David H Ackley, Geoffrey E Hinton, and Terrence J Sejnowski. 1985. A learning algorithm for boltzmann machines. *Cognitive science*, 9(1):147–169.
- Pranjal Aggarwal, Aman Madaan, Yiming Yang, and Mausam. 2023. [Let’s sample step by step: Adaptive-consistency for efficient reasoning and coding with LLMs](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12375–12396, Singapore. Association for Computational Linguistics.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. 2022. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*.
- Jacob Andreas. 2022. [Language models as agent models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5769–5779, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Anthropic. 2023. [Introducing claude](#).
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Leo Breiman. 1996. Bagging predictors. *Machine learning*, 24:123–140.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. In *The Eleventh International Conference on Learning Representations*.
- Justin Chih-Yao Chen, Swarnadeep Saha, and Mohit Bansal. 2023a. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. *arXiv preprint arXiv:2309.13007*.
- Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Ke-fan Xiao, Pengcheng Yin, Sushant Prakash, Charles Sutton, Xuezhi Wang, and Denny Zhou. 2023b. Universal self-consistency for large language model generation. *arXiv preprint arXiv:2311.17311*.
- Marc-Alexandre Côté, Akos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, et al. 2019. Textworld: A learning environment for text-based games. In *The 7th Computer Games Workshop at the 27th International Conference on Artificial Intelligence (IJCAI 2018)*.
- Yukun Ding, Jinglan Liu, Jinjun Xiong, and Yiyu Shi. 2020. Revisiting the evaluation of uncertainty estimation and its application to explore model complexity-uncertainty trade-off. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 4–5.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multi-agent debate. *arXiv preprint arXiv:2305.14325*.
- Jessica Fidler and Yoav Goldberg. 2017. Controlling linguistic style aspects in neural language generation. In *Proceedings of the Workshop on Stylistic Variation*, pages 94–104.
- Yoav Freund and Robert E Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139.
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. [ROSCOE: A suite of metrics for scoring step-by-step reasoning](#). In *The Eleventh International Conference on Learning Representations*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*.
- Zhenwen Li and Tao Xie. 2024. Using llm to select the right sql query from candidates. *arXiv preprint arXiv:2401.02115*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. *Transactions on Machine Learning Research*.
- Xi Victoria Lin, Chenglong Wang, Luke Zettlemoyer, and Michael D. Ernst. 2018. [NL2Bash: A corpus and semantic parser for natural language interface to the linux operating system](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. 2015. Obtaining well calibrated probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.
- OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Sundar Pichai. 2023. An important next step on our ai journey: Google; 2023 [updated 6 feb 2023].

- Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. 2023a. Adapt: As-needed decomposition and planning with language models. *arXiv preprint arXiv:2311.05772*.
- Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and Mohit Bansal. 2023b. [ReCEval: Evaluating reasoning chains via correctness and informativeness](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10066–10086, Singapore. Association for Computational Linguistics.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, J  r  my Rabin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.
- Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J. Maddison, and Tatsunori Hashimoto. 2023. Identifying the risks of lm agents with an lm-emulated sandbox. *arXiv preprint arXiv:2309.15817*.
- Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2023. Are emergent abilities of large language models a mirage? *arXiv preprint arXiv:2304.15004*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dess  , Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.
- Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 14.
- Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10740–10749.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre C  t  , Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. [ALFWorld: Aligning Text and Embodied Environments for Interactive Learning](#). In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2023. Prog-prompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11523–11530. IEEE.
- Elias Stengel-Eskin and Benjamin Van Durme. 2023a. [Calibrated interpretation: Confidence estimation in semantic parsing](#). *Transactions of the Association for Computational Linguistics*, 11:1213–1231.
- Elias Stengel-Eskin and Benjamin Van Durme. 2023b. [Did you mean...? confidence-based trade-offs in semantic parsing](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2621–2629, Singapore. Association for Computational Linguistics.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pi  ric Cistac, Tim Rault, R  mi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs. In *The Twelfth International Conference on Learning Representations*.
- John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2023. Intercode: Standardizing and benchmarking interactive coding with execution feedback. In *Advances in Neural Information Processing Systems*.



Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. In *Advances in Neural Information Processing Systems*.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*.

Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2023. Language agent tree search unifies reasoning acting and planning in language models. *arXiv preprint arXiv:2310.04406*.

## A Method and Dataset Details

### A.1 Hyperparameters

We select the threshold  $\tau$  on the dev set for both Adaptive-Consistency baseline and Adaptive SOFT-SC. For Adaptive-Consistency baseline, we set the threshold  $\tau$  of 0.8, 0.85, and 0.8 for Bash, WebShop, and ALFWorld respectively. For Adaptive SOFT-SC, we set the threshold  $\tau$  to 0.95, 3.0, and 3.5 for Bash, WebShop, and ALFWorld respectively. Because Adaptive SOFT-SC accumulates minimum probabilities over  $k$  samples for comparing with the threshold, the threshold may be  $\geq 1$ .

For greedy decoding, we use a temperature of 0.7 for all datasets. In case of sampling  $k > 1$  outputs from the model, we set the temperature of open-source models to 0.7 for Bash, 0.9 for WebShop, and 0.9 for ALFWorld, with top-p value of 0.9 and top-k value of 40, and with max\_tokens set to 100. For obtaining generations from the OpenAI API, we use a temperature of 0.7 for Bash, 0.9 for WebShop and ALFWorld and top-p value of 1 for all datasets.

### A.2 Model Checkpoints and Licenses

Webshop, Bash, and ALFWorld all have MIT licenses. CodeLlama is released under a custom permissive license available here: <https://github.com/facebookresearch/llama/blob/main/LICENSE>. Mistral uses an Apache License 2.0. For CodeLlama, we used the CodeLlama-\*b-Instruct checkpoints. For Mistral, we used the Mistral-7B-Instruct-v0.2

checkpoint. All open-source models were accessed via Huggingface Transformers (Wolf et al., 2019). For OpenAI models, we used the gpt-3.5-turbo-0613 and gpt-4 checkpoints. All models were run for inference only with int-8 quantization on Nvidia 40GB A100 GPUs. We will release our code under an MIT license.

### A.3 Bash

Yang et al. (2023) propose an interactive benchmark for evaluating LMs on a bash coding task, created by bootstrapping queries from NLP2Bash benchmark (Lin et al., 2018). The dataset has 200 user queries or instructions that can be completed via bash actions, which we split into 50 dev and 150 test. After each action is executed, the agent observes the corresponding output from the file system. The agent’s performance is measured via success rate, which is determined by a reward function based on modifications to the file system with respect to a gold command as well the latest execution output – a success means the reward is 1.0. For example, given a query “find files in the /workspace directory and sub-directories, that changed within last hour”, the agent generates a corresponding command `find /workspace -cmin -60`.

**Setup.** We focus on the single-turn setting instead of the multi-turn setting because we find the observation (i.e., the execution output of the action) from the Bash environment and the oracle reward rarely helps the agent generate correct commands. In our preliminary experiments, we observed that generating multiple commands using temperature-based sampling under the single-turn setting resulted in a success rate comparable to or even better than the multi-turn setting. Furthermore, in real-world scenarios, it is impossible to obtain oracle rewards to determine whether the generated commands are correct. Therefore, we prompt the LLM with a simple description of the task setting to sample  $k$  commands that would address the query. The final command selected by different methods is executed in the InterCode Bash environment and the response is scored to get the success rate.

**Metric.** After submitting the generated action, the environment returns a reward  $r \in [0, 1]$ . The reward function takes into account the differences in the file system resulting from executing the predicted command and the file system resulting from executing the gold command, as well as the latest



execution output. The *Success Rate* (SR) metric is defined as the proportion of tasks where  $r = 1$ .

#### A.4 WebShop

WebShop (Yao et al., 2022) is a simulated online shopping website environment with 1.18 million real-world products. The underlying task requires an agent to navigate a simulation of a shopping website via a series of commands and buy a suitable product as per the user’s instruction (e.g., 3oz bottle of natural citrus deodorant for sensitive skin under \$30). At the end of the trajectory, the environment returns a numeric score  $\in [0, 1]$  reflecting the degree to which the bought product matches the input criteria. Performance is measured based on the score as well as the success rate (i.e., a perfect score of 1). WebShop also has a large action space, as there are millions of products to select from. We use 30 user queries *not* in the test set to finalize our prompts and thresholds used for adaptive consistency as well as adaptive SOFT-SC.

