

Language Model as Planner and Formalizer under Constraints

Anonymous ACL submission

Abstract

LLMs have been widely used in planning, either as planners to generate action sequences end-to-end, or as formalizers to represent the planning domain and problem in a formal language that can derive plans deterministically. However, both lines of work rely on standard benchmarks that only include generic and simplistic environmental specifications, leading to potential overestimation of the planning ability of LLMs and safety concerns in downstream tasks. We bridge this gap by augmenting widely used planning benchmarks with manually annotated, fine-grained, and rich natural language constraints spanning four formally defined categories. Over 4 state-of-the-art reasoning LLMs, 4 formal languages, and 4 datasets, we show that the introduction of one-sentence constraints consistently halves performance, indicating current LLMs' lack of robustness and avenue for future research.¹

1 Introduction

Large language models (LLMs) have garnered attention in recent years for their capabilities in planning domains. An intuitive methodology, *LLM-as-planner*, directly generates a plan given a description of the domain and the problem in an end-to-end manner. While the algorithmic nature of formal planning tasks significantly challenge LLM planners (Valmeekam et al., 2024a; Kambhampati et al., 2024), the emergence of reasoning LLMs has shown improved performance (Valmeekam et al., 2024b; Huang and Zhang, 2025) but lacks interpretability and control. Another methodology, *LLM-as-formalizer*, instead translates the input description into some formal languages (Xie et al., 2023; Liu et al., 2023a; Zhang et al., 2024a,b; Zhu et al., 2024). This formal representation can then be passed into a solver which then outputs a plan

deterministically, providing more interpretability, control, and promising performance.

As a recent, nascent research direction, most previous work relies on standard planning benchmarks, some having existed for decades, where environments are described in a generic, homogeneous, simplistic manner. This is a critical issue because the simplicity of the data may lead to overestimation of the planning or formalizing ability of an LLM. While existing benchmarks enable proof of concept, they are distant from real-life planning instructions that commonly include idiosyncratic requirements inflicted by the user or the resources.

To bridge this gap, we introduce CoPE (Constrained Planning Environments), a benchmark that augments widely used planning environments such as BlocksWorld (IPC, 1998) and CoinCollector (Yuan et al., 2019) with **constraints**. Each example originally includes a domain description, a problem description, and ground-truth PDDL as a simulator to validate a predicted plan. In CoPE, we additionally annotate a natural language constraint that is not only linguistically rich but also formally categorized (Figure 1). We systematically evaluate both LLM-as-planner and LLM-as-formalizer methodologies, the latter including 4 formal languages including Planning Definition Language (PDDL) 1.2 and 3, Linear Temporal Logic (LTL), and Satisfiability Modulo Theories (SMT). Among 4 state-of-the-art LLMs on 4 datasets, we show that one-sentence constraints consistently halves performance, indicating lack of robustness of both LLM planning methodologies. To inspire future work, we show that different formalisms respond differently to categories of constraints and also domains. As we corroborate past work by showing the robustness of LLM-as-formalizer over LLM-as-planner to problem complexity and lexical shift, we clearly demonstrate that the introduction of constraints substantially negates such robustness. We pose our findings as a timely reflection among a

¹Our code and data are attached with the submission.

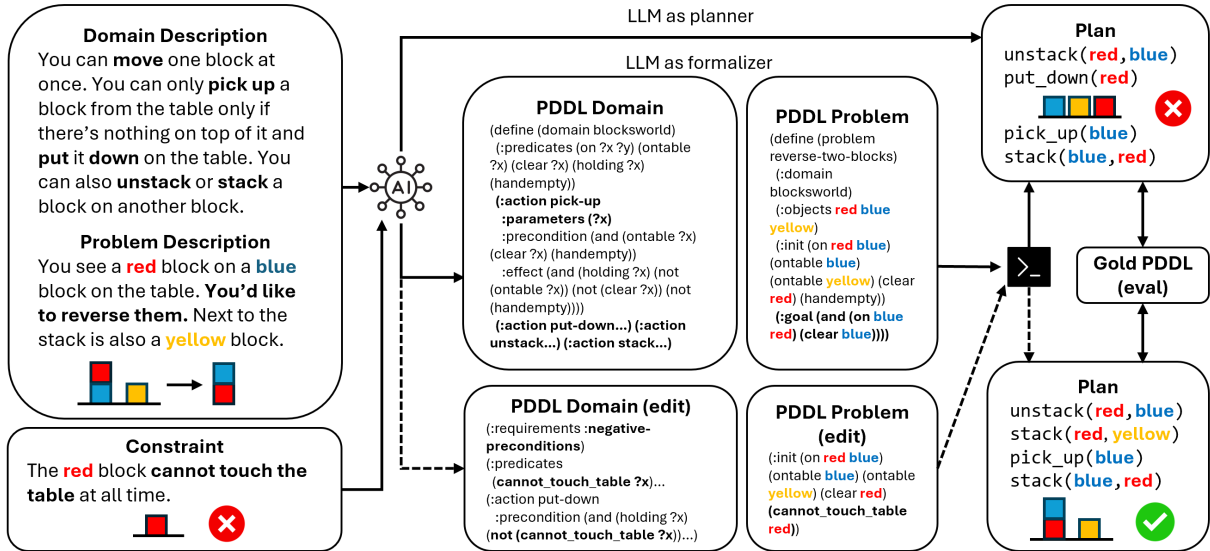


Figure 1: The flow of evaluation on CoPE. An LLM is given a description of the planning domain and problem, as well as a natural language constraint. As a planner, the LLM directly outputs a plan consisting of a sequence of actions. As a formalizer, the LLM attempts to construct a formal representation that can invoke an external solver to derive a plan. Alternatively, to address the constraint, the LLM as an editing formalizer (dashed arrows) would first predict PDDL without the constraint before generating edits to it.

recent slew of methodological work in LLM-based planning, highlighting limitations and opening up future directions.

2 Related Work

As comprehensive literature reviews of the LLM-as-planner and LLM-as-formalizer methodology have been accomplished by recent surveys (Tantakoun et al., 2025; Wei et al., 2025), we will instead focus on the benchmarks used in existing work.

Planning Benchmarks Most work on planning with LLMs experimented with International Planning Competition (IPC) domains including BlocksWorld, Gripper, Logistics, Barman, etc. (Kagitha et al., 2025; Hu et al., 2025; Zhu et al., 2024; Liu et al., 2023a; Shirai et al., 2024; Hao et al., 2025; Guo et al., 2024; Shojaee et al., 2025; Stein et al., 2025; Stechly et al., 2024; Valmeekam et al., 2024a; Yao et al., 2025; Guan et al., 2023). We provide a table of works and a non-exhaustive list the domains used, as well as examples from common domains in Appendix A. Beyond those domains and datasets that have existed for decades, few planning benchmarks have been constructed in recent years. Moreover, Zuo et al. (2024); Huang and Zhang (2025) recently pointed out the templated nature and linguistic homogeneity of these datasets and attempted to describe the domains and problems in a natural fashion. While this modifica-

tion added difficulty to the task, the semantics of the domains and problems remain standard, subject to LLMs’ memorization, and does not reflect the complexity of real-world planning.

Data Augmentation Much recent work augmented datasets for different tasks in order to make the problem more realistic, challenging, and less prone to data contamination. This has been done for question answering (Anantheswaran et al., 2024; Hong et al., 2024), common sense reasoning (Sun and Emami, 2024), instruction following (Hong et al., 2024), and so on. For planning, datasets have been augmented to increase robustness, for example by adding noise (Zheng et al., 2025) and lexical perturbation (Valmeekam et al., 2024a; Zheng et al., 2024), with no change to the actual semantics. In contrast, our proposed CoPE semantically modifies both the domain and the problems, leading to increased realism and challenge.

Planning With Constraints In the formal planning community, syntax that expresses some constraints were introduced to the PDDL3 language (Gerevini and Long, 2005; Edelkamp, 2006). However, few follow-up work in the NLP and AI community recognized and addressed the real-life challenge of planning with constraints. Yang et al. (2023) introduced a hybrid system that translates natural language constraints into LTL to assist LLM-as-planner methods, though their focus was

137 primarily on robotic applications and lacked a thor- 187
 138 ough discussion of diverse constraint patterns and 188
 139 formal representations. Guo et al. (2024) encoded 189
 140 constraints with code to interface a SMT solver 190
 141 and combined it with PDDL for task and motion 191
 142 planning though without any public data or code. 192
 143 Moreover, they only considered coarse-grained con- 193
 144 straints and oversimplified the task by assuming 194
 145 parts of the PDDL which is often not realistic. We 195
 146 closely consider these formalisms when defining 196
 147 the constraints. 197

148 3 Task Formulation 199

149 **Definition of Planning** We will first formally 200
 150 define the task of text-based formal planning in 201
 151 line with classical STRIPS formalism (Fikes and 202
 152 Nilsson, 1971) and recent conventions in the nat- 203
 153 ural language processing community (Huang and 204
 154 Zhang, 2025). The input I consists of a triplet 205
 155 $(D_d, D_p, \mathcal{DF}')$: 206

- 156 1. D_d is a textual description of the domain, spec- 207
 157 ifying the types (the type system for entities 208
 158 e_i), predicates (in relation to the types), and 209
 159 actions a_i including pre-conditions and effects 210
 160 (parameterized by conjunctive logical formu- 211
 161 las over the the predicates). 212
- 162 2. D_p is a textual description of the problem, 213
 163 specifying the objects, initial states and goal 214
 164 states in terms of predicates; 215
- 165 3. \mathcal{DF}'_G is the header of the ground-truth PDDL 216
 166 domain file \mathcal{DF}_G that only provides the names 217
 167 and parameters of the actions. This is required 218
 168 so that the eventual plan is grounded to an 219
 169 existing ontology and can be evaluated fairly. 220

170 The final objective is a plan $\pi =$ 221
 171 $[a_0(\bar{e}_0), \dots, a_m(\bar{e}_m)]$, where each $a_i \in A$ is 222
 172 an action provided by \mathcal{DF}'_G , and each \bar{e}_i is a 223
 173 tuple of grounded entities corresponding to the 224
 174 parameters of a_i . The plan π , when executed 225
 175 in the environment expressed by a ground-truth 226
 176 \mathcal{DF}_G and \mathcal{PF}_G , must (i) fulfill the pre-conditions 227
 177 of each action, and (ii) constitutes a successful 228
 178 transition from the initial state to the goal state. 229
 179 We will revisit the formalism of plan execution 230
 180 after defining the domain and problem files. 231

181 The domain file \mathcal{DF} is defined by three main 232
 182 components: types (T), predicates (R), and action 233
 183 semantics $(\alpha^{param}, \psi_a^{pre}, \psi_a^{eff})$: 234

- 184 1. Types: $T = t_1, \dots, t_n$ is the set of entity types 235
 185 present in the environment. 236
- 186 2. Predicates: $R = r_1, \dots, r_n$ is the set of re- 237

187 lational predicates, each in the form of $r(\bar{t})$, 188
 189 where \bar{t} is a tuple of types that participate in 190
 relation r . 191

- 192 3. Action semantics: given an action, the pa- 193
 194 rameters α^{param} are a set of types \bar{t} that par- 194
 195 ticipate in this action. The pre-conditions 195
 196 $\psi_a^{pre} : \mathcal{S} \rightarrow \mathbb{B}$ is a boolean formula over the 196
 197 current state which presides over a set of predi- 197
 198 cates \bar{r} , where \mathcal{S} is the state space. The action 198
 199 is executable iff $\psi_a^{pre} = \text{True}$. The effects 199
 200 $\psi_a^{eff} : \mathcal{S} \rightarrow \mathcal{S}$ is a formula that transition the 200
 201 current state to the next state after the execu- 201
 202 tion of this action. 202

203 The problem file \mathcal{PF} is defined by three main 204
 205 components: objects (E), the initial state (s_0), and 205
 206 the goal states (\mathcal{S}_g): 206

- 207 1. Objects: $E = e_1, \dots, e_n$ is the set of named 207
 208 entities present in the environment. Each ob- 208
 209 ject is an instance of an entity type t defined 209
 210 in \mathcal{DF}'_G . 210
- 211 2. Initial State: $s_0 \in \mathcal{S}$ is a state defined by set of 211
 212 relational facts of the form $r(\bar{e})$, where $r \in R$ 212
 213 is a relational predicate; \bar{e} is a tuple of entities 213
 214 from E that participate in relation r . 214
- 215 3. Goal States: $\mathcal{S}_g \subseteq \mathcal{S}$ specifies the set of goal 215
 216 states to be reached. A state s is a goal state if 216
 217 and only if $s \in \mathcal{S}_g$. 217

