

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

Anil B Murthy¹, Jaron Mink^{1*}, Lindsay Sanneman^{1*}

¹School of Computing & Augmented Intelligence, Arizona State University
abmurthy@asu.edu, jaron.mink@asu.edu, lindsay.sanneman@asu.edu

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 of real-world agentic systems. 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 multiple case study examples, we additionally demonstrate the feasibility of comparative analysis of different task-specific data distributions obtained through this metric.

Introduction

The ability to plan, perform sequential decision-making, and reason about action and change is 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 applications, from question answering and document summarization to code generation (Hagos, Battle, and Rawat 2024). Despite their impressive capabilities, LLMs have shown significant limitations in planning, reasoning, and decision-making, particularly in autonomous applications (Kokel et al. 2025b; Valmeekam et al. 2023b; Kambhampati et al. 2024; Handa et al. 2025). Such limitations in LLMs’ performance are noted through task evaluations that utilize binary success criteria metrics that involve comparison with ground truth answers obtained by automated solvers, planners, or validators. However, there exists useful information about the quality and precision of the models’ responses for

these task evaluations, which is not necessarily captured by standard binary metrics, that can help with comprehensive and domain/instance-specific diagnostic analyses, for developing real-world deployable agentic systems.

As LLMs are probabilistic models and generate tokens in an autoregressive manner, it is perhaps not surprising that they struggle to perform accurately on Reasoning about Action, Change (RAC), 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 standard binary success metrics. Specifically, our proposed metrics elicit more information about 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, and to make decisions about how best to utilize the model in architectural frameworks that are based on LLM-Modulo (Kambhampati et al. 2024), ReAct (Yao et al. 2022), and other finetuning or prompting setups to enhance performance.

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

Background & Related Works

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 (Valmeekam et al. 2023a; Handa et al. 2025; He et al. 2022; Kokel et al. 2025b). 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 (He et al. 2022). They pre-train and evaluate transformer models such as GPT-2 (Radford et al. 2019) 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 unclear if the task design maintains structural validity (measurement reflecting the internal structure of the construct) (Salaudeen et al. 2025).

Valmeekam et al. developed PlanBench, a PDDL-based planning benchmark suite with 8 planning-related tasks, such as plan generation, cost-optimal planning, plan verification, goal recognition, replanning, plan reuse, reasoning about actions and effects, and plan generalization (Valmeekam et al. 2023a). The PlanBench work evaluates LLMs like GPT-4 (Achiam et al. 2023) and Instruct-GPT-3 (Ouyang et al. 2022) 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 (Russell, Norvig, and Intelligence 1995; Ghallab, Nau, and Traverso 2025).

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 (competition 2024) like Blocksworld (Handa et al. 2025). 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 (Wang et al. 2023).

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 (Shridhar et al. 2021) and a novel 'swap' planning domain (Kokel et al. 2025b). 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 (Kokel et al. 2025a), which may lead to low or unclear construct validity (Salaudeen et al. 2025).

Domains

To demonstrate the utility of our proposed benchmarks, we utilize standard IPC planning domains (competition 2024) such as Blocksworld and Depots for our experiments to evaluate the planning and action reasoning abilities of LLMs. For these two domains, we create 500 problem instances in PDDL, for each of which we further create natural language templates for the initial and goal states, and questions for 6 different tasks, resulting in approximately 6000 questions

that we use to evaluate the Llama 8B and 70B models. For each problem, all the 6 task questions have the same object complexity, initial state, and goal state, only differing in the question prompt. A common natural language context containing the domain description, initial state description and goal state description (if necessary) is utilized for evaluating the LLMs, to ensure as 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 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.

Tasks: Reasoning about Action, Change, and Planning

Drawing from the above benchmarks in Section , 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 (Reuel-Lamparth et al. 2024; Salaudeen et al. 2025).

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 kind of 'set comparison'-based metric allows us to obtain more fine-grained information about the quality of LLMs' performance for each task.

$$\text{Task Metric} = \frac{\text{LLM Answers} \cap \text{Ground Truth}}{\text{LLM Answers} \cup \text{Ground Truth}} \quad (1)$$

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

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

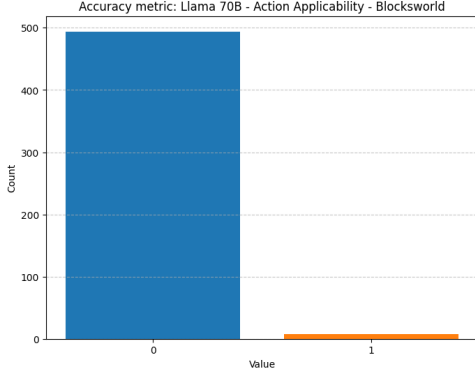


Figure 1: Llama 70B Performance with Standard binary success metric on Action Applicability Task in Blocksworld; Accuracy = 0.014%; The model’s responses are correct on only 7/501 problems.

