Metrics for Holistic Evaluation of LLM Reasoning about Action, Change, and Planning

Anonymous Author(s)

Affiliation Address email

Abstract

Planning, reasoning, and sequential decision-making have played a pivotal role in the development of AI systems. While Large Language Models (LLMs) have demonstrated impressive capabilities, their evaluation for planning and Reasoning about Action and Change (RAC) problems is performed using strict binary success criteria, which limits information for further analysis and development. Given the probabilistic and autoregressive nature of LLMs, this work proposes the use of simple non-binary task-specific metrics for the evaluation of LLM responses for planning and reasoning tasks that go beyond perfect matching with ground truth, by utilizing set comparison methods, while still maintaining rigid and non-malleable evaluation criteria. We demonstrate the utility and usefulness of this type of metric in obtaining richer data fidelity and information about the quality, precision, nature of LLMs' responses, and their closeness to the ground truth through evaluations on six different tasks across two domains. With two case study examples, we additionally demonstrate the feasibility of comparative analysis of different task-specific data distributions obtained through this metric.

1 Introduction

2

3

5

6

7

10

11

12

13

14

15

The ability to plan, perform sequential decision-making, and reason about action and change is 17 one of the fundamental tenets of human intelligence, and has been one of the cornerstones of AI. Today, modern generative AI and Large Language Models (LLMs) are useful for a plethora of 19 applications, from question answering and document summarization to code generation [5]. Despite 20 their impressive capabilities, LLMs have shown significant limitations in planning, reasoning, and 21 decision-making, particularly in autonomous applications [9, 20, 8, 6]. Such limitations in LLMs' 22 performance are noted through task evaluations that utilize binary success criteria metrics that 23 involve comparison with ground truth answers obtained by automated solvers, planners, or validators. 24 However, there could exist useful information about the quality and precision of the models' responses 25 for these task evaluations, which is not necessarily captured by standard binary metrics.

As LLMs are probabilistic models and generate tokens in an autoregressive manner, it is perhaps 27 not surprising that they struggle to perform accurately on Reasoning about Action, Change (RAC), 28 and planning problems. However, by considering intersection over union (IoU) metrics for task evaluations, we find a more nuanced picture of these models' task performance than is elicited by 30 standard binary success metrics. Specifically, our proposed metrics elicit more information about 31 LLMs' task performance, related to precision and quality, that is missed when applying standard binary success criteria as overviewed in Figure 3. Having information about how close a model is to optimal or expected task performance can be extremely useful for failure analysis, causal analysis, 34 and to make decisions about how best to utilize the model in architectural frameworks such as 35 Auto-ToS[2], LLM-Modulo [8], and other finetuning or prompting setups to enhance performance.

- 37 In the next section, we review benchmarks and related works that evaluate LLMs on Planning and
- 38 RAC tasks, briefly detailing the tasks and metrics used. Then, we outline our evaluation domains,
- 39 proposed metrics, and tasks. Finally, we discuss the results, utility, and usefulness of our metrics for
- 40 RAC and Planning tasks through two examples.