218 Here, the state space \mathcal{S} comprises all possi- 218
 219 ble configurations of relational facts over the ob- 219
 220 ject set E with the predicates R . Each state 220
 221 $s \in \mathcal{S}$ is a set of instantiated facts $\bar{r}(\bar{e})$, describing 221
 222 which relations hold among the objects. Success- 222
 223 fully executing a plan $\pi = [a_0(\bar{e}_0), \dots, a_m(\bar{e}_m)]$ 223
 224 from the initial state yields a trace of states $L =$ 224
 225 $[s_0, s_1, \dots, s_{m+1}]$, where $\psi_{a_i}^{pre}(s_i) = \text{True}$, $s_{i+1} =$ 225
 226 $\psi_{a_i}^{eff}(s_i)$, and $s_{m+1} \in \mathcal{S}_g$. 226

227 While LLM-as-planner maps the input triplet I 227
 228 directly to a plan L , LLM-as-formalizer generates 228
 229 the domain file \mathcal{DF}_P and \mathcal{PF}_P before inputting 229
 230 them into a PDDL solver to search for a plan. 230

231 **Definition of Constraint** While the definition 231
 232 of *constraint* varies by community, we take a lin- 232
 233 guistic and pragmatic approach and define it as 233
 234 supplementary, non-destructive information to the 234
 235 task description that restricts possible behaviors in 235
 236 a plan. We aim to power real-world systems that 236
 237 can adhere to any constraint expressed by human 237
 238 users, encompassing formal definitions from es- 238
 239 tablished fields (e.g., classical planning) without 239
 240 being limited by them. Specifically, CoPE includes 240
 241 formally categorized constraints as follows. 241

Category	Example constraint	Formal explanation	Plan
None			unstack(red,blue) put_down(red) pick_up(blue) stack(blue,red)
Initial	The red, blue and yellow blocks are on the table.	$s'_0 = \{\text{on-table}(\text{red}) \wedge \text{on-table}(\text{blue}) \wedge \text{on-table}(\text{yellow})\}$	pick_up(blue) stack(blue,red)
Goal	You'd like to have all three blocks on the table.	$S'_g \models \{\text{on-table}(\text{red}) \wedge \text{on-table}(\text{blue}) \wedge \text{on-table}(\text{yellow})\}$	unstack(red,blue) put_down(red)
Action	If you move the red block you must move the yellow block immediately after.	$\Pi - \Pi' = \{L \mid \exists a_i \in L, a_i = \text{move}(\text{red}) \wedge a_{i+1} \neq \text{move}(\text{yellow}), L \in \mathcal{L}\}$	unstack(red,blue) put_down(red) pick_up(yellow) put_down(yellow) pick_up(blue) stack(blue,red)
State	Once you start, the number of blocks in each stack should not exceed 1 at any time.	$\mathcal{L} - \mathcal{L}' = \{s \mid \forall b \in \bar{b}, \text{height}(b) > 1, s \in \mathcal{S}\}$	Plan does not exist

Table 1: Examples of constraints in each category categories, explanation of the assignment, and how they would affect the plan of an example problem in BlocksWorld (Figure 1): “You see a red block on a blue block on the table. You’d like to reverse them. Next to the stack is also a yellow block.”

We consider a short text \mathcal{C} as a constraint when a plan is valid if and only if it satisfies each of the conditions specified by \mathcal{C} . We introduce two concepts *primitive* and *modified*, where applying constraints modifies the primitive action or state space.

Definition 1. An action is a **primitive action** a^* of an action a , if $\psi_a^{pre} \models \psi_{a^*}^{pre}, \psi_a^{eff} \models \psi_{a^*}^{eff}$.

Definition 2. A state is a **primitive state** s^* of a state s , if $s \models s^*$.

Assuming constraints only introduce new predicates but do not remove any existing ones, every modified action will find exactly one primitive action, and every modified state will find exactly one primitive state from the original action and state space, respectively. We then define the constraint types as follows:

Definition 3. Let $L'^* = [s_0'^*, \dots, s_T'^*]$ denote a trace of primitive states after introducing \mathcal{C} , and s_0^* be initial states without \mathcal{C} , \mathcal{C} is an **initial constraint** iff $s_0^* = s_0'$ (no new predicate is introduced to the initial state) and $s_0'^* \neq s_0^*$ (the initial state changes after introducing \mathcal{C}).

Definition 4. Let S_g and S'_g be goal states before and after introducing \mathcal{C} , \mathcal{C} is a **goal constraint** iff $\{s \mid s \in S_g, s \in \mathcal{S}\} \neq \{s'^* \mid s' \in S'_g, s' \in \mathcal{S}'\}$.

Definition 5. Let $\pi'^* = [a_0'^*, \dots, a_T'^*]$ denote a sequence of primitive actions after introducing \mathcal{C} , and Π and Π' be all valid plans before and after

introducing \mathcal{C} , \mathcal{C} is an **action constraint** iff $\Pi'^* \subsetneq \Pi^*$, where Π^* is a collection of all $\pi'^*, \pi' \in \Pi'$.

Definition 6. Let $L' = [s_0^*, \dots, s_T^*]$ denote a trace of primitive states after introducing \mathcal{C} , and \mathcal{L} and \mathcal{L}' be all valid state traces before and after introducing \mathcal{C} , \mathcal{C} is an **state constraint** iff $\mathcal{L}'^* \subsetneq \mathcal{L}$, where \mathcal{L}'^* is a collection of all $L'^*, L' \in \mathcal{L}'$, and \mathcal{C} is not any other constraint defined above.

Our categorization is complete as the definition of state constraints subsumes all possible constraints. By design, each category of constraints is expressed differently for each formalism (PDDL, PDDL3, LTL, SMT), requiring different level of modification to the original code. Notably, the union of the four categories is not trivialized by any formalism. For example, simply using the the *constraint* syntax in PDDL3 cannot express action constraints or state constraints requiring invention of new predicates. While conjunction of multiple predicates is possible, we do not discuss it in this work. Table 1 shows examples of constraints and justification of their categories.

4 Data

To evaluate LLM planning in environments with constraints, we propose the CoPE benchmark that includes two widely used planning domains:

BlocksWorld (IPC, 1998) is a domain to rearrange stacks of blocks on a table using a robotic arm. It

involves four actions: *pick up*, *put down*, *stack*, and *unstack*. We build on the BlocksWorld-100 benchmark from Huang and Zhang (2025) with 100 problems, each including a domain and problem description as well as the headers of actions (D_d, D_p, \mathcal{DF}') as input to models. It also comes with ground-truth PDDL \mathcal{DF} and \mathcal{PF} , which are used to validate a predicted plan.

CoinCollector (Yuan et al., 2019) is a domain to navigate a house full of rooms in order to find a coin. It involves two actions: *move* and *pick up*. Originally designed to evaluate interactive agents, this is a partially observable environment where no complete information is given and so making a complete plan is impossible. To place it under the framework of CoPE, we convert it into a fully observable environment by Using a depth-first search, we find all room arrangements and create problem descriptions and ground-truth problem files. We sample 20 problems with 3, 5, 7, 9, and 11 rooms using the benchmark from Jansen and Côté (2022). In total, we also have 100 problem instances consisting of input and ground-truth PDDL.

While these two domains are likely included in modern LLMs’ training data, they remain standard in literature and meaningfully challenge state-of-the-art LLMs. CoPE itself is a step forward to address memorization by changing the semantics of the domains and problems. For each domain, we manually annotate 100 constraints spanning the four categories with ground-truth code of each formalism. Unlike the original datasets, we allow problems to only occasionally have no solutions. In the full set of CoPE, each constraint is paired with each problem, resulting in 10,000 constrained planning problems. To evaluate the most representative cases while keeping analysis tractable, we our following evaluation on a subset where we manually pair each constraint and one representative problem, leading to 100 problems per domain. Examples are shown in Appendix B.

5 Experimental Setup

Given the input with the natural language constraint ($D_d, D_p, \mathcal{DF}', \mathcal{C}$), we prompt the LLM to generate 4 types of outputs: a plan directly (**LLM-as-planner**), or formal programs including **PDDL** (1.2), **PDDL3**, **SMT** (more precisely, code to interact with the Python Z3 library), and **LTL**. The generation of all output types is done via one-shot prompting following cited work in Section 2. When

generating formal programs, we consider three techniques: **Generation**, where the model directly generates the code in that formal language, **editing**, where the model first generates the code without considering the constraint, suggests an edit to that code to address the constraint, and re-generates the code accordingly, and **revision**, where the model is given up to 3 tries to revise the code via re-generation based on any error message returned by executing the code, if any.

Following the majority of cited work in Section 2, the predicted plan resulting from any method is validated against the ground-truth PDDL provided by CoPE, instead of being compared against “ground-truth” plans (Lyu et al., 2023; Liu et al., 2023b; Pan et al., 2023) since there could be multiple correct plans. We evaluate the predicted plans using the *correctness* metric which indicates the percentage that are successfully validated. For PDDL output, we use the `dual-bfws-ffparser` planner implemented by Muise (2016) as the solver and VAL (Howey et al., 2004) as the validator. For PDDL3, the output is first compiled using TCORE (Bonassi et al., 2021) into standard PDDL, then solved and validated the same way as PDDL. For SMT output, the LLM is responsible for generating code that includes a call to invoke the Z3 solver. For LTL, the LLM generates formulas representing environment dynamics, initial and goal states, and constraints. A satisfying plan is then produced by searching for a Lasso-shaped accepting run over the conjunction of all formulas. See more details in Appendix D. We use Deepseek-R1, Deepseek-V3 (Guo et al., 2025), Qwen3-32B (Team, 2025), and Qwen2.5-Coder-32B-Instruct (Hui et al., 2024) as LLMs that span different sizes of the ability to scale reasoning tokens during inference. We consciously only consider open-source models in light of open science.² We query Deepseek models with their API and Qwen models using KANI (Zhu et al., 2023) with default temperature on 1 H100 GPU. Prompts can be found in Appendix C.

6 Results

With evaluation on CoPE, we answer an array of research questions, draw findings, and provide actionable recommendations for future work.

How challenging are the constraints? Figure 2 and 6 display the performance of all methods, lan-

²See requirement H14 in <https://aclrollingreview.org/reviewerguidelines>.

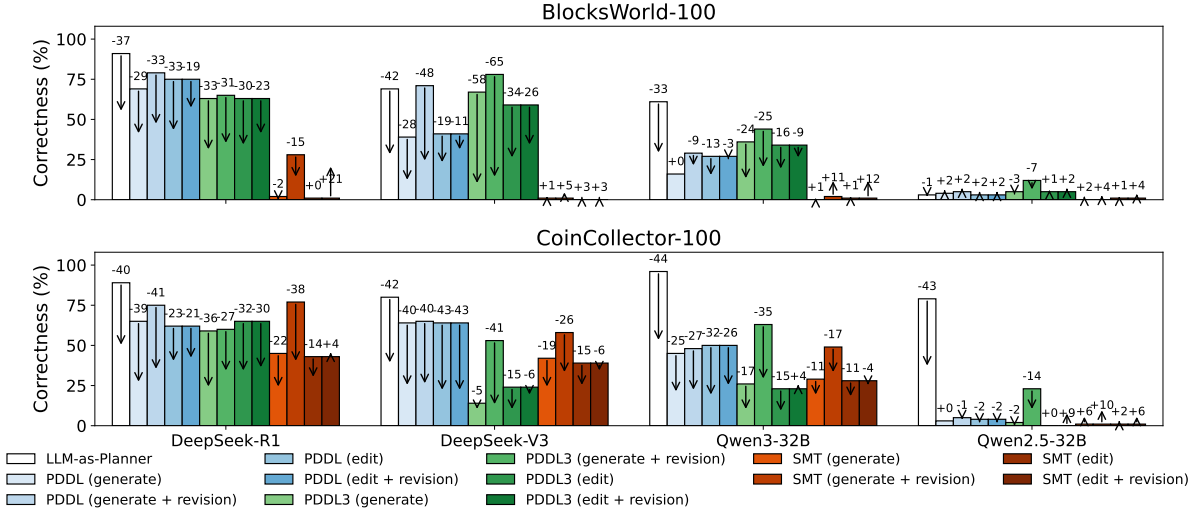


Figure 2: Performance of combinations of LLMs and methods on two domains. The arrows represent the performance difference once constraints are added. Revision is not included in this result to demonstrate that without extra help, this task is difficult. For all results, see Appendix E.