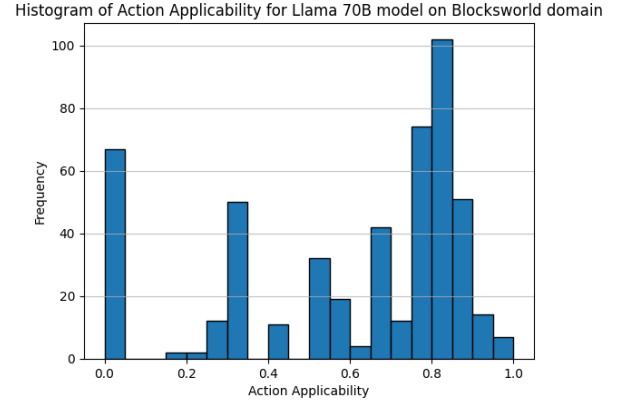


Figure 2: Llama 70B Performance with IoU metric on Action Applicability Task in Blocksworld; This right-skewed distribution provides information on the precision of the model’s responses. We can see that the model is close to correctness ($> 75\%$) on more than 200/501 problems.

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.

given state, provided the common context, as mentioned in the Domains section above, using the IoU evaluation metric shown in equation 1 and table 1.

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

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 of domain and initial state descriptions and the specified action.

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.

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. We prompt the LLMs for generating plans given the domain, state, and goal contexts. Evaluation is performed using the well-known set comparison metric called ‘Action Distance’ (Nguyen et al. 2012), as shown in Table 1 and the evaluation is further detailed in the Appendix .

Cost-Optimal Plan Generation: For a plan generation task, if actions have costs, then an optimal plan has the minimum possible cost. We prompt the LLMs to provide optimal

plans given the domain, state, and goal context. Evaluation is performed using the Action Distance metric (Nguyen et al. 2012).

Results and Discussion

In this work, we perform evaluations with 6 tasks (considering progression effects as two tasks) across two domains of 500 problems each, on two instruction-tuned pretrained LLMs, using informative task-specific IoU metrics. In figure 2, we can see that the data distribution obtained through the IoU metric provides us with substantial information on the precision, quality, and nature of models’ responses that are entirely missed by binary success metrics, as shown in figure 1.

In figure 2, the right-skewness of the distribution demonstrates that the model is much closer to being correct than the 0 values for 494 samples imply. In fact, the model’s performance is over 75% accurate for more than 200 samples. This information is extremely beneficial for compute-intensive and cost-incurring decisions such as finetuning procedures, and for inference-time decisions such as model-routing, repeated sampling or prompting setups. Additionally, figure 2 shows that over 70 instances have a low performance of $< 0.05\%$, indicating the need for instance-specific analysis of those samples. Further, this metric helps the design of future experiments to understand and improve specific atomic reasoning constructs or capabilities of generative AI models, such as reasoning about action applicability and state understanding.

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

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

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

Histogram of Action Applicability for Llama 8B model on Depots domain

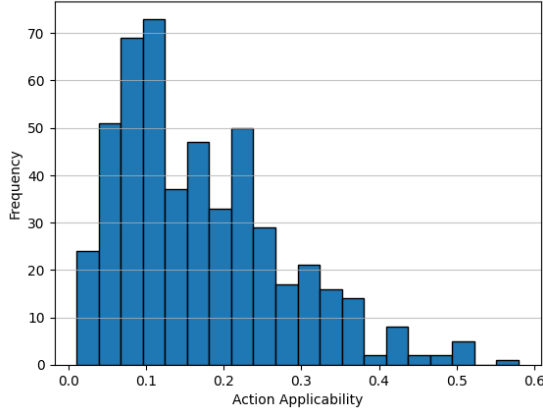


Figure 4: Llama 8B Performance with IoU metric on Action Applicability in Depots domain;

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

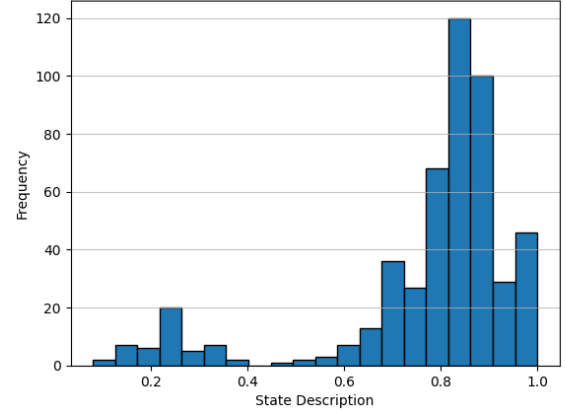


Figure 5: Llama 8B Performance with IoU metric on State Comprehension/ Description Task in Depots domain