2 Background & Related Works

42 2.1 Related Works

- Recognizing the importance of benchmarking and evaluating the planning, decision-making, and reasoning abilities of LLMs, various benchmarks have been proposed in the literature [19, 6, 7, 9]. He et al. propose the Textual Reasoning about Action and Change (TRAC) benchmark, with 4 Reasoning about Action and Change (RAC) tasks such as projection, action executability, plan verification, and goal recognition, evaluated in the Planning Domain Definition Language (PDDL) based Blocksworld planning domain [7]. They pre-train and evaluate transformer models such as GPT-2 [14] on TRAC, and find that they struggle to generalize to scaling of objects, action sequence lengths, and composite tasks. The evaluations are conducted in a standard binary (true/false) manner and the overall accuracies are computed. However, it is unclease if the task design maintains structural
- and the overall accuracies are computed. However, it is unclear if the task design maintains structural validity (measurement reflecting the internal structure of the construct) [17].
- Valmeekam et al. developed PlanBench, a PDDL-based planning benchmark suite with 8 planningrelated tasks, such as plan generation, cost-optimal planning, plan verification, goal recognition,
 replanning, plan reuse, reasoning about actions and effects, and plan generalization [19]. the
 PlanBench work evaluates LLMs like GPT-4 [1] and Instruct-GPT-3 [13] on their generated plans
 across Blocksworld and Logistics domains, with a primary focus on variants of planning tasks and a
 limited focus on RAC tasks. The evaluations are performed based on the standard binary plan success
 criteria, as has been used in automated planning [16, 4].
- Another notable benchmark is ActionReasoningBench, which evaluates multiple LLMs on RAC tasks such as state tracking, fluent tracking, action executability, and composite question combinations, on 8 different classical planning competition domains [3] like Blocksworld [6]. The evaluation is performed on binary and free-response answers of LLMs, for a few fixed sequence lengths of actions. However, it is important to note here that the free response questions were evaluated using a Llama-70B model in an LLM-as-a-judge framework in order to make the evaluation scalable, potentially leading to inaccurate reporting of performance statistics [21].
- More recently, Kokel et al. proposed ACP Bench that consists of binary and multiple-choice questions on 7 different atomic reasoning and planning tasks, such as reasoning about applicable actions, atom reachability, action reachability, plan verification, progression, landmarks, and plan justification. They perform comprehensive evaluations on various LLMs on multiple classical planning domains, including the Alfworld household domain [18] and a novel 'swap' planning domain [9]. Following this work, Kokel et al. performs evaluations on the generative response version of this dataset, where task-specific evaluations use binary success metrics with perfect matching criteria against stored ground truth answers [10], which may lead to low or unclear construct validity [17].

75 2.2 Domains

To demonstrate the utility of our proposed benchmarks, we utilize standard IPC planning domains [3] 76 77 such as Blocksworld and Depots for our experiments to evaluate the planning and action reasoning abilities of LLMs. For each of the 500 problems in the two domains, we create natural language 78 templates for the initial and goal states, and questions for each of the 6 tasks, resulting in approximately 79 6000 questions that we use to evaluate the Llama 8B and Llama 70B models. For each problem, all 80 the 6 task questions have the same object complexity, initial state, and goal state, only differing in the 81 question prompt. A common natural language context containing the domain description, initial state 82 description and goal state description (if necessary) is utilized for evaluating the LLMs, to ensure as 83 holistic an evaluation as possible.

Blocksworld: Blocksworld is a domain where blocks can be placed on top of each other or on the table. There is one robotic arm that can move the blocks. The goal is to rearrange the blocks from an initial configuration to a goal configuration. This can be challenging as there may be interactions between subgoals. For our evaluation, we design a challenging dataset of 500 problems with 3-12

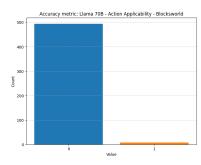


Figure 1: Llama 70B Performance with Standard binary success metric on Action Applicability in Blocksworld; Accuracy = 0.014%; Model's Responses are correct on only 7/501 problems.

97

98

100

101

102

103

104

105 106

107

108

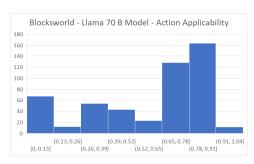


Figure 2: Llama 70B Performance with IoU metric on Action Applicability Task in Blocksworld; This right-skewed distribution provides information on the precision of the model's responses. We can see that the model is close to correctness on around 200/501 problems.

Figure 3: Comparison of IoU Metric vs Standard Binary Success metric. We get a lot more data fidelity and information about precision and quality of responses from the IoU metric compared to the binary success metric.

blocks, that have non-neutral initial states (A subset of blocks are in a stack, and the problems require unstacking and re-stacking), with an average optimal plan length of 18.7 actions.

Depots: The Depots domain is a combination of the blocksworld and logistics domains. In this domain, trucks can transport crates between places, the crates can be stacked onto pallets using hoists, and crates can be loaded into and unloaded from trucks using hoists. This domain inherits the challenges of subgoal interactions from Blocksworld, and reasoning about unreachable actions and states from Logistics. In this domain, we maintain the same object complexity (18) across all problems of the dataset, with an average optimal plan length 12 actions.

3 Tasks: Reasoning about Action, Change, and Planning

Drawing from the above benchmarks in Section 2, we select a set of key atomic tasks, such as action applicability, state tracking, progression of effects, and optimal plan generation, along with a new atomic task called State Comprehension (each task is detailed below). We focus on evaluating LLMs on free-response answers to task questions, instead of multiple-choice and binary responses, in order to obtain better construct validity and avoid construct confounds [15, 17].