**Setup.** Following Prasad et al. (2023a), we factorize the underlying agent into two modules: (i) selecting a suitable product, and (ii) buying a selected product. This simulates a “cart” functionality in online shopping. Given a user query, the agent first employs the search functionality and picks a few relevant products from the search page. It then explores the corresponding product page, matches its features, and determines if it can be added to the cart. We prompt the LLM to generate  $k$  such trajectories, potentially adding up to  $k$  products to the cart. In the end, we select a product by majority vote over product IDs and use a separate prompt to get the agent to buy the product while selecting relevant product options such as color, size, etc. The corresponding prompts are shown in Appendix C.

Note that due to the discrete and discontinuous nature of exact match (Schaeffer et al., 2023), SC can only perform selection over products. Given a description, SC navigates through the environment and selects multiple product pages, indexed by their IDs; these IDs can be aggregated via voting. However, within each product page, there are numerous follow-up options that must be selected, and which cannot be voted on as their selection happens across multi-step trajectories. Once a majority product is selected, SC uses a greedy action trajectory based on ReAct (Yao et al., 2023b) to specify the options for a selected product; this often results in suboptimal products being bought, as SC

often picks the default option.

In contrast, the scoring criterion in SOFT-SC allows us to score and select from trajectories to first select products as well as to specify their options and buy them, generating and scoring  $k$  trajectories overall. Thus, SOFT-SC accounts for diversity in each stage and yields higher performance. For example, for the user query “*natural looking long clip in extensions under \$40*” SC tallies votes for products IDs the cart after the product selecting stage: [B09QQLDJ93, B093BKWHFK, B09QQLDJ93], picking the B09QQLDJ93 as it forms a majority. It then uses a greedy ReAct trajectory to select the final options (e.g., the color) and to buy the item. SOFT-SC, on the other hand, can differentiate between action trajectories sampled for buying the *same* product ID, allowing it to distinguish between a final selection that has the default color “pink” and the correct product that uses the color “brown” – resulting in different scores from the environment.

**Metric.** When the LLM agent generates a **buy** action at the end of the trajectory, the environment returns a reward  $r \in [0, 1]$  reflecting the degree to which the bought product matches the input criteria. The *Success Rate* metric is defined as the portion of tasks where  $r = 1$ . The *Score* metric is defined as  $(100 \times \text{avg. reward})$ , which captures the average reward obtained across different task trajectories.

#### A.5 ALFWorld

ALFWorld (Shridhar et al., 2021) is a text-game adaption (Côté et al., 2019) of the embodied ALFRED benchmark (Shridhar et al., 2020). The underlying task requires the agent to perform basic household chores such as finding a mug, cleaning it, and putting it on a countertop via a series of low-level actions (e.g., “go to sink”). After each action, the environment provides textual feedback (e.g., the contents of the cabinet after it is opened). We evaluate on 134 unseen tasks spanning 6 task types and report the overall success rate. In Fig. 1, due to computational requirements of using a larger number of samples, we report performance on a subset of the test split consisting of a total of 30 tasks, picking 5 from each task type. For the dev set, we use a disjoint set of 12 tasks from the ‘valid seen’ split of ALFWorld. This is only used to select the scoring criteria, e.g., mean, min, or product, and the thresholds for the adaptive variants.

**Setup.** Unlike WebShop, tasks in ALFWorld cannot be decomposed uniformly such that each sub-

task is handled by an independent agent without significant planning and communication overhead (Prasad et al., 2023a). For instance, the sub-tasks involved in “putting a clean mug on a counter-top” vary considerably from the sub-tasks involved in “examining a spray-bottle under a desk lamp”. Therefore, in ALFWorld, at each step, we sample  $k$  actions, and for SC perform majority voting over these  $k$  actions. Note that both SOFT-SC and SC only score *actions*, not thoughts or comments generated by the agent to aid in problem-solving. We continue sampling responses until a valid action is reached, skipping “thought” actions (i.e., generations starting with “Think:”) as well as comments. We only allow the selection of actions, ignoring the reasoning generated before the action. Note that both SC and SOFT-SC are more computationally demanding in the case of ALFWorld, since we perform selection over actions at each step, as compared to WebShop, where selection is performed once at the end of the selection phase over products. Following Yao et al. (2023b), the prompt to the LLM includes one in-context trajectory corresponding to a query from the same task type as the test instance.