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.

clear that the quality and precision of the model’s responses for state comprehension are much better than its ability for reasoning about applicable actions. Also, the spread of the distribution for the Action applicability task, according to figure 7, indicates that the model’s responses are less precise and more fuzzy compared to those of State comprehension in the Depots domain. Thus, the IoU metric can potentially provide discriminant validity (Salaudeen et al. 2025), where the evaluation helps differentiate between constructs that should be distinct. Essentially, this distributional comparison indicates that the model is better at understanding a given initial state than at reasoning about what actions can be applied in that state in the Depots domain.

Also, these distributions can be 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 with reference to the domain-specific state properties. Com-

paring figures 4 and 11, the model seems to be more precise at tracking changes across states than at reasoning about applicable actions in the current state. However, further case-based analysis is required to examine the action sequence and the corresponding affected objects in high-state-tracking performance samples, to investigate whether any particular domain dynamics lead to higher state-tracking performance. Using the state tracking IoU metric, we have found preliminary evidence of specific domain dynamics acutely affecting the variance in state tracking performance in both domains, particularly with odd and even numbered action sequence lengths.

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, all of which are highly valuable for development decisions on finetuning and model utility in architectural frameworks. We have demonstrated the utility of the met-

ric through evaluations and comparative examples across two domains. A more in-depth correlational analysis across tasks, domain-specific and task-specific investigations that are beyond the scope of this project is left for future work.

References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- competition, I. 2024. ICAPS International Planning Competition (IPC). <http://www.icaps-conference.org/competitions/>. Accessed: 09/2025.
- Ghallab, M.; Nau, D.; and Traverso, P. 2025. *Acting, Planning, and Learning*. Cambridge University Press.
- Hagos, D.; Battle, R.; and Rawat, D. B. 2024. Recent Advances in Generative AI and Large Language Models: Current Status, Challenges, and Perspectives. *IEEE Transactions on Artificial Intelligence*, 5: 5873–5893.
- Handa, D.; Dolin, P.; Kumbhar, S.; Son, T. C.; and Baral, C. 2025. ActionReasoningBench: Reasoning about Actions with and without Ramification Constraints. In *The Thirteenth International Conference on Learning Representations*.
- He, W.; Huang, C.; Xiao, Z.; and Liu, Y. 2022. TRAC: A Textual Benchmark for Reasoning about Actions and Change. *arXiv preprint arXiv:2211.13930*.
- Kambhampati, S.; Valmeekam, K.; Guan, L.; Verma, M.; Stechly, K.; Bhambri, S.; Saldyt, L.; and Murthy, A. 2024. 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.
- Kokel, H.; Katz, M.; Srinivas, K.; and Sohrabi, S. 2025a. ACPBench Hard: Unrestrained Reasoning about Action, Change, and Planning. In *AAAI 2025 Workshop LM4Plan*.
- Kokel, H.; Katz, M.; Srinivas, K.; and Sohrabi, S. 2025b. Acpbench: Reasoning about action, change, and planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 26559–26568.
- Kulkarni, A.; Zha, Y.; Chakraborti, T.; Vadlamudi, S. G.; Zhang, Y.; and Kambhampati, S. 2016. Explicablility as minimizing distance from expected behavior. *arXiv preprint arXiv:1611.05497*.
- Nguyen, T. A.; Do, M.; Gerevini, A. E.; Serina, I.; Srivastava, B.; and Kambhampati, S. 2012. Generating diverse plans to handle unknown and partially known user preferences. *Artificial Intelligence*, 190: 1–31.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; Sutskever, I.; et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8): 9.
- Reuel-Lamparth, A.; Hardy, A.; Smith, C.; Lamparth, M.; Hardy, M.; and Kochenderfer, M. J. 2024. Betterbench: Assessing AI benchmarks, uncovering issues, and establishing best practices. *Advances in Neural Information Processing Systems*, 37: 21763–21813.
- Russell, S.; Norvig, P.; and Intelligence, A. 1995. A modern approach. *Artificial Intelligence*. Prentice-Hall, Egnlewood Cliffs, 25(27): 79–80.
- Salaudeen, O.; Reuel, A.; Ahmed, A.; Bedi, S.; Robertson, Z.; Sundar, S.; Domingue, B.; Wang, A.; and Koyejo, S. 2025. Measurement to Meaning: A Validity-Centered Framework for AI Evaluation. *arXiv preprint arXiv:2505.10573*.
- Shridhar, M.; Yuan, X.; Côté, M.-A.; Bisk, Y.; Trischler, A.; and Hausknecht, M. 2021. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Valmeekam, K.; Marquez, M.; Olmo, A.; Sreedharan, S.; and Kambhampati, S. 2023a. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. *Advances in Neural Information Processing Systems*, 36: 38975–38987.
- Valmeekam, K.; Sreedharan, S.; Marquez, M.; Olmo, A.; and Kambhampati, S. 2023b. On the planning abilities of large language models (a critical investigation with a proposed benchmark). *arXiv preprint arXiv:2302.06706*.
- Wang, P.; Li, L.; Chen, L.; Zhu, D.; Lin, B.; Cao, Y.; Liu, Q.; Liu, T.; and Sui, Z. 2023. Large Language Models are not Fair Evaluators. *ArXiv*, abs/2305.17926.
- Xie, J.; Zhang, K.; Chen, J.; Zhu, T.; Lou, R.; Tian, Y.; Xiao, Y.; and Su, Y. 2024. TravelPlanner: A Benchmark for Real-World Planning with Language Agents. In *Forty-first International Conference on Machine Learning*.
- Yao, S.; Zhao, J.; Yu, D.; Du, N.; Shafran, I.; Narasimhan, K. R.; and Cao, Y. 2022. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.