Additionally, we formulate a simple non-binary task-specific metric for evaluation of RAC and planning tasks: we compute the Intersection over Union (IoU) of LLM answers and ground truth answers as shown in equation 1, resulting in task-specific metrics as shown in Table 1. Unlike binary evaluation metrics that have a success/ failure criterion based on perfect matching with ground truth answers, this metric allows us to obtain information about the quality of LLMs' performance for each task.

$$Task \ Metric = \frac{LLM \ Answers \cap Ground \ Truth \ Answers}{LLM \ Answers \cup Ground \ Truth \ Answers} \tag{1}$$

The tasks are detailed as follows (with extended descriptions available in Appendix B):

Action Applicability: One of the fundamental atomic RAC tasks is the ability to reason about applicable actions at a given state. We evaluate the generative free responses of LLMs by asking the LLM to list the applicable actions in a given state, provided the common context, as mentioned in section 2.2, using the IoU evaluation metric shown in equation 1 and table 1.

State Comprehension: This task is on simply understanding the given state, such as all the objects, predicates associated with their properties, and the environment properties. Thus, this requires the LLM to provide all the predicates associated with a given state, given the common context 2.2.

Table 1: IoU Task Evaluation Metrics Summary. (GT: Ground Truth)

Task	Resulting Evaluated Formula
Action Applicability	# Correct LLM Answered Actions # LLM Answered Actions U# GT Applicable Actions
State Comprehension	# Correct LLM Answered Predicates Total LLM Answered Predicates UGT Predicates
Progression (Positive/ Negative)	# Correct LLM Answered Effects Total LLM Answered Effects UGT Effects
State Tracking	# Correct LLM Answered Predicates Total LLM Answered Predicates GT Predicates
Optimal Plan Generation	$1 - \frac{\text{\#Overlapping Unique Actions}}{\text{All Unique LLM Actions} \cup \text{Unique Actions from GT Plan}}$

- Progression: This task evaluates the LLMs' ability to understand the effects of an action on the state.
 We design two separate atomic tasks asking the LLM for the positive and negative effects of a single action, respectively, given the common context.
- State Tracking: State tracking is the ability to track entire states across multiple time steps after executing a sequence of actions. We design an atomic version of this task by asking LLMs to provide the complete set of predicates that represent the final state after performing a sequence of two actions.
- Optimal Plan Generation: Plan generation is a classical planning task where the task is to provide a valid sequence of actions that can be executed consecutively from a given state to reach the goal state. If actions have costs, then an optimal plan is one that has the minimum cost. We prompt the LLMs to provide optimal plans given the domain, state, and goal context. Evaluation is performed using the well-known Action Distance metric [12], as shown in Table 1 and detailed in Section B.5.

4 Results and Discussion

128

In this work, we perform evaluations with 6 tasks (considering progression effects as two tasks) across 129 two domains of 500 problems each, on two instruction-tuned pretrained LLMs, using informative 130 task-specific IoU metrics. In Figure 2, we can see that the data distribution obtained through the 131 IoU metric provides us with information on the precision, quality, and nature of models' responses 132 that are entirely missed by binary success metrics, as shown in Figure 1. The right-skewness of the 133 distribution demonstrates that the model is much closer to being correct than the 0 values for 494 134 samples imply. This information is extremely beneficial for compute-intensive and cost-incurring 135 decisions such as finetuning procedures, and the design of future experiments to understand and 136 improve specific atomic reasoning constructs such as action applicability. 137

In figure 6, we compare the IoU metric performance graphs of action applicability and state comprehension tasks of Llama 8B model from the Depots domain. From the stark contrast in the skewness of the distributions, it is pretty clear that the quality and precision of the model's responses for state 140 comprehension are much better than its ability for reasoning about applicable actions. Also, the 141 spread of the distribution for the Action applicability task, according to figure 7, indicates that the 142 model's responses are less precise and more fuzzy compared to those of State comprehension in the 143 Depots domain. Thus, the IoU metric can potentially provide discriminant validity [17], where the 144 evaluation helps differentiate between constructs that should be distributions can be 145 compared with those of State Tracking over 2 actions, shown in Figure 11, which has a slightly lesser height, but a more chaotic spread, which can provide information about the model's reasoning ability 147 with reference to the domain-specific state properties. 148

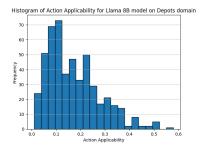
Thus, the IoU metric is beneficial in reasoning and planning tasks, to obtain information on the precision, quality, nature of models' responses, and their closeness to ground truth. We have demonstrated the utility of the metric through evaluations and comparative examples across two domains. A more in-depth correlational analysis across tasks and domain-specific investigations that are beyond the scope of this project is left for future work.

