

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 KNOWLEDGE MODEL PROMPTING INCREASES LLM PERFORMANCE ON PLANNING TASKS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Large Language Models (LLM) can struggle with reasoning ability and procedural tasks. Many prompting techniques have been developed to assist with LLM reasoning, notably Chain-of-Thought (CoT); however, these techniques, too, have come under scrutiny as LLMs' ability to reason at all has come into question. Borrowing from the domain of cognitive and educational science, this paper investigates whether the Task-Method-Knowledge (TMK) framework can improve LLM reasoning capabilities beyond its previously demonstrated success in educational applications. The TMK framework's unique ability to capture causal, teleological, and hierarchical reasoning structures, combined with its explicit task decomposition mechanisms, makes it particularly well-suited for addressing language model reasoning deficiencies, and unlike other hierarchical frameworks such as HTN and BDI, TMK provides explicit representations of not just what to do and how to do it, but also why actions are taken. The study evaluates TMK by experimenting on the PlanBench benchmark, focusing on the Blocksworld domain to test for reasoning and planning capabilities, examining whether TMK-structured prompting can help language models better decompose complex planning problems into manageable sub-tasks. Results also highlight significant performance inversion in reasoning models. TMK prompting enables the reasoning model to achieve up to an accuracy of 97.3% on opaque, symbolic tasks (Random versions of Blocksworld in PlanBench) where it previously failed (31.5%), suggesting the potential to bridge the gap between semantic approximation and symbolic manipulation. Our findings suggest that TMK functions not merely as context, but also as a mechanism that steers reasoning models away from their default linguistic modes to engage formal, code-execution pathways in the context of the experiments.

## 1 INTRODUCTION

Large Language Models (LLMs) have been criticized by recent research for their poor reasoning and planning abilities (Valmeekam et al., 2023b; Chan, 2024; Shojaee et al., 2025). Prompt engineering has been a source of study for attempts at improving these abilities, with mixed results (Stechly et al., 2024; Bhambri et al., 2025). Evaluations of reasoning, through notable methods that allow for problem decomposition such as Chain-of-Thought (CoT) (Wei et al., 2022), have shown modest improvement in some planning domains. Other research has called into question the validity of these findings and argues that LLMs are incapable of the in-context procedural reasoning from the n-shot examples that the CoT authors claim (Stechly et al., 2024; Bhambri et al., 2025). Critics argue that LLMs are incapable of planning but instead rely on pattern matching with similar prompts for performance increases.

We investigate whether prompting using a TMK framework can improve LLM performance on planning tasks by borrowing methods from the domain of cognitive and education science (Sushri et al., 2024; Dass et al., 2025; Lum et al., 2025), while still accounting for criticisms given to previous research. In the introduction and related works sections of CoT and ReACT papers (Wei et al., 2022; Yao et al., 2022), the cognitive thought process was similarly referred to as motivation.

We introduce TMK (Task, Method, Knowledge) (Murdock & Goel, 2008), a knowledge-based self-explanation framework used successfully to explain procedural learning to students (Sushri et al., 2024), into a language model prompt to determine if TMK can improve a language model's proce-

dural planning ability. The paper seeks to extend on previous work from other fields by focusing on TMK language, which was developed as part of research into cognitive architectures in intelligent agent systems to allow systems to reason about their own processing and make any necessary adjustments (Murdock, 2001; Murdock & Goel, 2008). TMK is a representationally agnostic model authored by domain experts; JSON is one serialization used when writing the model into the LLM prompt (Dass et al., 2025). Successful use in the education domain to supplement LLM in teaching procedural learning to students, offered motivation to explore whether the TMK framework could also qualitatively improve language models’ performance. Beyond raw scores, we investigate TMK’s potential as an additional inference steering mechanism. A structured prompt that can steer a model’s reliance on probabilistic linguistic patterns and activate its latent symbolic processing capabilities.

The hypothesis driving the experiments is assessing only raw scores. We hypothesize that a TMK-structured prompt acts as a **symbolic steering mechanism**, enabling language models to decouple logical reasoning from semantic interference. We test this by evaluating whether TMK improves planning performance specifically in the PlanBench Blocksworld domain, where semantic cues are absent (Random Blocksworld) or misleading (Mystery Blocksworld). This is to be supported by observing bias shifts from linguistic approximation to symbolic execution with TMK.