**Metric.** After each action generated by the LLM agent, the environment provides textual feedback (e.g., the contents of the cabinet after it is opened). The feedback “*You won!*” in addition to reward  $r = 1$  indicates that the agent has completed the task successfully. The *Success Rate* metric is the percentage of tasks where the agent succeeds.

## A.6 Aggregation Methods

For a given input  $\mathbf{x}$  containing the task description and a corresponding sampled action  $\mathbf{y}$  composed of tokens  $y_1, \dots, y_n$ , we can compute  $\text{score}(\mathbf{y})$  using the following probability aggregation methods:

- **Mean:**  $\text{score}(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n P_{\text{LM}}(y_i | y_{<i}, \mathbf{x})$
- **Min:**  $\text{score}(\mathbf{y}) = \min_{1 \leq i \leq n} P_{\text{LM}}(y_i | y_{<i}, \mathbf{x})$
- **Length-Normalized Product:**  $\text{score}(\mathbf{y}) = \exp\left(\frac{1}{n} \sum_{i=1}^n \log P_{\text{LM}}(y_i | y_{<i}, \mathbf{x})\right)$ .

For Bash and ALFWorld, we perform scoring and selection at the action level, where the mean probability serves as an effective measure of the overall confidence in an action being the correct response to a given query. WebShop involves trajectory-level evaluations, where the correctness of a sequence of actions (a trajectory) towards accomplishing a task is assessed. In the case of WebShop, the trajectory

Method	Bash	WebShop	ALFWorld
SC	20.0	22.0	6.70
min	18.0	<b>33.0</b>	10.0
mean	<b>24.0</b>	30.0	<b>16.7</b>
product	22.0	16.7	13.3

Table 3: Dev success rates for one seed across aggregation methods. For Bash and WebShop we use CodeLlama-34B and for ALFWorld we use Mistral-7B.

represents a sequence of actions to *select* a suitable product based on the user query by navigating through a series of webpages; this sequential nature makes min better-suited. We also demonstrate experimental results on dev set for all aggregation methods to validate our explanation in Table 3.

## A.7 Baselines

**Greedy Decoding.** We sample trajectories with greedy decoding on all datasets; prompts are given in Appendix C. For WebShop and ALFWorld, we follow a ReAct prompt format (Yao et al., 2023b) while for Bash we follow the standard format provided by Yang et al. (2023). This is equivalent to both SC or SOFT-SC when  $k = 1$  (since with a single sample, there is no selection needed, making the selection strategy irrelevant).

**Self-Consistency (SC).** We use self-consistency as described by Wang et al. (2023), with majority voting as the selection criterion. We tally multiple votes towards a response only if the model generates the *exact* response multiple times.

## A.8 Adaptive SOFT-SC

To improve sample efficiency, Aggarwal et al. (2023) introduce adaptive-consistency (AC), which reduces the number of samples ( $k$ ) needed for selection by approximating the final vote tally through sampling. Specifically, AC adds generations one at a time (i.e., it increments  $k$  starting from 1) and terminates when a stopping criterion is satisfied or the number of generations has reached the maximum allowed. The stopping criterion is based on samples from a discrete distribution over vote distributions, parameterized by the current vote counts; these samples represent likely future vote distributions given the current trends. If the samples have converged, then further generations are unnecessary. For example, if 5/10 samples have been generated and 4 are identical, then the probability that the next 5 will change the majority vote is vanishingly small, meaning that generating further solutions is