154 References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv
 preprint arXiv:2303.08774, 2023.
- Daniel Cao, Michael Katz, Harsha Kokel, Kavitha Srinivas, and Shirin Sohrabi. Automating thought of
 search: A journey towards soundness and completeness (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 29328–29330, 2025.
- [3] IPC competition. ICAPS International Planning Competition (IPC). http://www.icaps-conference.
 org/competitions/, 2024. Accessed: 09/2025.
- [4] Malik Ghallab, Dana Nau, and Paolo Traverso. Acting, Planning, and Learning. Cambridge University
 Press, 2025.
- [5] D. Hagos, Rick Battle, and Danda B. Rawat. Recent advances in generative ai and large language models:
 Current status, challenges, and perspectives. *IEEE Transactions on Artificial Intelligence*, 5:5873–5893,
 2024. URL https://api.semanticscholar.org/CorpusId:271329267.
- 168 [6] Divij Handa, Pavel Dolin, Shrinidhi Kumbhar, Tran Cao Son, and Chitta Baral. Actionreasoningbench: Reasoning about actions with and without ramification constraints. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=NUDO3NBDOE.
- [7] Weinan He, Canming Huang, Zhanhao Xiao, and Yongmei Liu. Trac: A textual benchmark for reasoning about actions and change. *arXiv preprint arXiv:2211.13930*, 2022.
- [8] Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri,
 Lucas Saldyt, and Anil Murthy. Position: Llms can't plan, but can help planning in llm-modulo frameworks.
 In Proceedings of the 41st International Conference on Machine Learning, ICML'24. JMLR.org, 2024.
- [9] Harsha Kokel, Michael Katz, Kavitha Srinivas, and Shirin Sohrabi. Acpbench: Reasoning about action,
 change, and planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages
 26559–26568, 2025.
- 179 [10] Harsha Kokel, Michael Katz, Kavitha Srinivas, and Shirin Sohrabi. Acpbench hard: Unrestrained reasoning about action, change, and planning. In AAAI 2025 Workshop LM4Plan, 2025. URL https://openreview.net/forum?id=cfsVixNuJw.
- 182 [11] Anagha Kulkarni, Yantian Zha, Tathagata Chakraborti, Satya Gautam Vadlamudi, Yu Zhang, and Sub-183 barao Kambhampati. Explicablility as minimizing distance from expected behavior. *arXiv preprint* 184 *arXiv:1611.05497*, 2016.
- Tuan Anh Nguyen, Minh Do, Alfonso Emilio Gerevini, Ivan Serina, Biplav Srivastava, and Subbarao Kambhampati. Generating diverse plans to handle unknown and partially known user preferences. Artificial Intelligence, 190:1–31, 2012. ISSN 0004-3702. doi: https://doi.org/10.1016/j.artint.2012.05.005. URL https://www.sciencedirect.com/science/article/pii/S0004370212000707.
- [13] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang,
 Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with
 human feedback. Advances in neural information processing systems, 35:27730–27744, 2022.
- 192 [14] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [15] Anka Reuel-Lamparth, Amelia Hardy, Chandler Smith, Max Lamparth, Malcolm Hardy, and Mykel J
 Kochenderfer. Betterbench: Assessing ai benchmarks, uncovering issues, and establishing best practices.
 Advances in Neural Information Processing Systems, 37:21763–21813, 2024.
- [16] Stuart Russell, Peter Norvig, and Artificial Intelligence. A modern approach. Artificial Intelligence.
 Prentice-Hall, Egnlewood Cliffs, 25(27):79–80, 1995.
- 199 [17] Olawale Salaudeen, Anka Reuel, Ahmed Ahmed, Suhana Bedi, Zachary Robertson, Sudharsan Sundar, Ben
 200 Domingue, Angelina Wang, and Sanmi Koyejo. Measurement to meaning: A validity-centered framework
 201 for ai evaluation. *arXiv preprint arXiv:2505.10573*, 2025.
- 202 [18] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew
 203 Hausknecht. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In
 204 Proceedings of the International Conference on Learning Representations (ICLR), 2021. URL https:
 205 //arxiv.org/abs/2010.03768.