To test this hypothesis, we will use OpenAI models available at time of writing against the PlanBench benchmark (Valmeekam et al., 2023b). The benchmark evaluates language models on a set of planning problems, which maintains a publicly available leader board (Valmeekam, 2023) of Blocksworld problems and their variations that we will use as a comparison for our hypothesis.

Our findings indicate that TMK inclusion does increase the success of the language model in planning tasks across the flagship models. We observe improvements hold for flagship models but not models that are specifically optimized for domains or expediency. Optimization techniques such as distillation, quantization, pruning and other techniques as shown in Zhu et al. (2024) often comes at a cost of performance and can cause catastrophic forgetting (McCloskey & Cohen, 1989) which may impact planning, we discuss this further in section 5.

Across models and the various Blocksworld datasets, the TMK framework generally increased performance, with significant gains across reasoning models for random Blocksworld problems and a maximum gain of 65.8% on the o1 model for Random Blocksworld (rising from 31.5% to 97.3%). We believe the performance gains are an indication that the TMK framework is likely to demonstrate similar gains in planning tasks for other domains suggested in section 6.

## 2 RELATED WORK

The paper’s research is at the intersection of two distinct areas: standardized planning benchmarks for language models and TMK prompting strategies for planning. While prior work has investigated these areas in isolation, this paper explores how integrating TMK into prompts can improve planning tasks and surpass state-of-the-art (SoTA) performance identified in recent literature for flagship models in publicly available planning benchmarks (Valmeekam, 2023).

### 2.1 PROMPTING FOR PLANNING

Prompt engineering for supporting a language model’s response is a well-researched area, with CoT and ReACT, being the most notable with regard to problem decomposition (Wei et al., 2022; Yao et al., 2022). A literature survey by Sahoo et al. (2024) found over 40 prompting techniques. However, due to the wide breadth of applications, few have claimed improvement in planning. Even fewer define planning formally, where every step leading to the final answer must be correct (as opposed to approximations that sound feasible but are incorrect).

**Chain-of-Thought and Its Limitations.** While CoT showed initial promise on reasoning tasks, recent systematic evaluations reveal severe limitations for classical planning. Stechly et al. (2024) tested CoT on 261 simpler Blocksworld problems, ”Table-To-Stack” configurations where every block starts on the table, without pre-defined stacks, with the goal of one single stack. Stechly et al. (2024) found that zero-shot CoT on average achieved only 1.1% improvement when compared to no CoT (baselined across GPT-4, GPT-4-Turbo, and Claude-3-Opus).

108 **Chain-of-Symbols (CoS)’s Informativity.** CoS (Hu et al., 2023) another in-literature candidate, uses  
 109 an extension of CoT using symbolic representations with 5 demonstrated (five-shot) examples. Brick  
 110 World, despite the name similarity, is a distinct, simpler domain from classical Blocksworld. It uses  
 111 Longest Common Sequence(LCS) to evaluate intermediate steps recorded as precision (LCS divided  
 112 by length of output) and recall (LCS divided by length of ground truth) percentages, and reports the  
 113 final answer match as accuracy. It lacks the formal PDDL specification and validation mechanisms  
 114 required for rigorous planning evaluation.

115 **ReACT’s Brittleness in Planning.** ReACT, which interleaves reasoning traces with action execu-  
 116 tion, faces similar challenges. A critical evaluation by Bhambri et al. (2025) found that ReACT’s  
 117 performance on AlfWorld and WebShop domains is minimally influenced by the reasoning trace  
 118 content. Instead, success stems from ”unreasonably high similarity between input example tasks  
 119 and queries,” requiring instance-specific examples that significantly increase human cognitive bur-  
 120 den. When examples differ even slightly from queries (e.g., different object names, goal locations,  
 121 or task types within the same domain), performance collapses.