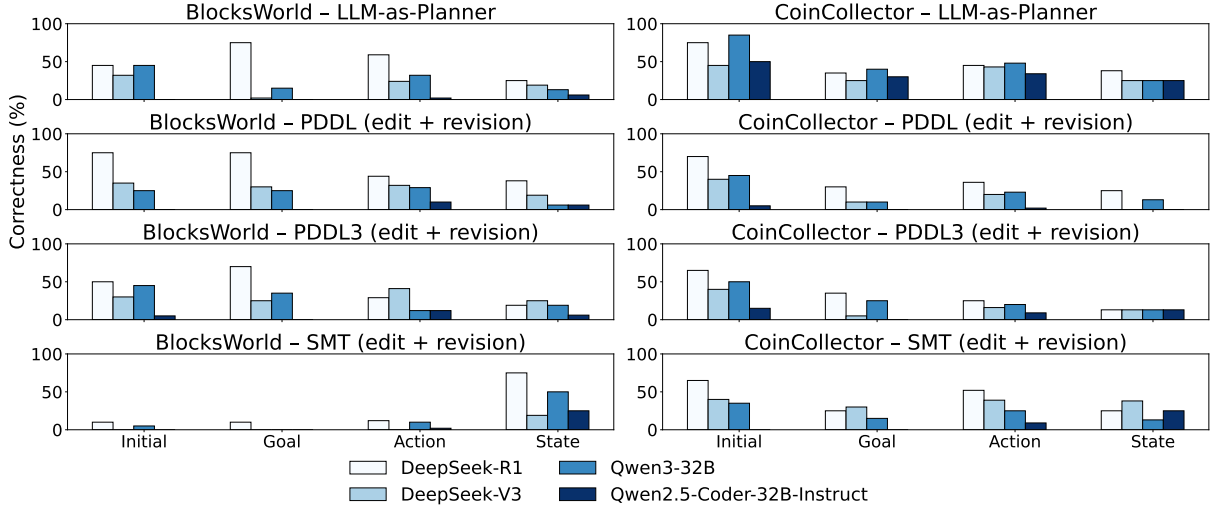


Figure 3: Best Performance of LLM-as-planner and LLM-as-formalizer (Generate) and (Edit) by constraint categories on BlocksWorld-100 and CoinCollector-100. Revision improved LLM-as-formalizer for all methods and SMT performed better than PDDL on CoinCollector-100.

393 guages, techniques, LLMs, and domains. It is
 394 evident that simple one-line **constraints greatly**
 395 **diminish performance in both planning and**
 396 **formalizing.** While all experiments maintain
 397 non-trivial correctness greater than 20%, intro-
 398 ducing constraints consistently degrades (often
 399 halves) the planning performances over all cases.
 400 Overall, LLM-as-planner outperforms its LLM-
 401 as-formalizer counterparts before and after the
 402 introduction of constraints. However, LLM-as-
 403 formalizer, with interpretability and formal guar-
 404 antee, shows competitive performance with larger
 405 DeepSeek models and using revision by solver er-

406 rors. With constraints, it benefits from the two-step
 407 generate-then-edit technique over simply generat-
 408 ing PDDL in some but not all cases. Appendix Fig-
 409 ure 7 shows that once constraints are introduced,
 410 the number of syntax errors increases for most
 411 methods. Comparing the two domains, Coin Collec-
 412 tor (2 actions) is much simpler than BlocksWorld (4
 413 actions) with regard to LLM-as-planner and LLM-
 414 as-SMT-formalizer but comparable with regard to
 415 other methods. Intuitively, this evidences LLM-as-
 416 PDDL-formalizer’s robustness against the complex-
 417 ity by the number of actions over other methods.
 418 Regarding the choice of formal language, PDDL

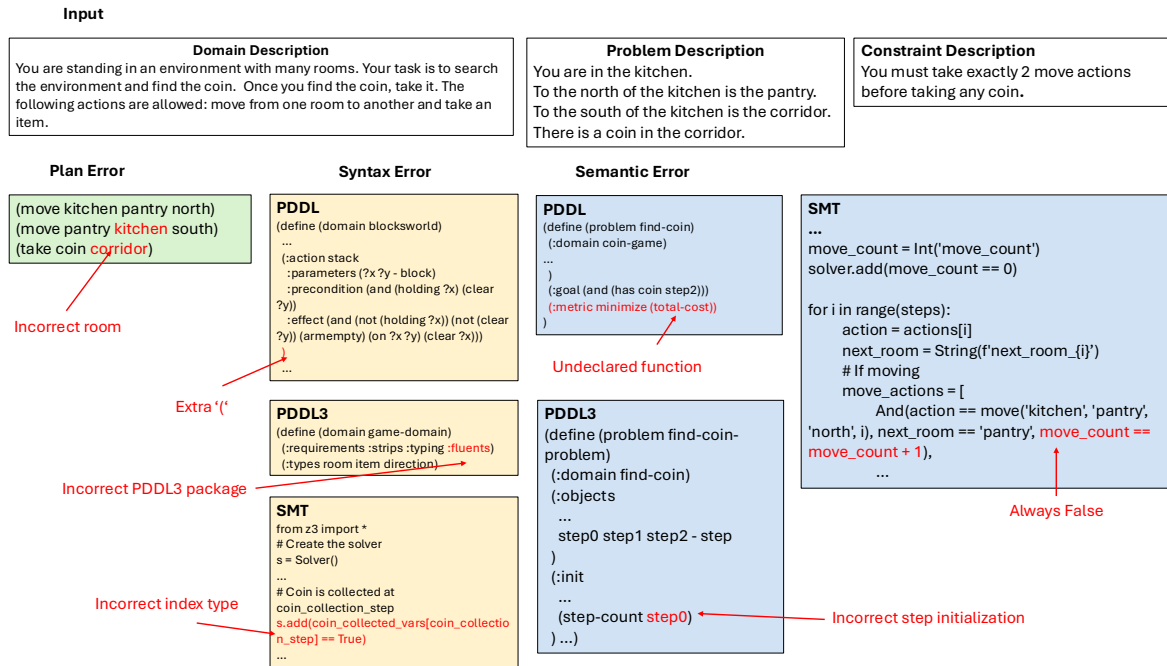


Figure 4: Examples of different errors found in model outputs.

greatly outperforms SMT in most cases, likely because PDDL more closely expresses the planning problem. However, while seemingly a more appropriate language which models constraints, PDDL3 generation is consistently worse than PDDL due to an elevated amount of compilation and syntax errors, likely because PDDL3 is even lower-resource for LLM pre-training.

When do models fail? Figure 4 qualitatively demonstrates an example set of errors for a particular constraint. A single constraint may cause various types of errors. For LLM-as-planner this set of errors includes hallucinating plans that do not exist, not satisfying newly added constraints and not following preconditions or effects set by the actions in the domain. The set of errors for LLM-as-formalizer includes syntax errors of PDDL, PDDL3 or Python, declaration of logically-incorrect functions, predicates, action parameters, preconditions in PDDL and PDDL3, or falsely encoding of constraints in SMT.

What category of constraints is hard? One may hypothesize that some categories of constraints may pose more challenge than others for a particular method. Recalling the definitions in Section 3, local edits may allow LLMs as PDDL formalizers to satisfy initial and goal constraints more easily compared to state constraints, since they do not require introducing and applying new predicates.

In contrast, state constraints may be much easier to model using SMT which provides analogous functions. Moreover, PDDL3 may do well on both initial and goal constraints (due to local edits), as well as state constraints since the syntax can now encode constraints. Figure 3 show some yet non-decisive evidence towards those hypotheses. Appendix Figure 6 demonstrates that state constraints are more difficult than action constraints for LTL using DeepSeek models. Furthermore, PDDL3 demonstrated similar results to standard PDDL and in contrast to initial assumptions, performed the worst on state constraints. The most common errors include syntax errors when compiling PDDL3 into PDDL and semantic errors when finding a plan after compilation.

Can LLMs scale with complexity? Recent work has cast doubt on the ability of LLMs, as planner and formalizers, to scale with the problem complexity (Shojaee et al., 2025; Lin et al., 2025; Amonkar et al., 2025). In domains considered in this work, problems from BlocksWorld-100 range from 2 to 15 blocks, while those from CoinCollector-100 range from 3 to 11 rooms. Although our results represented in Figure 2 pose a significant room for improvement in model performance, especially after the introduction of constraints, it would only be more meaningful to conduct studies on problems with higher complexity emulate more

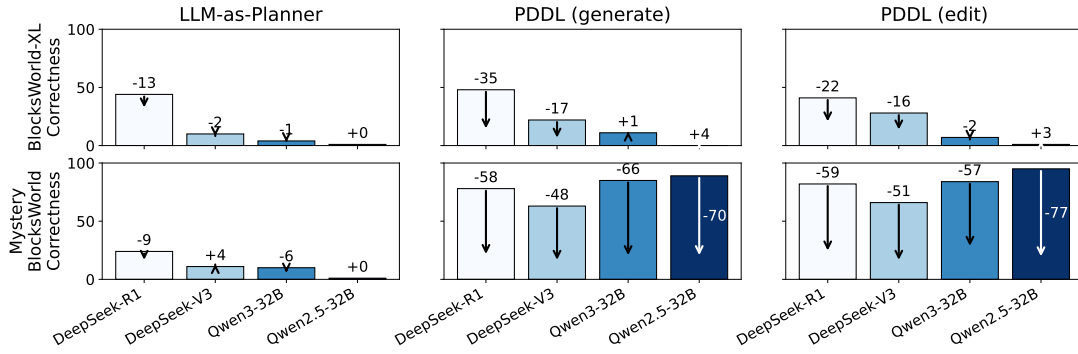


Figure 5: Performance of various methods on BlocksWorld-XL-100 and MysteryBlocksWorld-100 data to study models’ robustness to complexity and lexical perturbation respectively, under the influence of constraints. As before, the arrows represent the performance difference once constraints are added. Results for SMT are omitted due to near-zero performance on both datasets.

to real life applications. Therefore, we present BlocksWorld-XL-100, a dataset included in CoPE where all problems contain 50 blocks, in order to see the impact of problems with a larger entity space on planning, especially with constraints.

Comparing performance of the first row of Figure 5 to that on BlocksWorld-100 (Figure 2), performance for all methods sharply decreases as the problems become more complex even without the introduction of constraints, in line with cited literature. Even so, LLM-as-formalizer displays more robustness to complexity than LLM-as-planner. However, **such scaling behavior disappears with the constraints** as performance degrades by more than 50% for the well-performing DeepSeek models. This finding clearly indicates the challenge caused by the constraints aggravates the existing concern of scaling inability in LLM-assisted planning methodologies.

Do LLMs memorize the domains? Among the domains included in CoPE, BlocksWorld has existed for decades; and the simple semantics of navigation in CoinCollector have existed in many longstanding planning domains. As state-of-the-art LLMs are trained with most publicly available data, data contamination has been a stubborn challenge that hinders research (Sainz et al., 2023). As we have discussed in Section 2, we pose CoPE and the introduction of constraints as a remedy to shift the semantics in the data and bring models out of their comfort zone. We take one step further and apply the lexical obfuscation method of Mystery BlocksWorld (Valmeekam et al., 2024a) to BlocksWorld-100 to construct MysteryBlocksWorld-100, by replacing all the

names of the types, predicates, actions, and objects with nonsensical placeholders. We add our constraints by similarly swapping out the concepts in the constraints and the ground-truth PDDL.

From the second row of Figure 5, the performance without constraints is consistent with that of Huang and Zhang (2025), where LLM-as-planner is devastated compared to the original BlocksWorld-100, suggesting some extent of memorization, while LLM-as-formalizer remains robust to lexical perturbation. However, **such robustness disappears with the constraints**, as the performance for all 4 LLMs and 2 formalizing methods decreases by two thirds or more.