Extended Task Descriptions

Action Applicability

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 (Xie et al. 2024; Kokel et al. 2025b; Handa et al. 2025). 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 , using the IoU evaluation metric shown in equation 1 and table 1.

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 understanding the properties of the current state. We ask the LLM to provide the list of predicates that fully 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, based on the initial state. The ground truth predicates representing the state are stored and used for evaluating LLMs’ responses using the IoU metric mentioned in Table 1.

Progression

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 (Handa et al. 2025; Kokel et al. 2025b; Valmeekam et al. 2023b). Also, prior works have found that LLMs’ performance differs with positive and negative predicates (Handa et al. 2025). 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). For both tasks, the predicates representing the corresponding effects are stored as ground truth and used for evaluating LLMs’ responses using the progression IoU metric mentioned in Table 1.

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.

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. In this work, we provide a sequence of 2 actions, and prompt the LLM for the predicates of the final state, with domain and initial-state descriptions as context. The evaluation is performed in the same manner as State Comprehension, using the IoU metric in Table 1.

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. Given the domain description, initial state, and goal state, this task asks the LLM to provide a sequence of actions that constitute a plan to reach the goal state from the initial state. As there may be multiple possible satisfying plans from the initial state to reach the goal state, we store only the optimal plan as the ground truth reference for evaluation with the action distance metric.

Evaluation with the Action Distance Metric 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 (Nguyen et al. 2012; Kulkarni et al. 2016). 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 (Kambhampati et al. 2024), 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. Also, unlike the IoU metrics for other tasks, the action distance metric has an additive inverse with respect to 1. This means that an action distance of 1 represents that the model’s plan has an entirely different set of actions compared to the ground truth reference plan. And an action distance of 0 represents that the model’s plan has the same set of actions as the ground truth reference plan. However, as the action distance metric does not account for ordering of actions, a plan with action distance 0 may still be invalid and incorrect. This can be construed as ”the plan has all the right actions, but not in the right order”. From this perspective, the action distance metric can be useful to identify how far off generative AI models are at generating the correct set of actions.

For the plan generation task, although there may be numerous satisficing plans for a given pair of initial state and goal state, we evaluate the action distance metric with respect to an optimal plan as the reference. This provides us with information on the model’s ability to choose landmark actions (actions that are part of all plans for a given initial state and goal state).

Cost-Optimal Plan Generation

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 is performed similarly to plan generation using the action distance metric with optimal plan as the reference, which is also the ground truth for this task.

Evaluation Setup

In order to enable performance comparison between tasks and characterize model-specific and model-agnostic strengths, the evaluation is set up to be as holistic as possible. Thus, to prompt the models for each task, a common prompt comprising the domain description, action model description, description of the initial state conditions, and (if necessary) goal state conditions is provided to the models. Finally, the question that is specific to each task, such as "Generate a list of actions that are applicable in the given state," is provided to the model. Depending on whether the task involves output of actions or predicates, the system prompt is also task-specific, instructing the model to "provide actions within [] and separated by commas." As the models that we are evaluating are instruction-tuned, we believe the changes in the system prompt instructing the model to provide answers in certain formats, to enable robust and easier evaluation, should not affect the model performance significantly. The full prompts are provided in the Prompts section of the Appendix .