122 In our paper, we deal with such criticism by following the PlanBench benchmark, which has the fol-  
 123 lowing requirement: (1) The leader board was derived from zero-shot and one-shot examples, which  
 124 discounts highly identical example similarity as a factor to planning ability. (2) A correct answer is  
 125 not just the final state, but also the entirety of the given reasoning process. These two requirements  
 126 address criticisms given in this section. For this paper, we will only focus on Blocksworld and  
 127 OpenAI models; however, we believe that our results indicate a strong direction for future research.

## 128 2.2 PLANBENCH

131 Existing literature shows that language models struggle with classical planning tasks (Valmeekam  
 132 et al., 2023b; Chan, 2024; Shojaee et al., 2025). Planning requires every action to be formally valid  
 133 (“close enough” approximations fail verification) and is fundamentally distinct from the question-  
 134 answering and commonsense reasoning tasks in ReACT and CoT (Wei et al., 2022; Yao et al., 2022),  
 135 where next-token prediction patterns can often derive correct final answers through approximating  
 136 intermediate as described in section 2.1. This paper utilizes PlanBench (Valmeekam et al., 2023a),  
 137 an extensible benchmark based on International Planning Competition (IPC) domains that verifies  
 138 language model outputs using classical planning tools (Helmert, 2006; Howey et al., 2004) without  
 139 providing the models access to those tools at inference time. Rather, it uses those tools to ro-  
 140 bustly validate plans generated by language models. Unlike commonsense and arithmetic reasoning  
 141 benchmarks, PlanBench tests the planning ability of a language model and requires formal symbolic  
 142 correctness validated by automated planners and plan validators.

142 Crucially, PlanBench introduces domain obfuscation to decouple a language model’s reasoning abil-  
 143 ity from its pre-trained linguistic priors. As shown in Table 1, the benchmark includes three variants:  
 144

- 145 • **Classic:** Uses canonical English labels (e.g., pick up, stack), allowing models to rely on  
 146 semantic memory.
- 147 • **Mystery:** Maps actions to semantically distinct but unrelated words (e.g., attack, feast).  
 148 This tests if the model can reason using the provided rules rather than semantic associations.
- 149 • **Random:** Replaces labels with opaque alphanumeric strings. This steers the model to rely  
 150 more on symbolic manipulation.

152 For the scope of this paper, we focus on Blocksworld across these three to give some initial insight  
 153 into how TMK impacts planning in language models. This selection provides a baseline for language  
 154 models without TMK, allowing us to isolate whether the TMK framework improves planning tasks  
 155 in Blocksworld and its more challenging variations.

## 156 2.3 THE TMK FRAMEWORK

159 The TMK (Task, Method, Knowledge) framework is a formal knowledge representation language  
 160 originally developed to model the teleology of intelligent agents (Murdock, 2001). It hierarchi-  
 161 cally decomposes a system into three components: (1) **Tasks** (goals and conditions), (2) **Methods**  
 162 (procedural mechanisms), and (3) **Knowledge** (domain concepts and relationships). This explicit

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

Figure 1: Expansion of an example Task, Method and Knowledge structure, where goals are teleologically connected to methods by which they are achieved using their knowledge of the environment or domain (Blocksworld in this example). See OSF link for detailed expansion: (Anonymous, 2025)

210

211

212

213

214

215



216 decomposition enables the capture of not only *what* to do and *how* to do it, but crucially, *why* actions  
 217 are taken (teleology), which we suspect can help ground a model’s reasoning and planning.  
 218

219 While other hierarchical task decomposition frameworks exist, such as Belief-Desire-Intention  
 220 (BDI) (Georgeff et al., 1998) and Hierarchical Task Network (HTN) (Erol et al., 1994), TMK is  
 221 distinct in its emphasis on teleological and causal self-explanation. Like TMK, these frameworks  
 222 are candidates for symbolic interpreters and may offer future potential for dual use within both  
 223 prompts and systems with symbolic solvers. However, TMK’s specific architectural focus on the  
 224 “why” of actions makes it uniquely applicable to our study.  
 225