Our results jointly show that LLM-based planning and formalization are far from achieving robust performance on complex problems in niche domains. They also emphasize the importance of future benchmarks like CoPE that transcend simple problems in typical domains.

7 Conclusion

We introduce CoPE, a benchmark that augments existing planning domains with linguistically rich and formally categorized constraints to increase complexity and real-world applicability. Results show that CoPE presents significant challenges and consistently degrades performance. We provide a detailed analysis of behavioral differences of methods and planning languages, suggesting future research on tooling for LLM-based planning. Importantly, our constraints amplify models’ lack of robustness with increasing problem complexity and lexical perturbation, calling for attention to the gap between benchmarks and real-world planning.

8 Limitation

Despite our best effort in defining, formalizing, and categorizing the constraints, we inevitably fall short of the goal of addressing all possible constraints (or more generally, requests) expressed in natural language. Even within our formalism, interesting phenomena such as the conjunction, negation, or ambiguity of the constraints are left for future work.

Our evaluation metric, plan correctness, may induce false positives where the plan may happen to be correct, but the generated code does not actually correctly describe the environment satisfy the constraint, leading to concerns in interpretability and faithfulness. However, to the best of our knowledge there is no feasible alternative. While there exists methods to evaluate generated PDDL problem files (Zuo et al., 2024) or domain files (Coulter et al., 2022) against a ground-truth, fundamentally there exists non-unique ways to formulate the same constraint.

While CoPE augments 4 datasets (all under the MIT License) in the BlocksWorld and CoinCollector domains, they are still oversimplified proof of concept and not representative to problems in the real world. This may pose a risk to users using this code on real world problems. Still, the challenge to state-of-the-method as we have demonstrated argues for the necessity of incremental evaluation.

References

Rikhil Amonkar, May Lai, Ronan Le Bras, and Li Zhang. 2025. *Are llms better formalizers than solvers on complex problems?* *Preprint*, arXiv:2505.13252.

Ujjwala Anantheswaran, Himanshu Gupta, Kevin Scaria, Shreyas Verma, Chitta Baral, and Swaroop Mishra. 2024. *Cutting through the noise: Boosting llm performance on math word problems.*

Luigi Bonassi, Alfonso Emilio Gerevini, Francesco Percassi, and Enrico Scala. 2021. *On planning with qualitative state-trajectory constraints in pddl3 by compiling them away.* *Proceedings of the International Conference on Automated Planning and Scheduling*, 31(1):46–50.

Alex Coulter, Teo Ilie, Renee Tibando, and Christian Muise. 2022. *Theory alignment via a classical encoding of regular bisimulation.* In *Workshop on Knowledge Engineering for Planning and Scheduling (KEPS)*.

Alexandre Duret-Lutz, Etienne Renault, Maximilien Colange, Florian Renkin, Alexandre Gbaguidi Aisse,

Philipp Schlehuber-Caissier, Thomas Medioni, Antoine Martin, Jérôme Dubois, Clément Gillard, and Henrich Lauko. 2022. *From Spot 2.0 to Spot 2.10: What’s new?* In *Proceedings of the 34th International Conference on Computer Aided Verification (CAV’22)*, volume 13372 of *Lecture Notes in Computer Science*, pages 174–187. Springer.

Stefan Edelkamp. 2006. *On the compilation of plan constraints and preferences.* In *Proceedings of the Sixteenth International Conference on International Conference on Automated Planning and Scheduling*, ICAPS’06, page 374–377. AAAI Press.

Richard E. Fikes and Nils J. Nilsson. 1971. *Strips: a new approach to the application of theorem proving to problem solving.* In *Proceedings of the 2nd International Joint Conference on Artificial Intelligence*, IJCAI’71, page 608–620, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. *Pal: Program-aided language models.* In *International Conference on Machine Learning*, pages 10764–10799. PMLR.

Alfonso Gerevini and Derek Long. 2005. *Plan constraints and preferences in pddl3.* Technical report, Technical Report 2005-08-07, Department of Electronics for Automation

Liancheng Gong, Wang Zhu, Jesse Thomason, and Li Zhang. 2025. *Zero-shot iterative formalization and planning in partially observable environments.* *Preprint*, arXiv:2505.13126.

Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. 2023. *Leveraging pre-trained large language models to construct and utilize world models for model-based task planning.* *Advances in Neural Information Processing Systems*, 36:79081–79094.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. *Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.* *arXiv preprint arXiv:2501.12948*.

Weihang Guo, Zachary Kingston, and Lydia E. Kavraki. 2024. *Castl: Constraints as specifications through llm translation for long-horizon task and motion planning.* *Preprint*, arXiv:2410.22225.

Yilun Hao, Yang Zhang, and Chuchu Fan. 2025. *Planning anything with rigor: General-purpose zero-shot planning with LLM-based formalized programming.* In *The Thirteenth International Conference on Learning Representations*.

Xingyun Hong, Yan Shao, Zhilin Wang, Manni Duan, and Xiongnan Jin. 2024. *Towards building a robust knowledge intensive question answering model with large language models.* *ArXiv*, abs/2409.05385.

651	R. Howey, D. Long, and M. Fox. 2004. Val: automatic plan validation, continuous effects and mixed initiative planning using pddl . In <i>16th IEEE International Conference on Tools with Artificial Intelligence</i> , pages 294–301.	706
652		707
653		708
654		709
655		710
656	Mengkang Hu, Tianxing Chen, Yude Zou, Yuheng Lei, Qiguang Chen, Ming Li, Yao Mu, Hongyuan Zhang, Wenqi Shao, and Ping Luo. 2025. Text2world: Benchmarking large language models for symbolic world model generation . <i>Preprint</i> , arXiv:2502.13092.	711
657		712
658		713
659		714
660		
661		
662	Cassie Huang and Li Zhang. 2025. On the limit of language models as planning formalizers . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4880–4904, Vienna, Austria. Association for Computational Linguistics.	715
663		716
664		717
665		718
666		719
667		720
668	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. 2024. Qwen2. 5-coder technical report. <i>arXiv preprint arXiv:2409.12186</i> .	721
669		722
670		723
671		724
672	IPC. 1998. International planning competition. https://www.icaps-conference.org/competitions .	725
673		726
674	Adam Ishay and Joohyung Lee. 2025. Llm+al: Bridging large language models and action languages for complex reasoning about actions . <i>Preprint</i> , arXiv:2501.00830.	727
675		728
676		729
677		730
678	Peter A. Jansen and Marc-Alexandre Côté. 2022. Textworldexpress: Simulating text games at one million steps per second . <i>arXiv</i> .	731
679		732
680		733
681	Prabhu Prakash Kagitha, Andrew Zhu, and Li Zhang. 2025. Addressing the challenges of planning language generation . <i>Preprint</i> , arXiv:2505.14763.	734
682		735
683		736
684	Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Saldyt, and Anil Murthy. 2024. Llms can’t plan, but can help planning in llm-modulo frameworks. <i>arXiv preprint arXiv:2402.01817</i> .	737
685		738
686		739
687		740
688		
689	Manling Li, Shiyu Zhao, Qineng Wang, Kangrui Wang, Yu Zhou, Sanjana Srivastava, Cem Gokmen, Tony Lee, Li Erran Li, Ruohan Zhang, et al. 2024. Embodied agent interface: Benchmarking llms for embodied decision making. In <i>NeurIPS 2024</i> .	741
690		742
691		743
692		744
693		745
694	Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter Clark, and Yejin Choi. 2025. ZebraLogic: On the scaling limits of LLMs for logical reasoning . In <i>Forty-second International Conference on Machine Learning</i> .	746
695		747
696		748
697		749
698		750
699	Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Prithviraj Ammanabrolu, Yejin Choi, and Xiang Ren. 2023. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	751
700		752
701		753
702		754
703		755
704		756
705		757
		758
		759
		760
	Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. 2023a. Llm+ p: Empowering large language models with optimal planning proficiency. <i>arXiv preprint arXiv:2304.11477</i> .	
	Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023b. Evaluating the logical reasoning ability of chatgpt and gpt-4. <i>arXiv preprint arXiv:2304.03439</i> .	
	Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning . In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 305–329, Nusa Dua, Bali. Association for Computational Linguistics.	
	Bodhisattwa Prasad Majumder, Bhavana Dalvi Mishra, Peter Jansen, Oyvind Tafjord, Niket Tandon, Li Zhang, Chris Callison-Burch, and Peter Clark. 2023. Clin: A continually learning language agent for rapid task adaptation and generalization . <i>Preprint</i> , arXiv:2310.10134.	
	Christian Muise. 2016. Planning.Domains . In <i>The 26th International Conference on Automated Planning and Scheduling - Demonstrations</i> .	
	Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-LM: Empowering large language models with symbolic solvers for faithful logical reasoning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 3806–3824, Singapore. Association for Computational Linguistics.	
	Mihir Parmar, Xin Liu, Palash Goyal, Yanfei Chen, Long Le, Swaroop Mishra, Hossein Mobahi, Jindong Gu, Zifeng Wang, Hootan Nakhost, Chitta Baral, Chen-Yu Lee, Tomas Pfister, and Hamid Palangi. 2025. Plangen: A multi-agent framework for generating planning and reasoning trajectories for complex problem solving . <i>Preprint</i> , arXiv:2502.16111.	
	Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10776–10787, Singapore. Association for Computational Linguistics.	
	Keisuke Shirai, Cristian C. Beltran-Hernandez, Masashi Hamaya, Atsushi Hashimoto, Shohei Tanaka, Kento Kawaharazuka, Kazutoshi Tanaka, Yoshitaka Ushiku, and Shinsuke Mori. 2024. Vision-language interpreter for robot task planning . <i>Preprint</i> , arXiv:2311.00967.	