Tasks Performance Graphs for IoU metric on Depots Domain

Histogram of Action Applicability for Llama 8B model on Depots domain

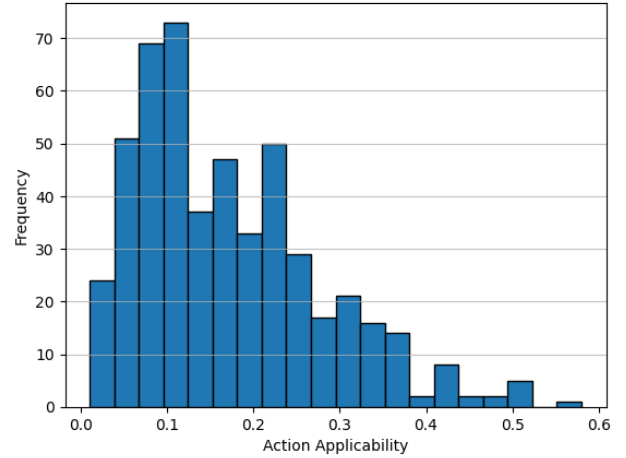


Figure 7: Llama 8B Performance on Action Applicability in Depots Domain

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

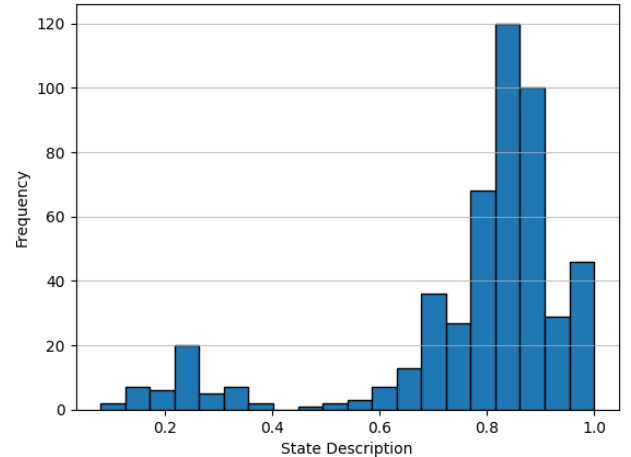


Figure 8: Llama 8B Performance on State Comprehension in Depots Domain

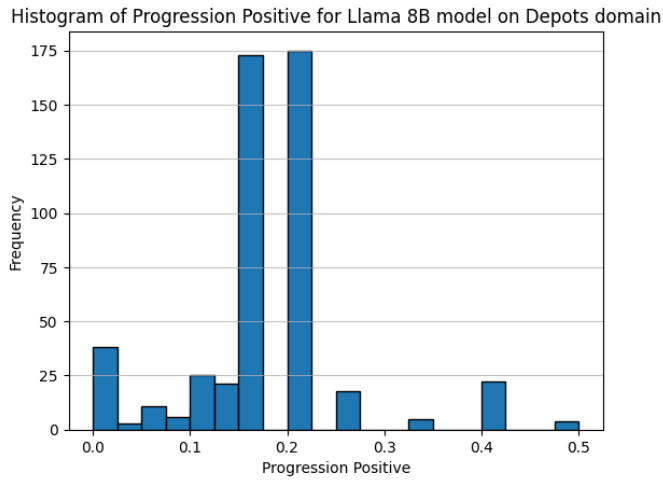


Figure 9: Llama 8B Performance on Identifying Positive Effects of Action progression in Depots Domain

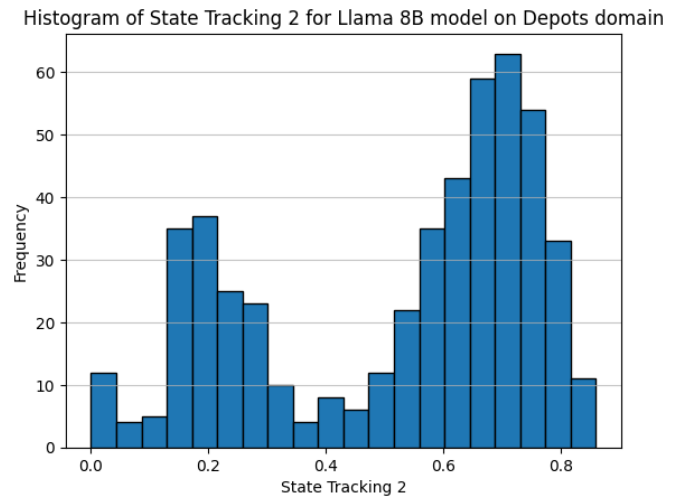


Figure 11: Llama 8B Performance on State tracking with 2 Actions in Depots Domain