226 The framework provides a structured basis for reasoning by breaking available actions into simpler,  
 227 manageable components, a process that mirrors cognitive decomposition (Correa et al., 2023). Con-  
 228sequently, TMK’s unique ability to capture causal and teleological structures makes it well-suited  
 229 for addressing the planning deficiencies often observed in language models. It achieves this by ex-  
 230plicitly defining clear steps (methods), the relationships between them (knowledge), and the goals  
 231 driving them (tasks).  
 232

233 Prior research has demonstrated that language model outputs often favor declarative knowledge over  
 234 procedural steps (Dass et al., 2025). This preference contributes to the poor planning performance  
 235 described in section 2.2, evidence in Valmeekam (2023). Because TMKs are authored by domain  
 236 experts specifically to facilitate procedural learning through causal reasoning, we hypothesize that  
 237 prompting with TMK can equip a language model with the procedural and causally grounded context  
 238 preferred by experts (Lum et al., 2025). Validating this hypothesis is the primary investigation of this  
 239 paper. Furthermore, task decomposition was essential for reasoning in pre-LLM systems (Murdock,  
 240 2001; Georgeff et al., 1998; Erol et al., 1994). Given the documented struggles of modern LLMs  
 241 with autonomous reasoning (Valmeekam et al., 2023a; Chan, 2024; Shojaee et al., 2025), applying  
 242 the TMK framework to prompt engineering offers a promising avenue for enhancing LLM planning  
 243 capabilities.  
 244

### 245 3 METHODS

#### 246 3.1 DESIGNING A MODEL WITH THE TMK FRAMEWORK

247 The TMK knowledge representation framework has three components: (1) Task describes the goal  
 248 of a system, (2) Method describes mechanisms by which the goal is achieved, and (3) Knowledge  
 249 defines the domain ontology required to interpret tasks, mechanisms and support them (Murdock,  
 250 2001; Murdock & Goel, 2008; Goel & Rugaber, 2017). For the purpose of this paper, the authors  
 251 converted the domain knowledge from plain text to TMK by breaking down each possible step  
 252 within Blocksworld into methods. The authors kept to a three-layer (as shown in Fig. 1) hierarchy  
 253 decomposition to maintain the conciseness of the TMK structure to keep within effective context  
 254 windows (Liu et al., 2024); however, there is no limit to the number of layers for more complex  
 255 procedural goals to be decomposed or the level of intricacy to be included as problem complex-  
 256 ity and context windows grows. For instance, there can be a plan verification (i.e. verifyPlan and  
 257 CheckGoalAchieved goals) that themselves can be further broken down much deeper down to actu-  
 258 ally verifying the sequence of correct steps using known algorithms, similar to verifiers and solvers  
 259 such as VAL (Howey et al., 2004) and fast downward (Helmert, 2006) respectively. TMK enables  
 260 calibration of the level of abstraction depending on the use case.  
 261

##### 262 3.1.1 TASK, GIVEN BY ITS GOALS (THE WHY)

263 Within each goal, fields such as *name* and *description* provide human-oriented identifiers and sum-  
 264maries, while input parameters and output parameters make data flow explicit. There are also the  
 265 *given* and *makes* clauses to encode pre- and post-conditions of the goals, such as “holding” or “is  
 266 clear”. Finally, there is also a *Mechanism* field that links the goal to a method, connecting “why”  
 267 (goal) and “how” (mechanism). These fields in the TMK framework help guide designers and pro-  
 268 vide a structure as they build out their TMK.  
 269