761	Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. 2025. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity . <i>Preprint</i> , arXiv:2506.06941.	814
762		815
763		816
764		817
765		818
766		819
767	Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. 2024. Chain of thoughtlessness: An analysis of cot in planning. <i>arXiv preprint arXiv:2405.04776</i> .	820
768		821
769		822
770		823
771	Katharina Stein, Daniel Fišer, Jörg Hoffmann, and Alexander Koller. 2025. Automating the generation of prompts for llm-based action choice in pdpl planning . <i>Preprint</i> , arXiv:2311.09830.	824
772		825
773		826
774		827
775	Jing Han Sun and Ali Emami. 2024. Evograd: A dynamic take on the winograd schema challenge with human adversaries . In <i>International Conference on Language Resources and Evaluation</i> .	828
776		829
777		830
778		831
779	Hao Tang, Darren Key, and Kevin Ellis. 2024. World-coder, a model-based llm agent: Building world models by writing code and interacting with the environment . <i>Preprint</i> , arXiv:2402.12275.	832
780		833
781		834
782		835
783	Marcus Tantakoun, Christian Muise, and Xiaodan Zhu. 2025. LLMs as planning formalizers: A survey for leveraging large language models to construct automated planning models . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 25167–25188, Vienna, Austria. Association for Computational Linguistics.	836
784		837
785		838
786		839
787		840
788		841
789		842
790	Qwen Team. 2025. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388.	843
791		844
792	Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2024a. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. <i>Advances in Neural Information Processing Systems</i> , 36.	845
793		846
794		847
795		848
796		849
797		850
798	Karthik Valmeekam, Kaya Stechly, and Subbarao Kambhampati. 2024b. Llms still can't plan; can lrms? a preliminary evaluation of openai's o1 on planbench . <i>Preprint</i> , arXiv:2409.13373.	851
799		852
800		853
801		854
802	Ruoyao Wang, Graham Todd, Eric Yuan, Ziang Xiao, Marc-Alexandre Côté, and Peter Jansen. 2023. Byte-sized32: A corpus and challenge task for generating task-specific world models expressed as text games . <i>Preprint</i> , arXiv:2305.14879.	855
803		856
804		857
805		858
806		859
807	Hui Wei, Zihao Zhang, Shenghua He, Tian Xia, Shijia Pan, and Fei Liu. 2025. PlanGenLLMs: A modern survey of LLM planning capabilities . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 19497–19521, Vienna, Austria. Association for Computational Linguistics.	860
808		861
809		862
810		863
811		864
812		865
813		866
		867
		868
		869
		870
	Lionel Wong, Jiayuan Mao, Pratyusha Sharma, Zachary S Siegel, Jiahai Feng, Noa Korneev, Joshua B Tenenbaum, and Jacob Andreas. 2023. Learning adaptive planning representations with natural language guidance . <i>arXiv preprint arXiv:2312.08566</i> .	814
		815
		816
		817
		818
		819
	Yaqi Xie, Chen Yu, Tongyao Zhu, Jinbin Bai, Ze Gong, and Harold Soh. 2023. Translating natural language to planning goals with large-language models . <i>arXiv preprint arXiv:2302.05128</i> .	820
		821
		822
		823
	Ziyi Yang, Shreyas S. Raman, Ankit Shah, and Stefanie Tellex. 2023. Plug in the safety chip: Enforcing constraints for llm-driven robot agents . <i>Preprint</i> , arXiv:2309.09919.	824
		825
		826
		827
	Jianzhu Yao, Kevin Wang, Ryan Hsieh, Haisu Zhou, Tianqing Zou, Zerui Cheng, Zhangyang Wang, and Pramod Viswanath. 2025. Spin-bench: How well do llms plan strategically and reason socially? <i>Preprint</i> , arXiv:2503.12349.	828
		829
		830
		831
		832
	Xingdi Yuan, Marc-Alexandre Côté, Alessandro Sordani, Romain Laroche, Remi Tachet des Combes, Matthew Hausknecht, and Adam Trischler. 2019. Counting to explore and generalize in text-based games . <i>Preprint</i> , arXiv:1806.11525.	833
		834
		835
		836
		837
	Li Zhang, Peter Jansen, Tianyi Zhang, Peter Clark, Chris Callison-Burch, and Niket Tandon. 2024a. PDDLEGO: Iterative planning in textual environments . In <i>Proceedings of the 13th Joint Conference on Lexical and Computational Semantics (*SEM 2024)</i> , pages 212–221, Mexico City, Mexico. Association for Computational Linguistics.	838
		839
		840
		841
		842
		843
		844
	Tianyi Zhang, Li Zhang, Zhaoyi Hou, Ziyu Wang, Yuling Gu, Peter Clark, Chris Callison-Burch, and Niket Tandon. 2024b. PROC2PDDL: Open-domain planning representations from texts . In <i>Proceedings of the 2nd Workshop on Natural Language Reasoning and Structured Explanations (@ACL 2024)</i> , pages 13–24, Bangkok, Thailand. Association for Computational Linguistics.	845
		846
		847
		848
		849
		850
		851
		852
	Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang, Xinyun Chen, Minmin Chen, Azade Nova, Le Hou, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and Denny Zhou. 2024. Natural plan: Benchmarking llms on natural language planning .	853
		854
		855
		856
		857
	Yinan Zheng, Ruiming Liang, Kexin ZHENG, Jinliang Zheng, Liyuan Mao, Jianxiong Li, Weihao Gu, Rui Ai, Shengbo Eben Li, Xianyuan Zhan, and Jingjing Liu. 2025. Diffusion-based planning for autonomous driving with flexible guidance . In <i>The Thirteenth International Conference on Learning Representations</i> .	858
		859
		860
		861
		862
		863
	Andrew Zhu, Liam Dugan, Alyssa Hwang, and Chris Callison-Burch. 2023. Kani: A lightweight and highly hackable framework for building language model applications . In <i>Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023)</i> , pages 65–77, Singapore. Association for Computational Linguistics.	864
		865
		866
		867
		868
		869
		870

871 Wang Zhu, Ishika Singh, Robin Jia, and Jesse Thoma-
 872 son. 2024. Language models can infer action seman-
 873 tics for classical planners from environment feedback.
 874 *arXiv preprint arXiv:2406.02791*.

875 Max Zuo, Francisco Piedrahita Velez, Xiaochen Li,
 876 Michael L Littman, and Stephen H Bach. 2024. Plan-
 877 etarium: A rigorous benchmark for translating text
 878 to structured planning languages. *arXiv preprint*
 879 *arXiv:2407.03321*.

880 A Benchmark Examples

881 Listings 1 and 2 display a benchmark example of
 882 BlocksWorld from Valmeekam et al. (2024a) and
 883 (Huang and Zhang, 2025). Listings 3 and 4 display
 884 an example of the Logistics domain from Huang
 885 and Zhang (2025). Language seen in this bench-
 886 mark is simplistic and not not realistic to the real
 887 world.

888 B Data Examples

889 We augment constraint descriptions to a ground-
 890 truth domain and problem file and create new
 891 ground-truth PDDL files that now encode the
 892 ground-truth. Here we include an example con-
 893 straint description and the ground-truth PDDL files
 894 for the constraint + problem pair. Listings 5, 6 and
 895 7 display an example constraint description and the
 896 annotated ground-truth \mathcal{DF} and \mathcal{PF} for a paired
 897 BlocksWorld problem. Listings 8, 9 and 10 display
 898 an example Domain, Problem and Constraint De-
 899 scription for CoinCollector-100 as well as their
 900 groundtruth \mathcal{DF} and \mathcal{PF} in PDDL. Listing 11
 901 shows an example ground-truth \mathcal{PF} for the XL
 902 dataset, the \mathcal{DF} remains the same. Listings 12,
 903 13 and 14 display the corresponding constraint de-
 904 scription, \mathcal{DF} and \mathcal{PF} that have now been obfus-
 905 cated.

906 C Prompts

907 Listing 15 is the prompt for LLM-as-Planner. List-
 908 ings 16, 17, 18 and 19 display the prompts for all
 909 experiment settings given to all models for LLM-as-
 910 PDDL-Formalizer. Listings 20, 21, 22 and 23 dis-
 911 play the prompts for all experiment settings given
 912 to all models for LLM-as-PDDL3-Formalizer. List-
 913 ings 24, 25 and 26 display prompts for experiments
 914 using SMT. Listing 27 displays all prompts for LTL.
 915 For PDDL and SMT, we asked the model to return
 916 the output in a JSON object whenever possible for
 917 easier parsing.

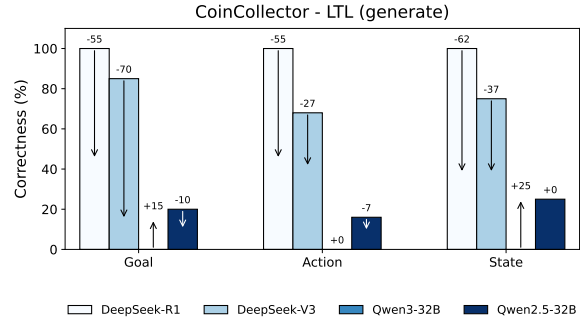


Figure 6: Performance of LTL on a subset of CoinCollector-100 problems. The arrows represent the performance change with constraints. Details on design choice are discussed in Appendix D.3.

918 D Experimental Setup Details

919 D.1 Planner

920 We use a dual-bfws-ffparser planner to find
 921 plans from the generated PDDL. This is a best
 922 first width search planner used with a fast forward
 923 planner.

924 D.2 VAL

925 The VAL library takes in a ground-truth PDDL \mathcal{DF} ,
 926 \mathcal{PF} and a plan and tries to execute the plan in the
 927 environment by checking whether each action in
 928 the found plan can be executed based on the pre-
 929 conditions written in the ground-truth files. If the plan
 930 is executable, VAL then checks whether the final
 931 state after executing the plan matches the goal
 932 state found in the ground-truth \mathcal{PF} . If either the plan is
 933 not executable or the final state does not match the
 934 goal state, VAL will return an error, therefore the
 935 plan found by the planner is not correct.

936 D.3 LTL

937 Given problem description, the language model
 938 is asked to generate LTL formulas for the initial
 939 state, goal state, environment dynamics, and action
 940 dynamics. All formulas are concatenated into a
 941 conjunctive form as the formal representation for
 942 planning. When a constraint are introduced, we
 943 make a separate call to the language model for
 944 the formula representing the constraint and con-
 945 catenate the formula to the existing conjunctive
 946 formula. We employ the spot library (Duret-Lutz
 947 et al., 2022) for encoding LTL formulas into Büchi
 948 automata. Then, the state traces satisfying the
 949 problem specification can be found by searching
 950 for a Lasso-shaped accepting run. To make the
 951 results compatible with our evaluation pipeline, we

Domains Used		
LLM-as-formalizer	Huang and Zhang (2025)	BlocksWorld, Logistics, Barman
	Kagitha et al. (2025)	BlocksWorld, Logistics, Barman
	Gong et al. (2025)	CoinCollector, ALFWorld
	Hu et al. (2025)	Over 100 domains including Grid and BlocksWorld
	Li et al. (2024)	BEHAVIOR, VirtualHome
	Zuo et al. (2024)	BlocksWorld, Gripper
	Zhu et al. (2024)	Barman, BlocksWorld, Floortile, Grippers, Storage, Termes, TyreWorld
	Zhang et al. (2024a)	CoinCollector, CookingWorld
	Zhang et al. (2024b)	Proc2PDDL
	Liu et al. (2023a)	Barman, BlocksWorld, Floortile, Grippers, Storage, Termes, TyreWorld
	Xie et al. (2023)	BlocksWorld, ALFRED-L
	Wong et al. (2023)	Mini Minecraft, ALFRED
	Guan et al. (2023)	Household, Logistics, TyreWorld
	Shirai et al. (2024)	Cooking, BlocksWorld, Tower of Hanoi
	Ishay and Lee (2025)	MCP
	Hao et al. (2025)	Coffee, Workforce, Facility, Task Allocation, Warehouse, BlocksWorld, Mystery BlocksWorld, Movie, Gripper
	Guo et al. (2024)	HouseChip, Kitchen, BlocksWorld
	Tang et al. (2024)	Sokoban, Minigrad, ALFWorld
	Wang et al. (2023)	ByteSized32
	Amonkar et al. (2025)	NaturalPlan, ZebraLogic
Lyu et al. (2023)	SayCan	
Gao et al. (2023)	GSM8K	
LLM-as-planner	Shojaee et al. (2025)	Tower of Hanoi, Checker Jumping, BlocksWorld, River Crossing
	Parmar et al. (2025)	NaturalPlan, GPQA, OlympiadBench
	Majumder et al. (2023)	ScienceWorld
	Lin et al. (2023)	ScienceWorld
	Stein et al. (2025)	BlocksWorld, Depot, Logistics, Gripper, Ferry, Floortile, Goldminer, Grid, Gripper, Movie, Rovers, Satellite, Visitall
	Zheng et al. (2024)	NaturalPlan
	Stechly et al. (2024)	BlocksWorld
	Valmeekam et al. (2024a)	BlocksWorld, Logistics, Mystery BlocksWorld
	Yao et al. (2025)	BlocksWorld, Gripper, Barman, Logistics and more

Table 2: Comparison with related work.

map the traces back to sequences of PDDL actions by inferring actions from abstract state changes. In this process, the language model is also asked to generate environmental information, which is the room connectivities in CoinCollector domain.

We prompt the LLM to generate formulas for goal, action, and state constraints on CoinCollector-100 only, since overriding initial states by editing is practically infeasible when all components of the environment and the problem are represented uniformly in a conjunctive formula. Models exhibit a similar performance degradation between 30-50% just as PDDL and SMT (Figure 6). Notably, different from PDDL and SMT, LTL is restricted to domains with simple action dynamics, e.g., navigational domains like CoinCollector, so while it performs well without constraints, it is not as generalizable as other languages.

E Full Results

Tables 3, 4, 5, 6 and 7 display results for all methodologies on BlocksWorld-100

and CoinCollector-100. Tables 8 and 9 display results for BlocksWorld-XL-100 and MysteryBlocksWorld-100 on PDDL.