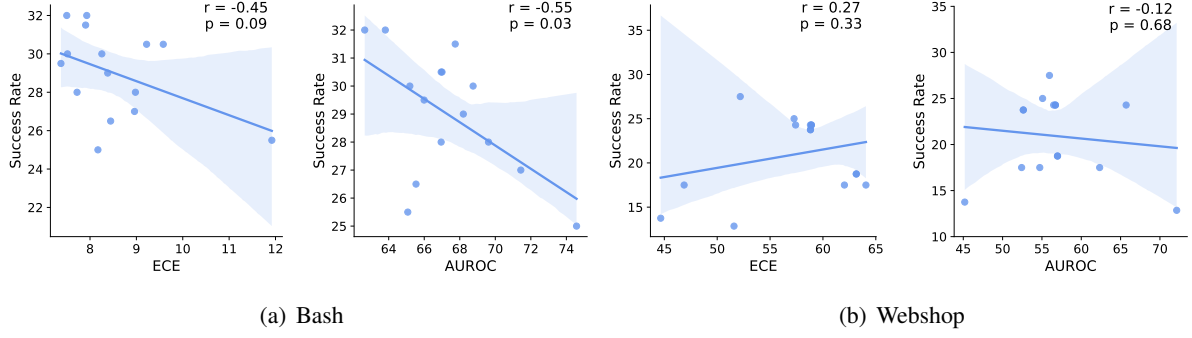


Figure 4: The Pearson correlations between two standard calibration metrics – ECE and AUROC – and SOFT-SC performance for CodeLlama-34B across seeds and values of  $k$  on Bash and Webshop test set.

wasteful. On the other hand, if there is no clear majority winner after 5 samples, further solutions would be needed.

We can apply a similar methodology to SOFT-SC. However, instead of estimating  $k$  by sampling from a discrete vote distribution, we estimate the stopping criterion for sampling by aggregating likelihood scores until a sufficient score threshold  $\tau$  is reached. While we use average probability across tokens for selection, we find that this score is poorly calibrated. Following Stengel-Eskin and Van Durme (2023a), who found minimum token probabilities to be better calibrated, we use the minimum probability for comparing with the threshold. Therefore, we sample actions one-at-a-time and stop when the number of samples  $k$  is such that  $\sum_{j=1}^k \min_{i=1}^{|y_j|} P_{\theta}(y_i | y_{<i}, \mathbf{x}) \geq \tau$ . The threshold  $\tau$  is a domain-specific hyperparameter that we select based on a dev set (discussed in Appendix A.1). Specifically, we set the threshold  $\tau$  to 0.95, 3.0, and 3.5 for Bash, WebShop, and ALFWorld respectively. Note that in this case, the threshold can be  $> 1$  as it represents a threshold on cumulative confidence values, rather a threshold on true probability distribution. This differs from adaptive-consistency, for which the threshold is over a normalized probability, i.e., it must be less than  $\leq 1$ .

## B Calibration

Following past work (Kuhn et al., 2023; Stengel-Eskin and Van Durme, 2023a), we use Expected Calibration Error (ECE) and Area Under the Receiver Operator Characteristic curve (AUROC) to check the calibration of scores used in SOFT-SC:

**Expected Calibration Error (ECE)** (Naeini et al., 2015) is used to quantify how well a model is calibrated. It computes the difference between the accuracy and confidence of the model, where accuracy is averaged across examples falling into

confidence bins. A well-calibrated model will have a low ECE, as it will have a smaller difference between the predicted rate of success (the average confidence) and the actual rate of success (the average accuracy) of a given set of predictions. While ECE is a standard metric, it suffers from sensitivity to the number of confidence bins used (Ding et al., 2020). To mitigate this, we use Stengel-Eskin and Van Durme (2023a)’s implementation of Ding et al. (2020)’s adaptive binning approach, which dynamically adjusts bin sizes to reduce bias in the confidence estimate.

**Area Under the Receiver Operator Characteristic curve (AUROC)** assesses the ability of the estimated confidence to distinguish correct and incorrect samples. AUROC measures the area under the curve formed by comparing the true positive rate to the false positive rate. If a model is well-calibrated, then there is some threshold for which we can separate predictions into correct predictions (above the threshold) and incorrect ones (below the threshold). In general, as we adjust the threshold there will be a tradeoff between true positives and false positives (e.g., a low threshold will result in a large number of false positives, while a high threshold will reduce the number of true positives). A higher AUROC score is better, with a perfect classifier achieving an AUROC of 1 while a random estimator would score 0.5.