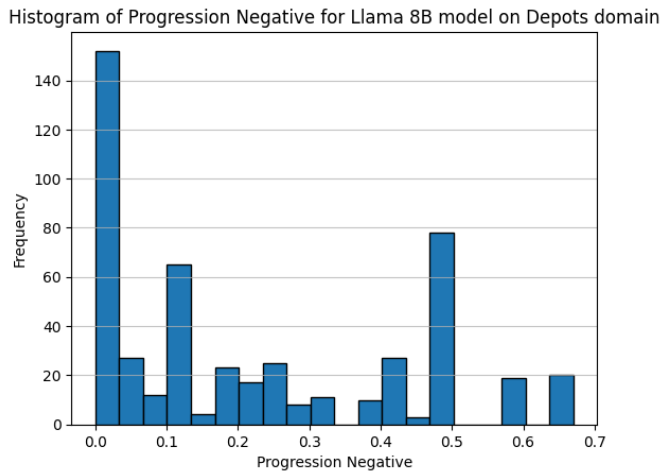


Figure 10: Llama 8B Performance on Identifying Negative Effects of Action Progression in Depots Domain

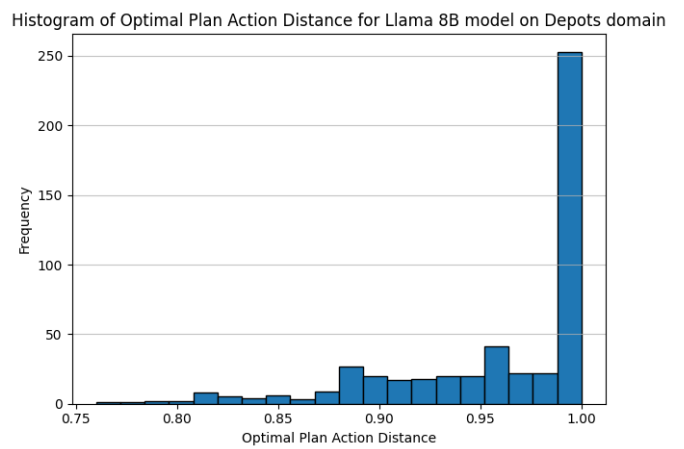


Figure 12: Llama 8B Optimal Plan Responses' Action Distance Histogram

Tasks Performance Graphs for IoU metric on Blocksworld Domain

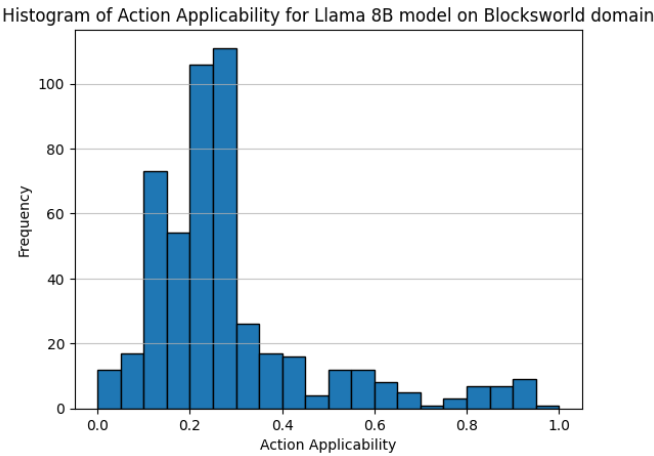


Figure 13: Llama 8B Action Applicability Histogram on Blocksworld Domain

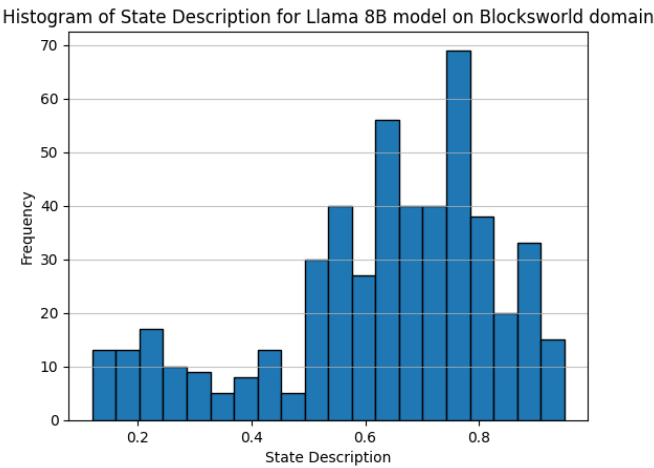


Figure 14: Llama 8B State Comprehension Histogram on Blocksworld domain

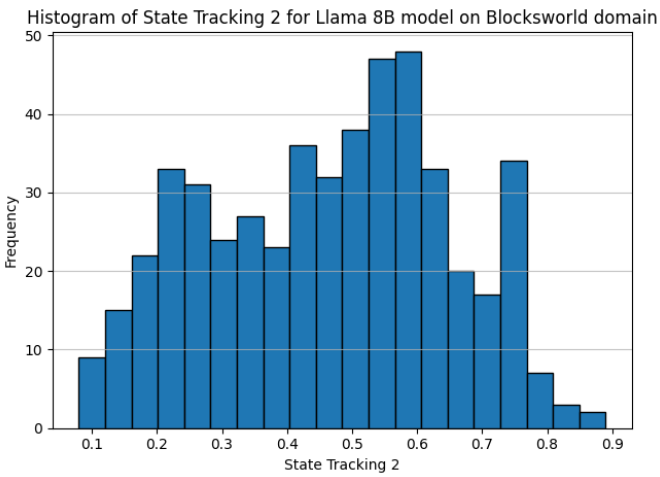


Figure 15: Llama 8B Performance Histogram for State tracking with 2 actions

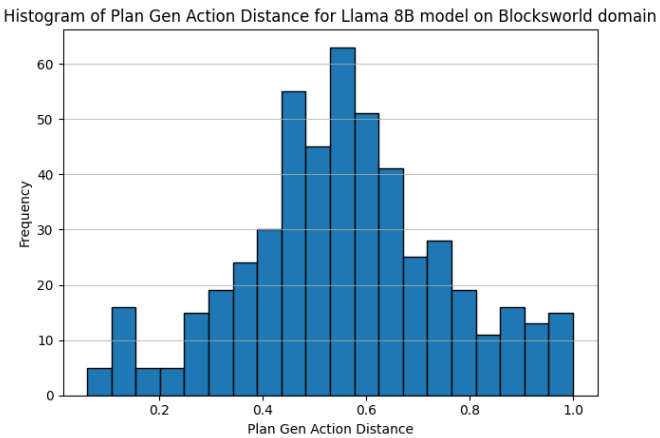


Figure 16: Llama 8B Plan Generation Action Distance Histogram