CoinCollector-100

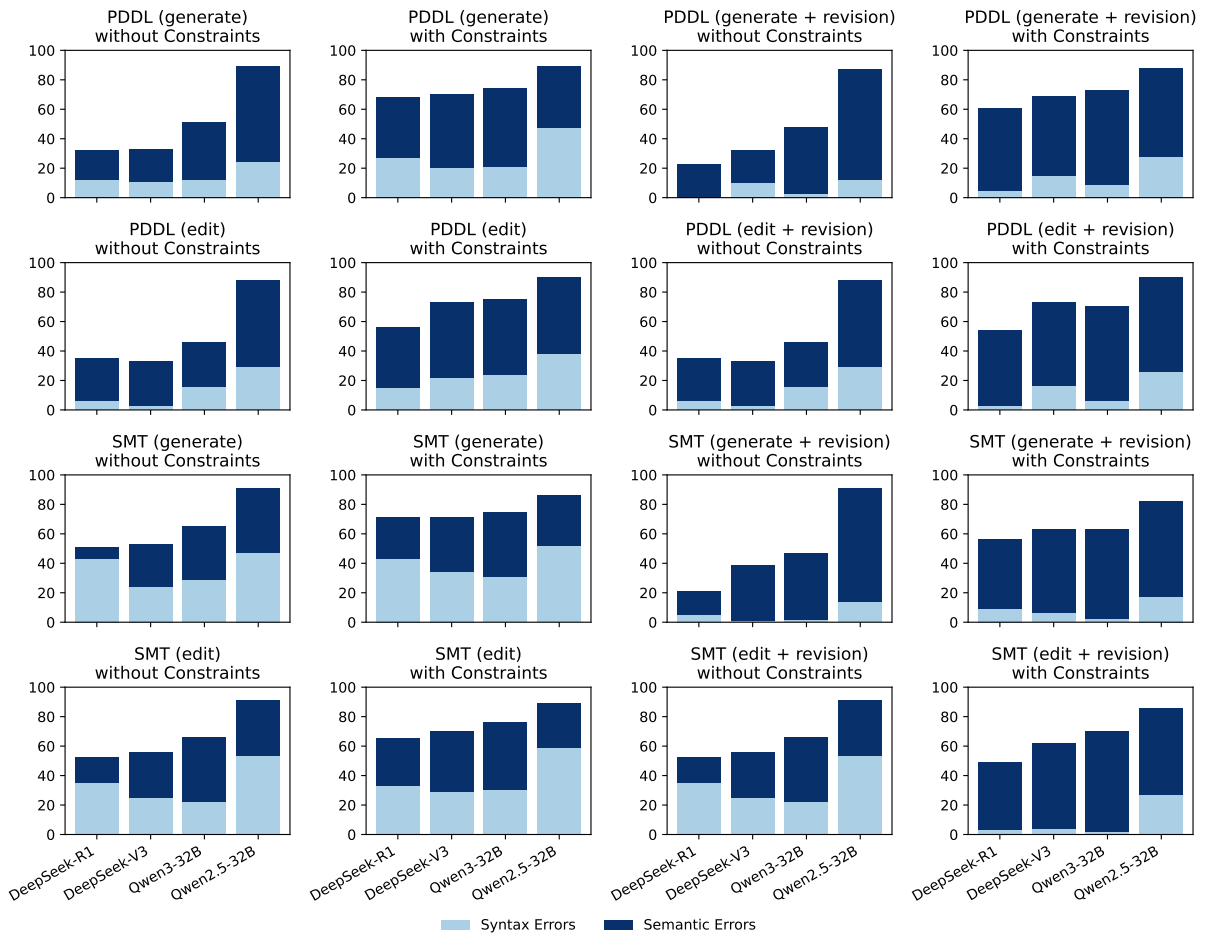


Figure 7: Counts of Syntax vs. Semantic Errors with and without Constraints on CoinCollector-100 using both PDDL and SMT.

Listing 1: Domain Description for Moderately Templated BlocksWorld-100

I am playing with a set of blocks where I need to arrange the blocks into stacks. Here are the actions I can do

Pick up a block
 Unstack a block from on top of another block
 Put down a block
 Stack a block on top of another block

I have the following restrictions on my actions:
 I can only pick up or unstack one block at a time.
 I can only pick up or unstack a block if my hand is empty.
 I can only pick up a block if the block is on the table and the block is clear. A block is clear if the block has no other blocks on top of it and if the block is not picked up.
 I can only unstack a block from on top of another block if the block I am unstacking was really on top of the other block.
 I can only unstack a block from on top of another block if the block I am unstacking is clear.
 Once I pick up or unstack a block, I am holding the block.
 I can only put down a block that I am holding.
 I can only stack a block on top of another block if I am holding the block being stacked.
 I can only stack a block on top of another block if the block onto which I am stacking the block is clear.
 Once I put down or stack a block, my hand becomes empty.
 Once you stack a block on top of a second block, the second block is no longer clear.

976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
~~1000~~

Listing 2: Problem Description for Moderately Templated BlocksWorld-100

As initial conditions I have that, the blue block is clear, the yellow block is clear, the hand is empty, the blue block is on top of the red block, the yellow block is on top of the green block, the red block is on the table, and the green block is on the table.
 My goal is to have that the blue block is on top of the yellow block, the green block is on top of the red block, the yellow block is on top of the green block, and the red block is on the table.

1001
 1002
 1003
 1004
 1005
~~1006~~

Listing 3: Domain Description for Moderately Templated Logistics-100

I need to move packages between locations. Here are the actions I can do

Load an package onto a truck at a location (load-truck package truck location)
 Load an package onto an airplane at a location (load-airplane package airplane location)
 Unload an package from a truck at a location (unload-truck package truck location)
 Unload an package from an airplane at a location (unload-airplane package airplane location)
 Drive a truck from location1 to location2 in a city (drive-truck truck location1 location2 city)
 Fly an airplane from airport1 to airport2 (fly-airplane airplane airport1 airport2)

I have the following restrictions on my actions:
 I can only load a package onto a truck or airplane if both the package and airplane are at the location.
 Once I load the package in the truck or airplane, it is no longer at the location.
 I can only unload a package from a truck or airplane if the truck or airplane is at the location and the package is in the truck or airplane.
 Once I unload the truck or airplane, the object is at the location and no longer in the truck or airplane.
 I can only drive a truck between locations if the truck is at the first location and both the first and second locations are in the same city. Once I drive a truck, the truck is in the second city and no longer in the first city.
 I can only fly an airplane between two airports and the airplane is at the first airport.
 Once I fly an airplane, the airplane is at the second airport and no longer at the first airport.

1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025
 1026
 1027
 1028
 1029
~~1030~~

Listing 4: Problem Description for Moderately Templated Logistics-100

As initial conditions, I have that, obj11 is a package, obj12 is a package, obj13 is a package, obj21 is a package, obj22 is a package, obj23 is a package, tru1 is a truck, tru2 is a truck, cit1 is a city, cit2 is a city, pos1 is a location, apt1 is a location, pos2 is a location, apt2 is a location, apt1 is an airport, apt2 is an airport, apn1 is an airplane, apn1 is at apt2, tru1 is at pos1, obj11 is at pos1, obj12 is at pos1, obj13 is at pos1, tru2 is at pos2, obj21 is at pos2, obj22 is at pos2, obj23 is at pos2, pos1 is in cit1, apt1 is in cit1, pos2 is in cit2, and apt2 is in cit2.

1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039

1049 My goal is to have that obj11 is at apt1, obj23 is at pos1, obj13 is at apt1, and obj21 is at pos1.

Listing 5: Example constraint description for an action constraint

1042 Do not stack block1 on top of block2
1044

Listing 6: Annotated ground-truth Domain File for BlocksWorld-100 with action constraint description

```
1045 (define (domain blocksworld)
1046 ;; CONSTRAINT Do not stack block1 on top of block2.
1047 (:requirements :strips)
1048 (:predicates (clear ?x)
1049              (on-table ?x)
1050              (arm-empty)
1051              (holding ?x)
1052              (on ?x ?y)
1053 ;; BEGIN ADD
1054              (do-not-stack ?x ?y)
1055 ;; END ADD
1056 )
1057
1058 (:action pickup
1059 :parameters (?ob)
1060 :precondition (and (clear ?ob) (on-table ?ob) (arm-empty))
1061 :effect (and (holding ?ob) (not (clear ?ob)) (not (on-table ?ob))
1062           (not (arm-empty))))
1063
1064 (:action putdown
1065 :parameters (?ob)
1066 :precondition (holding ?ob)
1067 :effect (and (clear ?ob) (arm-empty) (on-table ?ob)
1068           (not (holding ?ob))))
1069
1070 (:action stack
1071 :parameters (?ob ?underob)
1072 :precondition (and (clear ?underob) (holding ?ob)
1073 ;; BEGIN ADD
1074                 (not (do-not-stack ?ob ?underob))
1075 ;; END ADD
1076 )
1077 :effect (and (arm-empty) (clear ?ob) (on ?ob ?underob)
1078           (not (clear ?underob)) (not (holding ?ob))))
1079
1080 (:action unstack
1081 :parameters (?ob ?underob)
1082 :precondition (and (on ?ob ?underob) (clear ?ob) (arm-empty))
1083 :effect (and (holding ?ob) (clear ?underob)
1084           (not (on ?ob ?underob)) (not (clear ?ob)) (not (arm-empty))))
1086
```

Listing 7: Annotated ground-truth Problem File for BlocksWorld-100 with action constraint description

```
1087 (define (problem blocksworld-p50)
1088 ;; CONSTRAINT Do not stack block1 on top of block2.
1089 (:domain blocksworld)
1090 (:objects block1 block2 block3 block4 block5 block6 block7 )
1091 (:init
1092   (on-table block3)
1093   (clear block3)
1094   (on-table block5)
1095   (clear block5)
1096   (on-table block7)
1097   (on block1 block7)
1098   (clear block1)
1099   (on-table block2)
1100   (clear block2)
1101   (on-table block6)
1102   (clear block6)
1103
```

```

(on-table block4)
(clear block4)
(arm-empty)
;; BEGIN ADD
(do-not-stack block1 block2)
;; END ADD
)
(:goal (and
(on-table block5)
(on block2 block5)
(on-table block4)
(on block7 block4)
(on block6 block7)
(on block1 block6)
(on-table block3)
))
)

```

Listing 8: Example Domain, Problem and Action Constraint Description for CoinCollector-100

```

;; Domain Description
You are standing in an environment with many rooms. Your task is to search the environment and find
the coin. Once you find the coin, take it.

The following actions are allowed: move from one room to another and take an item.

;; Problem Description
You are in the kitchen.
To the east of the kitchen is the pantry.
To the west of the kitchen is the corridor.
There is a coin in the pantry.

;; Constraint Description
You must visit the corridor before the pantry.

```

Listing 9: Annotated ground-truth Domain File for CoinCollector problem with action constraint

```

(define (domain coin-collector)
;; CONSTRAINT You must visit the corridor before the pantry.
(:requirements :strips :typing)
(:types
room
direction
item
)
(:predicates
(at ?room - room)
(connected ?room1 - room ?room2 - room ?direction - direction)
(closed-door ?room1 - room ?room2 - room ?direction - direction)
(location ?item - item ?room - room)
(taken ?item - item)
(is-reverse ?direction - direction ?reverse - direction)
)
;; BEGIN ADD
(must-visit-before-next ?room - room)
(next-room ?room - room)
(visited-room)
;; END ADD
)

(:action move
:parameters (?room1 - room ?room2 - room ?direction - direction)
:precondition (and (at ?room1) (connected ?room1 ?room2 ?direction))
;; BEGIN ADD
(not (and (next-room ?room2) (not (visited-room))))
;; END ADD
)
:effect (and (not (at ?room1)) (at ?room2))
;; BEGIN ADD

```

```

1170   (when (must-visit-before-next ?room2) (visited-room))
1171 ;; END ADD
1172   )
1173   )
1174
1175   (:action take
1176    :parameters (?item - item ?room - room)
1177    :precondition (and (at ?room) (location ?item ?room) (not (taken ?item)))
1178    :effect (and (taken ?item) (not (location ?item ?room)))
1179   )
1180 )

```

Listing 10: Annotated ground-truth Problem File for CoinCollector problem with action constraint

```

1182 (define (problem coin_collector_numLocations3_numDistractorItems0_seed26)
1183 ;; CONSTRAINT You must visit the corridor before the pantry.
1184 (:domain coin-collector)
1185 (:objects
1186  kitchen pantry corridor - room
1187  north south east west - direction
1188  coin - item
1189 )
1190 (:init
1191  (at kitchen)
1192  (connected kitchen pantry east)
1193  (connected kitchen corridor west)
1194  (connected pantry kitchen west)
1195  (connected corridor kitchen east)
1196  (location coin pantry)
1197  (is-reverse north south)
1198  (is-reverse south north)
1199  (is-reverse east west)
1200  (is-reverse west east)
1201 )
1202 ;; BEGIN ADD
1203 (must-visit-before-next corridor)
1204 (next-room pantry)
1205 ;; END ADD
1206 )
1207 (:goal
1208  (taken coin)
1209 )
1210 )