270 3.1.2 METHOD, GIVEN BY ITS MECHANISMS (THE HOW)  
271272 Mechanisms correspond to goals, and they include similar *description*, *input*, and *output* fields but  
273 differ in having *require* and *provides* clauses to differentiate logical requirements of concepts and  
274 relations. The *require* and *provides* clauses alongside an operation within the process field document  
275 the logical causal updates, allowing the mechanism's post-condition to satisfy a parent goal, which  
276 in our case is picking up a block.277 3.1.3 KNOWLEDGE, DEFINES THE DOMAIN-SPECIFIC CONCEPTS AND RELATIONSHIPS,  
278 GOALS, AND MECHANISMS USED  
279280 Knowledge defines the ontology that parameter names and logical conditions reference in the TMK,  
281 ensuring consistent interpretation across goals and mechanisms. This combination is used in the  
282 experiment to define the domain knowledge of a procedural knowledge and reasoning question,  
283 giving the AI models a description of the decomposed task as input.284 3.1.4 TMK DESIGN, PLACEMENT AND IMPLEMENTATION  
285286 Our TMK prompt replaces the domain portion of the PlanBench prompt shown in Valmeekam et al.  
287 (2023b) with a TMK formatted in JSON, using the one-shot example and querying it for consistency.  
288 This can also be compared in the prompt examples in appendix A.289 It is important to note that the process of designing and building a TMK is iterative and can be  
290 decomposed deeper or abstracted further, as evident in Murdock (2001). In this paper, it is kept  
291 to three layers for ease of discussion and experimentation. The level of abstraction is left to the  
292 designer, who can also utilize the framework to structure their thinking.293 There are also differences in TMK in classic blocksworld, mystery blocksworld, and random  
294 blocksworld (see in OSF link, Anonymous (2025)). All prompts and system prompts are also in-  
295 cluded in the results file. Efforts are made to keep it as similar as possible, given constraints required  
296 by the different types of Blocksworld datasets. This difference is also consistent in the PlanBench  
297 prompts, where 'random' and 'mystery' Blocksworld were using true and false descriptions, but not  
298 the classical Blocksworld prompt, as Valmeekam et al. (2023a)'s prompt evolved.300 3.2 EVALUATION ON PLANBENCH  
301302 As part of this paper, we evaluated the TMK against the public PlanBench benchmarks found on  
303 Valmeekam (2023). There are a few key differences between our method and the one described as  
304 part of the benchmark. One of which is that our testing is one-shot while the Valmeekam (2023) is  
305 zero-shot. This difference is inconsequential to our findings for the following reasons:306  
307 1. **Baseline Rigor and Formatting:** The first is that the original PlanBench paper  
308 (Valmeekam et al., 2023a) consist entirely of one-shot questions while the publicly avail-  
309 able leader board, Valmeekam (2023) is of zero-shot. Given this mix, we made the decision  
310 to compare against the higher value (zero-shot) for rigor. For TMK, one-shot prompts also  
311 allowed us to conform TMK outputs exactly to what Valmeekam (2023) expects while  
312 reducing formatting instructions, minimizing non-TMK differences in prompts.  
313  
314 2. **Precedence and Performance:** The papers listed on Valmeekam (2023) reported only the  
315 zero-shot case, we suspect it is because as observed in the data, zero-shot did better than  
316 the one-shot for those listed papers (for plain text) (Valmeekam, 2023; Valmeekam et al.,  
317 2023c). Literature (Kojima et al., 2022) also confirms that LLMs are zero shot reasoners  
318 and this is additionally backed by our sample testing of one-shot examples in plain text (see  
319 results and code in Anonymous (2025), OSF link).  
320  
321 3. **Nature of the Example:** Lastly, although we offer a one-shot example, that example is  
322 random and not tailored to the problem at hand as is often the case with similar research.323 The Valmeekam (2023) code at the time of writing required update in the extracting random  
324 blocksworld to be comparable with the ground truth. Our experiment thus also added new code (in-  
325 cluded in Anonymous (2025)) to the extraction criteria which was applied for random blocksworld

324 data set. It focuses on the sequence of execution and order of objects after the actions, which are the  
 325 key set of criteria to determine the ability of the LLM to plan.  
 326

327 The authors also made the choice to not get perturbed by the stochastic language model responses  
 328 that can sometimes include: (1) symbols (like "-", "\_" ect.), (2) words (like "o", "obj") instead of  
 329 "object"(blocks) and (3) plan steps like "stack object b from object a", (translated to "Overcome  
 330 object b from object a" and "2ijg9q8swj2shjel stack object b from object a" for mystery and random  
 331 domains respectively). Our enhanced extraction function prevents these from getting evaluated as  
 332 incorrect because even if words like "object", "from" and symbols like "-", "\_", show up in the  
 333 response of a model, their sequence and ground truth action along with their order match.  
 334