- [19] Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati.
 Planbench: An extensible benchmark for evaluating large language models on planning and reasoning
 about change. Advances in Neural Information Processing Systems, 36:38975–38987, 2023.
- [20] Karthik Valmeekam, Sarath Sreedharan, Matthew Marquez, Alberto Olmo, and Subbarao Kambhampati.
 On the planning abilities of large language models (a critical investigation with a proposed benchmark).
 arXiv preprint arXiv:2302.06706, 2023.
- 212 [21] Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang
 213 Sui. Large language models are not fair evaluators. *ArXiv*, abs/2305.17926, 2023. URL https://api.
 214 semanticscholar.org/CorpusID:258960339.
- [22] Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su.
 Travelplanner: A benchmark for real-world planning with language agents. In *Forty-first International Conference on Machine Learning*, 2024.

218 A Task Evaluation Comparison Example



Histogram of State Description for Llama 8B model on Depots domain

Figure 4: Llama 8B Performance with IoU metric on Action Applicability in Depots domain;

Figure 5: Llama 8B Performance with IoU metric on State Comprehension/ Description Task in Depots domain

Figure 6: Comparison of IoU Metric evaluation of Action Applicability and State Comprehension tasks. It is evident from the left-skewed distribution of Figure 4 and the right-skewed distribution of Figure 5 that Llama 8B model's responses and performance is more precise and of higher quality for state comprehension than for reasoning about applicable actions.

219 B Extended Task Descriptions

220 B.1 Action Applicability

221 222

223

224 225

227

228

229

230

231

232

One of the fundamental atomic RAC tasks is the ability to reason about applicable actions at a given state. Previous works have shown that LLMs fall short of this ability and tend to provide invalid or hallucinated actions [22, 9, 6]. For actions to be valid in a given state, specific preconditions required by those actions must hold. We evaluate the generative free responses of LLMs by asking the LLM to list the applicable actions in a given state, provided the common context, as mentioned in section 2.2, using the IoU evaluation metric shown in equation 1 and table 1.

B.2 State Comprehension

A fundamental requirement of reasoning about actions, change, and planning is to simply understand the given state, such as all the objects, predicates associated with their properties, and the environment properties. It is impossible to accurately perform any higher-level reasoning task, such as state tracking, action applicability, or planning, without fully comprehending the properties of the current state. We ask the LLM to provide the list of predicates that represent the current state, giving the domain and state description, and available predicate information as context. Note that the task still involves some basic inferences about state properties from the generic domain description.

B.3 Progression

235

- This task evaluates the LLMs' ability to understand the effects of an action on the state. Keeping track of effects and changes through multiple states and action sequences is an important aspect of sequential decision-making and planning. LLMs have been shown to struggle with tracking changes across sequences of actions and states [6, 9, 20]. Also, prior works have found that LLMs' performance differs with positive and negative predicates [6]. We design two atomic tasks for tracking the positive and negative effects of a single action, given the domain description, current state description, and the available predicates (that can be used to represent effects on states).
- Positive Effects Positive effects are those that are not true in the current state and become true in the following state after the action is performed. These are also called add effects. Identifying positive effects is important as emerging effects can be preconditions to future actions along a plan.
- Negative Effects Negative effects are those that are true in the current state and become false in the following state after the action is performed. These are also called delete effects. Identifying negative effects is extremely important to avoid dead loops, inconsistent states, and ruling out invalid actions.

248 B.4 State Tracking

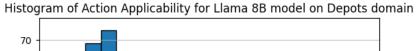
State tracking is the ability to track entire states across multiple time steps after executing a sequence of actions.

State tracking is a fundamental ability required for planning, as it involves generating valid successor states and actions at every visited state. Similar to Handa et al.'s ActionReasoningBench, we design an atomic version of this task by asking LLMs to provide the complete set of predicates that represent the final state after performing an action or sequence of actions.