```

Listing 11: Annotated ground-truth Problem File for the XL BlocksWorld-100 problems with action constraint description

```

1212 (define (problem blocksworld-p50)
1213 ;; CONSTRAINT Do not stack block1 on top of block2.
1214 (:domain blocksworld)
1215 (:objects block1 block2 block3 block4 block5 block6 block7 block8 block9 block10 block11 block12
1216  block13 block14 block15 block16 block17 block18 block19 block20 block21 block22 block23
1217  block24 block25 block26 block27 block28 block29 block30 block31 block32 block33 block34
1218  block35 block36 block37 block38 block39 block40 block41 block42 block43 block44 block45
1219  block46 block47 block48 block49 block50 )
1220 (:init
1221  (on-table block10)
1222  (on block41 block10)
1223  (on block50 block41)
1224  (on block29 block50)
1225  (on block46 block29)
1226  (on block43 block46)
1227  (on block5 block43)
1228  (on block38 block5)
1229  (on block48 block38)
1230  (on block8 block48)
1231  (on block11 block8)
1232  (on block22 block11)
1233  (on block13 block22)
1234 )

```

(on block35 block13)	1235
(on block49 block35)	1236
(on block20 block49)	1237
(on block31 block20)	1238
(on block34 block31)	1239
(on block17 block34)	1240
(on block28 block17)	1241
(on block14 block28)	1242
(on block47 block14)	1243
(on block26 block47)	1244
(on block6 block26)	1245
(on block4 block6)	1246
(on block25 block4)	1247
(on block9 block25)	1248
(on block23 block9)	1249
(on block15 block23)	1250
(on block21 block15)	1251
(on block18 block21)	1252
(on block39 block18)	1253
(on block33 block39)	1254
(on block1 block33)	1255
(on block3 block1)	1256
(on block2 block3)	1257
(on block44 block2)	1258
(on block16 block44)	1259
(on block45 block16)	1260
(on block7 block45)	1261
(on block40 block7)	1262
(on block24 block40)	1263
(on block27 block24)	1264
(on block36 block27)	1265
(on block19 block36)	1266
(clear block19)	1267
(on-table block32)	1268
(clear block32)	1269
(on-table block42)	1270
(on block37 block42)	1271
(on block30 block37)	1272
(on block12 block30)	1273
(clear block12)	1274
(arm-empty)	1275
;; BEGIN ADD	1276
(do-not-stack block1 block2)	1277
;; END ADD	1278
)	1279
(:goal (and	1280
(on-table block46)	1281
(on block45 block46)	1282
(on block37 block45)	1283
(on-table block30)	1284
(on block39 block30)	1285
(on-table block34)	1286
(on-table block49)	1287
(on-table block19)	1288
(on-table block38)	1289
(on-table block16)	1290
(on-table block25)	1291
(on-table block18)	1292
(on block5 block18)	1293
(on block17 block5)	1294
(on-table block42)	1295
(on-table block3)	1296
(on block32 block3)	1297
(on block35 block32)	1298
(on block26 block35)	1299
(on-table block12)	1300
(on-table block33)	1301
(on-table block14)	1302
(on-table block2)	1303
(on-table block23)	1304

1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1333

```
(on block1 block23)
(on block6 block1)
(on-table block4)
(on block20 block4)
(on-table block48)
(on-table block27)
(on-table block31)
(on-table block44)
(on block47 block44)
(on-table block43)
(on-table block36)
(on-table block7)
(on-table block22)
(on-table block28)
(on-table block9)
(on-table block11)
(on-table block50)
(on-table block24)
(on-table block41)
(on-table block8)
(on-table block21)
(on-table block40)
(on-table block13)
(on-table block15)
(on-table block29)
(on-table block10)
))
)
```

Listing 12: Example constraint description for a action constraint for Mystery BlocksWorld

1334
1336

```
Do not perform action3 with object1 and object2.
```

Listing 13: Annotated ground-truth Domain File for Mystery BlocksWorld-100 with action constraint description

1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370

```
(define (domain mystery_blocksworld)
;; CONSTRAINT Do not perform action3 with object1 and object2.
(:requirements :strips)
(:predicates (predicate1 ?x)
              (predicate2 ?x)
              (predicate3)
              (predicate4 ?x)
              (predicate5 ?x ?y)
;; BEGIN ADD
              (not-allowed ?x ?y)
;; END ADD
)

(:action action1
:parameters (?object1)
:precondition (and (predicate1 ?object1) (predicate2 ?object1) (predicate3))
:effect (and (predicate4 ?object1) (not (predicate1 ?object1)) (not (predicate2 ?object1))
            (not (predicate3))))

(:action action2
:parameters (?object1)
:precondition (predicate4 ?object1)
:effect (and (predicate1 ?object1) (predicate3) (predicate2 ?object1)
            (not (predicate4 ?object1))))

(:action action3
:parameters (?object1 ?object2)
:precondition (and (predicate1 ?object2) (predicate4 ?object1)
;; BEGIN ADD
              (not (not-allowed ?object1 ?object2))
;; END ADD
)
:effect (and (predicate3) (predicate1 ?object1) (predicate5 ?object1 ?object2))
```

```

        (not (predicate1 ?object2)) (not (predicate4 ?object1))))
(:action action4
 :parameters (?object1 ?object2)
 :precondition (and (predicate5 ?object1 ?object2) (predicate1 ?object1) (predicate3))
 :effect (and (predicate4 ?object1) (predicate1 ?object2)
              (not (predicate5 ?object1 ?object2)) (not (predicate1 ?object1)) (not (predicate3))))))

```

Listing 14: Annotated ground-truth Problem File for Mystery BlocksWorld-100 with action constraint description

```

(define (problem mystery_blocksworld-p50)
;; CONSTRAINT Do not perform action3 with object1 and object2.
 (:domain mystery_blocksworld)
 (:objects object1 object2 object3 object4 object5 object6 object7 )
 (:init
  (predicate2 object3)
  (predicate1 object3)
  (predicate2 object5)
  (predicate1 object5)
  (predicate2 object7)
  (predicate5 object1 object7)
  (predicate1 object1)
  (predicate2 object2)
  (predicate1 object2)
  (predicate2 object6)
  (predicate1 object6)
  (predicate2 object4)
  (predicate1 object4)
  (predicate3)
;; BEGIN ADD
  (not-allowed object1 object2)
;; END ADD
)
 (:goal (and
  (predicate2 object5)
  (predicate5 object2 object5)
  (predicate2 object4)
  (predicate5 object7 object4)
  (predicate5 object6 object7)
  (predicate5 object1 object6)
  (predicate2 object3)
))
)

```

Listing 15: Prompt for LLM-as-planner

```

You are a PDDL expert. Here is a game we are playing.
{domain_description}
{problem_description}
{constraint_description}

Write the plan that would solve this problem.

These are the available actions:
{available_actions}

Here is what the output should look like:
{example_answer}

```

Listing 16: Prompt for LLM-as-PDDL-formalizer, while generating the entire code

```

You are a PDDL expert. Here is a game we are playing.
{domain_description}
{problem_description}
{constraint_description}

Write the domain and problem files in minimal PDDL.

```

1435
1436
~~1438~~
These are the available actions:
{available_actions}

Listing 17: Prompt for LLM-as-PDDL-formalizer, while first generating the non-constrained PDDL for the Edit setting

1439
1440
1441
1442
1443
1444
1445
~~1446~~
You are a PDDL expert. Here is a game we are playing.
{domain_description}
{problem_description}
Write the domain and problem files in minimal PDDL.

These are the available actions:
{available_actions}

Listing 18: Prompt for LLM-as-PDDL-formalizer, while modifying the non-constrained PDDL for the Edit setting

1448
1449
1450
1451
1452
1453
1454
~~1455~~
You are a PDDL expert. Here is a PDDL domain and problem file.
{original_domain_file}
{original_problem_file}

Modify the PDDL files so that it satisfies the following constraint: {constraint_description}
These are the available actions:
{available_actions}

Listing 19: Prompt for LLM-as-PDDL-formalizer, including revision

1457
1458
1459
1460
1461
1462
~~1463~~
You are a PDDL expert. The following domain and problem files have the error: {error}

{domain_file}
{problem_file}

Revise the PDDL to remove the error.

Listing 20: Prompt for LLM-as-PDDL3-formalizer, while generating the entire code

1465
1466
1467
1468
1469
1470
1471
1472
1473
~~1474~~
You are a PDDL3 expert. Here is a game we are playing.
{domain_description}
{problem_description}
{constraint_description}

Write the domain and problem files in minimal PDDL3 code.

These are the available actions:
{available_actions}

Listing 21: Prompt for LLM-as-PDDL3-formalizer, while first generating the non-constrained PDDL for the Edit setting

1476
1477
1478
1479
1480
1481
1482
~~1483~~
You are a PDDL3 expert. Here is a game we are playing.
{domain_description}
{problem_description}
Write the domain and problem files in minimal PDDL3 code.

These are the available actions:
{available_actions}

Listing 22: Prompt for LLM-as-PDDL3-formalizer, while modifying the non-constrained PDDL for the Edit setting

1485
1486
1487
1488
1489
You are a PDDL3 expert. Here is a PDDL3 domain and problem file.
{original_domain_file}
{original_problem_file}

```
Modify the files using PDDL3 so that it satisfies the following constraint: {constraint_description}
These are the available actions:
{available_actions}
```

1490
1491
1493

Listing 23: Prompts for LLM-as-PDDL3-formalizer, including revision

```
;; Error when running Compiler
You are a PDDL3 expert. The following PDDL3 were compiled and resulted in an error. Revise the PDDL3
to remove the error. Return a JSON object in the following format:
{{
  \"domain file\": ...,
  \"problem file\":...
}}

Domain File:
{original_domain_file}

Problem File:
{original_problem_file}

Error Message:
{compilation_error}

;; Error when running Solver
You are a PDDL3 expert. The following PDDL3 were compiled and ran through a solver, which resulted in
an error. Revise the original PDDL3 to remove the error. Return a JSON object in the following
format:
{{
  \"domain file\": ...,
  \"problem file\":...
}}

Original Domain File:
{original_domain_file}

Original Problem File:
{original_problem_file}

Compiled Domain File:
{compiled_domain_file}

Compiled Problem File:
{compiled_problem_file}

Error Message:
{solver_error}
```

1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534

Listing 24: Prompt for LLM-as-SMT-formalizer (generate)

```
You are a Z3 expert. Here is a domain and problem instance for {domain}.
{domain_description}
{problem_description}
{constraint_description}

Generate Python code that uses the Z3 Python API to solve this instance. These are the available
actions:
{available_actions}

The output of your Python code should be a plan in the following format:{example_plan}

Set the number of steps allowed to 100.
```

1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548

Listing 25: Prompts for LLM-as-SMT-formalizer (edit)

```
;; Step 1
You are a Z3 expert. Here is a domain and problem instance for {domain}.
{domain_description}
```

1550
1551
1552
1553

```

1554 {problem_description}
1555 Generate Python code that uses the Z3 Python API to solve this instance. These are the available
1556 actions:
1557 {available_actions}
1558
1559 The output of your Python code should be a plan in the following format:
1560 {example_plan}
1561
1562 Set the number of steps allowed to 100.
1563
1564 ;; Step 2
1565 You are a Z3 expert. Here is a {domain} instance written in Z3.
1566 {original_python_code}
1567
1568 Modify the Z3 code so that it satisfies the following constraint: {constraint_description}
1569 These are the available actions:
1570 {available_actions}
1571
1572 The output of your Python code should be a plan in the following format:
1573 {example_plan}
1574 Set the number of steps allowed to 100.