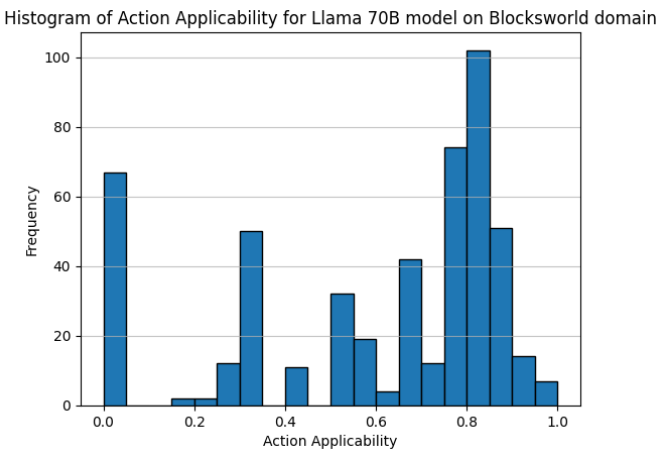


Figure 17: Llama 70B Action Applicability Histogram

Histogram of Optimal Plan Action Distance for Llama 70B model on Blocksworld domain

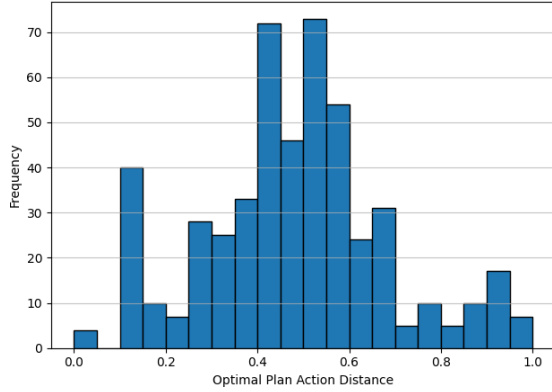


Figure 18: Llama 70B Optimal Plan Action Distance Histogram

Prompts

Prompt: Plan Generation

”This is a blocksworld domain where blocks can be placed on top of each other or on the table. There is one robotic arm hand that can move the block. Here are the actions that can be performed:

(pick-up block) to Pick up a block from the table,
(unstack block another_block) to Unstack a block from on top of another block,
(put-down block) to put down a block on the table,
(stack block another_block) to stack a block on top of another block.