Figure 4 illustrates Pearson correlations between two standard calibration metrics – ECE and AUROC – with SOFT-SC performance. For Bash, we find no significant correlation with ECE and a moderate negative correlation with AUROC. For Webshop, neither metric is significantly correlated. Therefore, we conclude that a well-calibrated model is not a prerequisite for SOFT-SC. This may be because calibration metrics do not measure ranking performance, which is central to our approach.

## C Prompts

We provide the prompts along with in-context examples supplied to the LLM for sampling trajectories for Bash and WebShop in Fig. 5, Fig. 6, and Fig. 7. As mentioned in Appendix A.5, for ALF-World, we use the prompts and in-context examples provided in Yao et al. (2023b).



## Bash

```
System: You are a helpful assistant expert specializing in BASH.
User: ## TASK DESCRIPTION
You are a BASH code generator helping me answer a question using BASH.
I will ask you a question, and your task is to interact with a Bourne Shell system using BASH commands
to come up with the answer.

## RESPONSE FORMAT
Your response should be a BASH command. Format your BASH command as follows:
```BASH
Your BASH code here
```

DO NOT WRITE ANYTHING EXCEPT FOR CODE in your response.
Try ```sql
SHOW TABLES``` or ```sql
DESCRIBE <table_name> to learn more about the database```.

## OUTPUT DESCRIPTION
Given your BASH command input, the system will then give back output formatted as follows:

Output: <string>
Reward: [0, 1]

The output is the standard output from executing your BASH command.
The reward is a decimal value between 0 and 1, which tells you how close your BASH command is to the
correct answer.
The closer the reward is to 1, the closer your BASH command is to the correct answer.

You have to try to maximize the reward.

Query: "{query}".
Do not generate any output or reward.
Assistant: {Model Completion}
```

Figure 5: Prompt for Bash tasks.

## WebShop (adding a product to cart or selection)

```
Instruction: Your task is to select a product that matches the user criteria and add it to your 'cart'.
You can use the following commands with the environment:
- search: To search for a query provided you are on a search page with a [search] button
- click: To click on a link enclosed within '[]'
- think: To add comments or thoughts without affecting the state of the webpage
- cart: To add an item to the cart from which it can be bought later.

Here are a few examples:
Criteria: I want to buy a 3 ounce bottle of citrus deodorant for sensitive skin and
clinically tested that is priced lower than $30 and natural.

Action: search[3 ounce bright citrus deodorant sensitive skin]
Observation:
[Back to Search]
Page 1 (Total results: 50)
[Next >]
[B078GWRC1J]
Bright Citrus Deodorant by Earth Mama | Natural and Safe for Sensitive Skin, Pregnancy and
Breastfeeding, Contains Organic Calendula 3-Ounce
$10.99
[B078GTKVXY]
Ginger Fresh Deodorant by Earth Mama | Natural and Safe for Sensitive Skin, Pregnancy and
Breastfeeding, Contains Organic Calendula 3-Ounce
$10.99
[B08KBVJ4XN]
Barrel and Oak - Aluminum-Free Deodorant, Deodorant for Men, Essential Oil-Based Scent, 24
Hour Odor Protection, Cedar & Patchouli Blend, Gentle on Sensitive Skin (Mountain Sage, 2.7
oz, 2-Pack)
$15.95

Action: think[B078GWRC1J and B078GTKVXY are bright citrus deodorant less then 50 dollars. I
can check B078GWRC1J first to check if it matches all my requirements.]
Observation: OK.

Action: click[B078GWRC1J]
Observation:
[Back to Search]
[< Prev]
scent [assorted scents][bright citrus][calming lavender][ginger fresh][simply non-scents]
size [travel set (4-pack)][3 ounce (pack of 1)][3-ounce (2-pack)]
Bright Citrus Deodorant by Earth Mama | Natural and Safe for Sensitive Skin, Pregnancy and
Breastfeeding, Contains Organic Calendula 3-Ounce
Price: $10.99
Rating: N.A.
[Description]
[Features]
[Reviews]
[Attributes]
[Buy Now]

Action: think[I should explore the features of this product to verify it matches all my
requirements.]
Observation: OK.

Action: click[Features]
Observation:
[Back to Search]
[< Prev]
NEW from Earth Mama (formerly Earth Mama Angel Baby), formulated especially for pregnancy,
breastfeeding and sensitive skin
Contains organic grapefruit, tangerine and calendula
NO propylene glycol, artificial fragrance, parabens or aluminum
Dermatologist tested and clinically tested for irritation
Better than natural organic! NSF/ANSI 305 Certified by Oregon Tilth

Action: think[B078GWRC1J matches my requirements: natural, citrus deodrnt, for sensitive
skin, clinincally tested, and price less than $30. Thus, I should put it in my cart.]
Observation: OK.

Action: cart[B078GWRC1J]
Observation: OK.

Criteria: {}
```