```

Listing 26: Prompt for revision step for LLM-as-SMT-formalizer

```

1576 You are a Z3 expert. The following file have the error: {error}
1577
1578 {python_file}
1579
1580 Revise the Z3 code to remove the error.
1581

```

Listing 27: Prompt for LLM-as-LTL-Formalizer

```

1583 ;; coin_collector_formula
1584 You are given a description of a grid world with rooms, adjacencies (north/south/east/west
1585 relations),
1586 and possibly items (e.g. a coin) located in some room.
1587
1588 Your task is to convert this description into a set of LTL formulas for planning.
1589 Decompose the problem into 5 parts:
1590 1. Initial condition
1591 2. Occupancy (exactly one room at a time)
1592 3. Move dynamics (allowed moves)
1593 4. Pickup dynamics (pickup conditions)
1594 5. Goal (eventual condition)
1595
1596 ### Rules:
1597 - Each room should be an atomic proposition in its name (e.g., `kitchen`, `bedroom`, etc.).
1598 - Exactly one room must be true at any time.
1599 - Legal moves: if you are in room A, the next state must be in one of the adjacent rooms or staying
1600 at the current room.
1601 - Legal pickups: if there is a coin in a room R, define an atomic proposition `has_coin`, and add
1602 these two rules to pickup dynamics for `has_coin`:
1603 - If the agent is in R and does not have the coin, then in the next state it may either gain the
1604 coin or stay without it.
1605 - If the agent is not in R and does not have the coin, then in the next state it must still not
1606 have the coin.
1607 - If the agent has the coin, it continues to have it in all future states.
1608 - if the coin is in room R, then picking it up must only change `has_coin` but not coincide with
1609 moving.
1610 - Initial condition: mark the starting room. `has_coin` is always false initially.
1611 - Goal: eventually `has_coin` must be true.
1612
1613 ### Output format:
1614 Produce 4 formulas in Python strings:
1615 - `init = "...`
1616 - `occupancy = "...`
1617 - `move_dynamics = "...`
1618 - `pickup_dynamics = "...`
1619

```

```

- `goal = "...`
1620
---
1621
[DESCRIPTION]
1622
---
1623
Now generate the formulas:
1624
1625
;; coin_collector_adjacency
1626
You are given a description of a world with rooms and directional relations.
1627
1628
Your task is to output the adjacency dictionary in Python format.
1629
1630
### Rules:
1631
- Use room names exactly as given in the text, lowercased and underscores instead of spaces.
1632
- Each key in the dictionary is a room.
1633
- Each value is another dictionary: {neighbor: direction}.
1634
- The direction must be from the perspective of the key room (north, south, east, west).
1635
- Always include both directions if they are implied.
1636
Example: "To the north of the kitchen is the pantry" means:
1637
adjacency["kitchen"]["pantry"] = "north"
1638
adjacency["pantry"]["kitchen"] = "south"
1639
1640
### Output format:
1641
adjacency = {
1642
"room1": {"neighbor1": "direction", "neighbor2": "direction", ...},
1643
"room2": {...},
1644
...
1645
}
1646
1647
---
1648
[DESCRIPTION]
1649
---
1650
Now generate only the adjacency dictionary in the output format:
1651
1652
;; coin_collector_constraints
1653
You are given a set of atomic propositions and an environment description that represent a grid
1654
world planning problem.
1655
- Each room name is an atomic proposition, e.g., [ROOMS]
1656
- The agent may pick up a coin, tracked by proposition "has_coin".
1657
Environment description:
1658
[PROBELM_DESCRIPTION]
1659
1660
Your task: Given an additional natural language constraint, output the corresponding LTL formula
1661
that enforces it.
1662
1663
### Rules:
1664
- Always output formulas in the form of Python strings, e.g., "G(...)" or "F(...)".
1665
- Use the existing atomic propositions exactly as defined.
1666
- If the constraint mentions a room, use the room's name directly.
1667
- You should only output formulas for the new constraints.
1668
1669
### Examples:
1670
- If a constraint says "never enter ROOM", output `G(!ROOM)`
1671
- If a constraint says "only move south", expand it into formulas of the form:
1672
For each room r, if its southern neighbor is s, then:
1673
G( r -> X(s) )
1674
and forbid all other neighbors in the X operator.
1675
- If a constraint says "only move south or west", expand similarly, listing all allowed successors
1676
explicitly.
1677
- If a constraint says "must visit ROOM_A before ROOM_B", encode it as `(!ROOM_B) U ROOM_A`
1678
- If a constraint says "visit ROOM_A right after pick up coin", be careful about the number of X
1679
operators:
1680
If you are picking up the coin in the next timestep, visiting room a will happen in the
1681
following timestep: "G( (!has_coin & X(has_coin)) -> XX(ROOM_A) )"
1682
- Always output Python string formulas, one per constraint.
1683
- Do not introduce undefined macros.
1684
1685
### Formatting Examples:
1686
Natural language: "The agent must never enter the backyard."
1687
Output: "G(!backyard)"
1688
1689

```

1690
1691
1692
1693
1694
1695
1696
1697
1698

Natural language: "The agent must visit the bathroom before visiting the bedroom."
Output: " $(\neg \text{bedroom}) \text{ U bathroom}$ "

Constraint Description:

[CONSTRAINTS_DESCRIPTION]

Now, given these natural language constraints, output only the corresponding LTL formulas, without explanations:

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-planner	N/A	deepseek-reasoner	91	54
		deepseek-chat	69	27
		Qwen3-32B	61	28
		Qwen2.5-Coder-32B-Instruct	3	2
llm-as-pddl-formalizer	generate	deepseek-reasoner	69	40
		deepseek-chat	39	11
		Qwen3-32B	16	16
		Qwen2.5-Coder-32B-Instruct	4	6
	generate + revision	deepseek-reasoner	79	46
		deepseek-chat	71	23
		Qwen3-32B	29	20
		Qwen2.5-Coder-32B-Instruct	5	7
	edit	deepseek-reasoner	75	42
		deepseek-chat	41	22
		Qwen3-32B	27	14
		Qwen2.5-Coder-32B-Instruct	3	5
	edit + revision	deepseek-reasoner	75	56
		deepseek-chat	41	30
		Qwen3-32B	27	24
		Qwen2.5-Coder-32B-Instruct	3	5
llm-as-smt-formalizer	generate	deepseek-reasoner	2	0
		deepseek-chat	1	2
		Qwen3-32B	0	1
		Qwen2.5-Coder-32B-Instruct	0	2
	generate + revision	deepseek-reasoner	28	13
		deepseek-chat	1	6
		Qwen3-32B	2	13
		Qwen2.5-Coder-32B-Instruct	0	4
	edit	deepseek-reasoner	1	1
		deepseek-chat	0	3
		Qwen3-32B	1	2
		Qwen2.5-Coder-32B-Instruct	1	2
	edit + revision	deepseek-reasoner	1	22
		deepseek-chat	0	3
		Qwen3-32B	1	13
		Qwen2.5-Coder-32B-Instruct	1	5

Table 3: Correctness with and without constraints on BlocksWorld-100 using all methods.

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-planner	N/A	deepseek-reasoner	89	49
		deepseek-chat	80	38
		Qwen3-32B	96	52
		Qwen2.5-Coder-32B-Instruct	79	36
llm-as-pddl-formalizer	generate	deepseek-reasoner	65	26
		deepseek-chat	64	24
		Qwen3-32B	45	20
		Qwen2.5-Coder-32B-Instruct	3	3
	generate + revision	deepseek-reasoner	75	34
		deepseek-chat	65	25
		Qwen3-32B	48	21
		Qwen2.5-Coder-32B-Instruct	5	4
	edit	deepseek-reasoner	62	39
		deepseek-chat	64	21
		Qwen3-32B	50	18
		Qwen2.5-Coder-32B-Instruct	4	2
edit + revision	deepseek-reasoner	62	41	
	deepseek-chat	64	21	
	Qwen3-32B	50	24	
	Qwen2.5-Coder-32B-Instruct	4	2	
llm-as-smt-formalizer	generate	deepseek-reasoner	45	23
		deepseek-chat	42	23
		Qwen3-32B	29	18
		Qwen2.5-Coder-32B-Instruct	1	7
	generate + revision	deepseek-reasoner	77	39
		deepseek-chat	58	32
		Qwen3-32B	49	32
		Qwen2.5-Coder-32B-Instruct	1	11
	edit	deepseek-reasoner	43	29
		deepseek-chat	39	24
		Qwen3-32B	28	17
		Qwen2.5-Coder-32B-Instruct	1	3
edit + revision	deepseek-reasoner	43	47	
	deepseek-chat	39	33	
	Qwen3-32B	28	24	
	Qwen2.5-Coder-32B-Instruct	1	7	

Table 4: Correctness with and without constraints on CoinCollector-100 using all LLM-as-Planner, LLM-as-PDDL-Formalizer, LLM-as-SMT-Formalizer.

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-ltl-formalizer	generate	deepseek-reasoner	71	27
		deepseek-chat	50	18
		Qwen3-32B	0	15
		Qwen2.5-Coder-32B-Instruct	7	12

Table 5: Correctness with and without constraints on CoinCollector-100 using LLM-as-LTL-Formalizer.

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-pddl3-formalizer	generate	deepseek-reasoner	61	29
		deepseek-chat	64	9
		Qwen3-32B	35	12
		Qwen2.5-Coder-32B-Instruct	5	2
	generate + revision	deepseek-reasoner	63	33
		deepseek-chat	76	13
		Qwen3-32B	43	18
		Qwen2.5-Coder-32B-Instruct	12	5
	edit	deepseek-reasoner	61	32
		deepseek-chat	57	24
		Qwen3-32B	33	17
		Qwen2.5-Coder-32B-Instruct	5	6
	edit + revision	deepseek-reasoner	61	39
		deepseek-chat	57	32
		Qwen3-32B	33	24
		Qwen2.5-Coder-32B-Instruct	5	7

Table 6: Correctness with and without constraints on BlocksWorld-100 using PDDL3.

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-pddl3-formalizer	generate	deepseek-reasoner	54	21
		deepseek-chat	13	8
		Qwen3-32B	24	8
		Qwen2.5-Coder-32B-Instruct	2	0
	generate + revision	deepseek-reasoner	55	30
		deepseek-chat	49	11
		Qwen3-32B	58	26
		Qwen2.5-Coder-32B-Instruct	21	8
	edit	deepseek-reasoner	60	30
		deepseek-chat	22	8
		Qwen3-32B	21	7
		Qwen2.5-Coder-32B-Instruct	0	0
	edit + revision	deepseek-reasoner	60	32
		deepseek-chat	22	17
		Qwen3-32B	21	25
		Qwen2.5-Coder-32B-Instruct	0	8

Table 7: Correctness with and without constraints on CoinCollector-100 using PDDL3.

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-planner	N/A	deepseek-reasoner	43	30
		deepseek-chat	10	8
		Qwen3-32B	4	3
		Qwen2.5-Coder-32B-Instruct	1	1
llm-as-pddl-formalizer	generate	deepseek-reasoner	47	13
		deepseek-chat	21	5
		Qwen3-32B	11	12
		Qwen2.5-Coder-32B-Instruct	0	4
	generate + revision	deepseek-reasoner	49	18
		deepseek-chat	54	13
		Qwen3-32B	22	14
		Qwen2.5-Coder-32B-Instruct	0	7
	edit	deepseek-reasoner	40	18
		deepseek-chat	27	12
		Qwen3-32B	7	5
		Qwen2.5-Coder-32B-Instruct	1	4
edit + revision	deepseek-reasoner	40	25	
	deepseek-chat	27	22	
	Qwen3-32B	7	15	
	Qwen2.5-Coder-32B-Instruct	1	5	

Table 8: Correctness with and without constraints on BlocksWorld-XL-100 using PDDL.

formalizer type	method	model	correctness (non-constrained)	correctness (constrained)
llm-as-planner	N/A	deepseek-reasoner	23	15
		deepseek-chat	11	6
		Qwen3-32B	10	4
		Qwen2.5-Coder-32B-Instruct	1	1
llm-as-pddl-formalizer	generate	deepseek-reasoner	76	19
		deepseek-chat	61	15
		Qwen3-32B	82	18
		Qwen2.5-Coder-32B-Instruct	86	18
	generate + revision	deepseek-reasoner	95	24
		deepseek-chat	91	23
		Qwen3-32B	93	26
		Qwen2.5-Coder-32B-Instruct	86	18
	edit	deepseek-reasoner	80	22
		deepseek-chat	64	15
		Qwen3-32B	81	26
		Qwen2.5-Coder-32B-Instruct	92	17
edit + revision	deepseek-reasoner	80	34	
	deepseek-chat	64	23	
	Qwen3-32B	81	27	
	Qwen2.5-Coder-32B-Instruct	92	17	

Table 9: Correctness with and without constraints on MysteryBlocksWorld-100 using all PDDL.