There are the following restrictions on actions:

Only one block at a time can be picked up or unstacked. A block can only be picked up or unstacked if the hand is empty. A block can only be picked up if it is on the table and is clear. A block is clear if no other blocks are on top of it and if it is not picked up. A block can only be unstacked from on top of another block if it is truly on top of the other block. A block can only be unstacked from on top of another block if it is clear. Once a block is picked up or unstacked, it is being held and is no longer clear. Once a block is being held, that block can either be put-down or stacked on top of another block. A block can only be stacked on top of another block by me if the block onto which it is being stacked is clear. Once a block is put down or stacked, the hand becomes empty. Once a block is stacked on top of a second block, the second block is no longer clear.

There are 3 blocks. Currently, the robotic arm is empty. The following blocks are on the table: i, f. The following blocks are stacked on top of another block: block g is on block i. The goal is to reach a state where the following facts hold: The following blocks are on the table: f. The following blocks are stacked on top of another block: block g is on block i and block i is on block f.

Provide a plan as a list of actions that can be executed consecutively from the current state to reach the goal state. The available actions are: (pick-up ?ob) - pick up block ?ob; (put-down ?ob) - put down block ?ob; (stack ?ob ?underob) - stack ?ob on top of ?underob; (unstack ?ob ?underob) - unstack ?ob from on top of ?underob;”

Prompt: Action Applicability

”This is a blocksworld domain where blocks can be placed on top of each other or on the table. There is one robotic arm hand that can move the block. Here are the actions that can be performed:

(pick-up block) to Pick up a block from the table,
(unstack block another_block) to Unstack a block from on top of another block,
(put-down block) to put down a block on the table,
(stack block another_block) to stack a block on top of another block.

There are the following restrictions on actions:

Only one block at a time can be picked up or unstacked. A block can only be picked up or unstacked if the hand is empty. A block can only be picked up if it is on the table and is clear. A block is clear if no other blocks are on top of it and if it is not picked up. A block can only be unstacked from on top of another block if it is truly on top of the other block. A block can only be unstacked from on top of another block if it is clear. Once a block is picked up or unstacked, it is being held and is no longer clear. Once a block is being held, that block can either be put-down or stacked on top of another block. A block can only be stacked on top of another block by me if the block onto which it is being stacked is clear. Once a block is put down or stacked, the hand becomes empty. Once a block is stacked on top of a second block, the second block is no longer clear.

There are 3 blocks. Currently, the robotic arm is empty. The following blocks are on the table: i, f. The following blocks are stacked on top of another block: block g is on block i.

Generate a list of ground actions that are applicable in this state. The available actions are: (pick-up ?ob) - pick up block ?ob; (put-down ?ob) - put down block ?ob; (stack ?ob ?underob) - stack ?ob on top of ?underob; (unstack ?ob ?underob) - unstack ?ob from on top of ?underob;”

Histogram of Action Applicability for Llama 8B model on Blocksworld domain

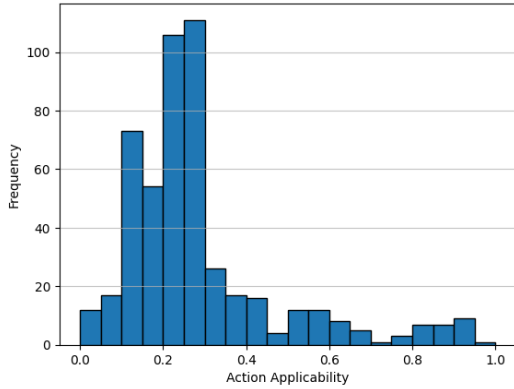


Figure 19: Llama 8B Performance with IoU metric on Action Applicability in Blocksworld domain;

Histogram of Action Applicability for Llama 70B model on Blocksworld domain

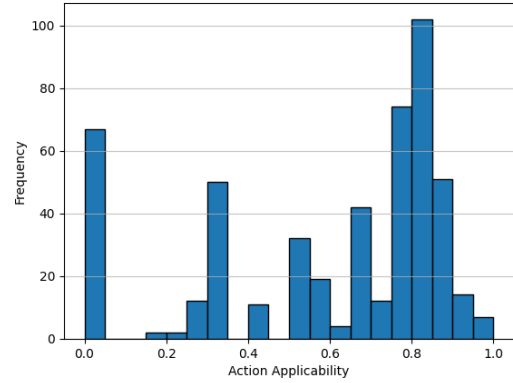


Figure 20: Llama 70B Performance with IoU metric on Action Applicability in Blocksworld domain

Figure 21: Comparison of IoU Metric evaluation of Llama 8B and 70B models on the Action Applicability Task. It is evident from the left-skewed distribution of figure 19 and the right-skewed distribution of figure 20 that Llama 70B model's responses and performance is more precise and of higher quality than those of the Llama 8B Model.