254 B.5 Optimal Plan Generation

- Plan generation is a classical planning task where the task is to provide a valid sequence of actions that can be executed consecutively from a given initial state to reach the goal state. If actions have costs, then an optimal plan is one that has the minimum cost. Unlike the other RAC tasks, the expected answer here is an ordered and optimal set of actions. This inherently implies a stricter evaluation criterion and, hence, is also more complex, as it requires coming up with optimal, goal-reaching actions, in addition to generating valid plans.
- Evaluation Unlike for previous tasks, there are already various proposed metrics in the planning literature to measure plan quality, such as Action Distance, Causal-Link Distance, and State Sequence Distance [12, 11]. These metrics have been used to measure the quality of plans compared to an optimal plan. As LLMs are probabilistic models and fare poorly at generating valid plans [8], utilizing such metrics can shed some light on their performance at generating plans that would not be available with perfect accuracy measures. Hence, we utilize the action distance metric for our evaluation. However, it is important to note that action distance is a set comparison metric between unique action sets and does not account for the ordering of actions.

267 C Tasks Performance Graphs for IoU metric on Depots Domain



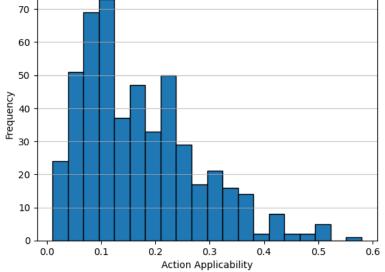


Figure 7: Llama 8B Performance on Action Applicability in Depots Domain

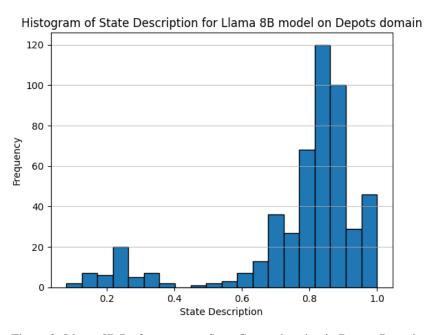


Figure 8: Llama 8B Performance on State Comprehension in Depots Domain

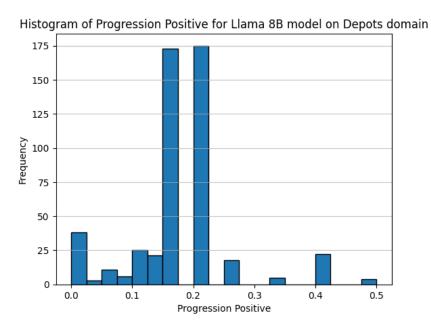


Figure 9: Llama 8B Performance on Identifying Positive Effects of Action progression in Depots Domain

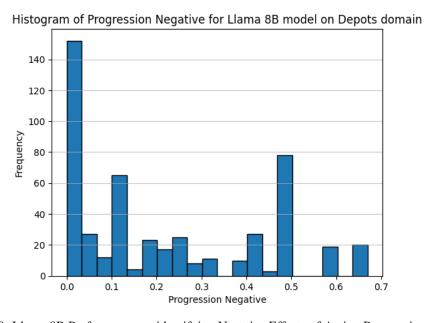


Figure 10: Llama 8B Performance on identifying Negative Effects of Action Progression in Depots Domain

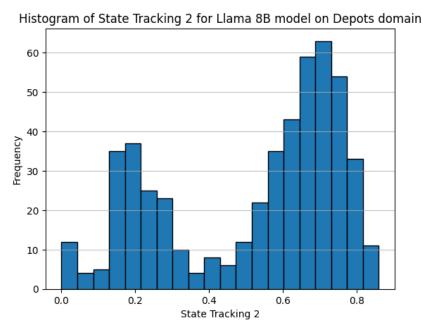


Figure 11: Llama 8B Performance on State tracking with 2 Actions in Depots Domain

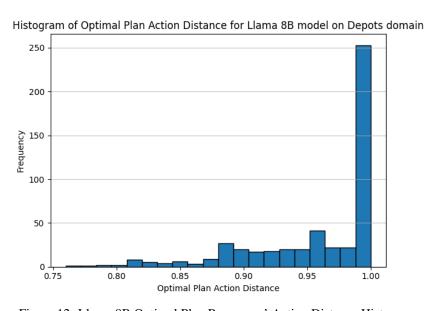


Figure 12: Llama 8B Optimal Plan Responses' Action Distance Histogram