Figure 6: Prompt for selecting and adding products to the cart in WebShop.

## WebShop (buying a product)

```
Instruction: Buy product [B078GWRC1J] that matches criteria: 3 ounce bottle of bright
citrus deodorant for sensitive skin, and price lower than 50.00 dollars
You can use the following commands with the environment:
- search: To search for a query provided you are on a search page with a [search] button
- click: To click on a link enclosed within '[]'. E.g., to buy use 'click[Buy Now]'
- think: To add comments or thoughts without affecting the state of the webpage

Action: think['I should be on the search page where I can see a list of products, I should
verify that.']
Observation: OK.

Action: load
Observation:
[Back to Search]
Page 1 (Total results: 50)
[Next >]
[B078GWRC1J]
Bright Citrus Deodorant by Earth Mama | Safe for Sensitive Skin, Pregnancy and
Breastfeeding, Contains Organic Calendula 3-Ounce
$10.99
[B08KBVJ4XN]
Barrel and Oak - Aluminum-Free Deodorant, Deodorant for Men, Essential Oil-Based Scent, 24
Hour Odor Protection, Cedar & Patchouli Blend, Gentle on Sensitive Skin (Mountain Sage, 2.7 oz, 2-Pack)
$35.95
[B078GTKVXY]
Ginger Fresh Deodorant by Earth Mama | Natural and Safe for Sensitive Skin, Pregnancy and
Breastfeeding, Contains Organic Calendula 3-Ounce
$10.99

Action: think['I can spot product [B078GWRC1J] and it is a clickable link since its
enclosed in square brackets. I should click on it.]
Observation: OK.

Action: click[B078GWRC1J]
Observation:
[Back to Search]
[< Prev]
scent [assorted scents][bright citrus][calming lavender][ginger fresh][simply non-scents]
size [travel set (4-pack)][3 ounce (pack of 1)][3-ounce (2-pack)]
Bright Citrus Deodorant by Earth Mama | Natural and Safe for Sensitive Skin, Pregnancy and
Breastfeeding, Contains Organic Calendula 3-Ounce
Price: $10.99
Rating: N.A.
[Description]
[Features]
[Reviews]
[Buy Now]

Action: think['I'm already on the product page. I see multiple possible options enclosed in
'[]' for scent and size for this product. To match the input criteria, I should click on
options '[bright citrus]' for scent and '[3 ounce (pack of 1)]' for size one by one and
then buy in the end.]
Observation: OK.

Action: click[bright citrus]
Observation: You have clicked bright citrus.

Action: click[3 ounce (pack of 1)]
Observation: You have clicked 3 ounce (pack of 1).

Action: think[My task is to buy the product, for it should to click 'buy now']
Observation: OK.

Action: click[Buy Now]
Observation: You have clicked buy now.

Action: think[I finished buying the product. Task completed!]

Here is another task in which you need to buy a product. When you finish buying the product
with the most relevant choices, use 'think[Task completed]'. If you cannot find the
matching options or proceed, think['Task failed']. Note that you can only click on text
enclosed in '[]' on the webpage. Everything else is only a description, not valid with
"click" action.

Instruction: Buy product [{}] that matches the criteria: {}
```

Figure 7: Prompt for buying products in WebShop.