335 This is rare in classic blocksworld, but seems to be an artifact evident within random blocksworld  
 336 domains since language models often have a bias towards generating grammatically correct, human  
 337 understandable, english sentences (Hu et al., 2024). The authors believe these stochastic errors do  
 338 not take away from the ability to assess if language models can plan and reason, which is consistent  
 339 within ICAPS (International Conference on Automated Planning and Scheduling) language refer-  
 340 ence documents (Fern, 2008) and importantly, similar to extraction examples listed in papers that  
 341 investigated PlanBench further (Valmeekam et al., 2023c).  
 342

343 We further confirm these findings by running the PlanBench benchmark for newer models for which  
 344 it has not been reported. In these cases, similar to the paper by the other authors, the one-shot case  
 345 was worse than the zero-shot case (shown in results file in OSF, Anonymous (2025)).  
 346

## 347 4 EXPERIMENT

348 The results of the experiment is consolidated in Table 2. Interestingly, while TMK improvements are  
 349 evident there is also evidence that under the methods of this experiment, more robust benchmarks  
 350 for planning are needed for LRM (large reasoning models) such as o1 and GPT-5, as they perform  
 351 particularly well on the blocksworld domain with TMK prompt significantly outperforming SoTA  
 352 in the random blocksworld domain.  
 353

### 354 4.1 TMK OUTPERFORMED PLAIN TEXT

355 At a high level, the experiments show that all flagship models with domain prompts replaced by  
 356 TMK structures perform better. The authors posit the that optimization processes (i.e. quantization,  
 357 pruning, distillation) that resulted in o1-mini's vastly reduced inference time, cost, and training  
 358 with a mathematical skew (OpenAI, 2025) may have caused it to be an outlier (further discussed in  
 359 section 5). Though without access to information beyond what was released by OpenAI, it remains  
 360 difficult to ascertain why. The majority of the models have also shown a marked improvement (in  
 361 bold) in different blocksworld sub-domains when prompted using the TMK framework with gains  
 362 up to 65.8% in o1 on the random dataset.  
 363

### 364 4.2 COMPARING BETWEEN MYSTERY AND RANDOM BLOCKSWORLD

365 The experiments reveal insight on LRM that justifies further investigation. Notably, we observe a  
 366 strong *performance inversion* in the reasoning models.  
 367

368 In the plain text baseline, o1 follows the typical pattern of language models, performing significantly  
 369 better on Mystery (74.3%) than on Random (31.5%). This suggests that without structure, the  
 370 model relies on semantic cues (even misleading ones) to guide its reasoning. However, under TMK  
 371 prompting, this relationship flips: o1 scores on Random (97.33%), surpassing its Mystery (83.3%).  
 372

373 This suggests that the TMK framework acts as a symbolic scaffold. By imposing a rigid, code-like  
 374 structure, TMK appears to shift the model's inference strategy away from linguistic approximation  
 375 and toward formal symbolic manipulation, a mode that operates most effectively when the tokens  
 376 (e.g., "1jpkithdyjmlikck", aka "pick up") are devoid of pre-trained semantic associations.  
 377

378 Conversely, the smaller o1-mini model does not show this inversion benefit in the Mystery domain  
 379 (dropping to 16.83%). We hypothesize that unlike flagship reasoning models, o1-mini lacks the  
 380 capacity to resolve the conflict between the rigid TMK structure and the semantic interference of the  
 381 Mystery vocabulary, resulting in the degradation observed in Table 2.  
 382

Table 1: Blocksworld predicate and action name correspondences across PlanBench domain variants: Classic (canonical names), Mystery (semantically meaningful obfuscations), and Random (opaque identifiers) used to evaluate planning under reduced linguistic cues while preserving a one to one mapping of planning semantics (Valmeekam et al., 2023a).

Classic	Mystery	Random
Predicates		
Actions		
empty hand	Harmony	3covmuy4yrjthijd
holding	Pain	gk5asm3f7u1fekpj
on table	Planet	51nbwlachmfartjn
on	Object Craves	4dmf1cmtyxgsp94g
clear	Province	aqcjuuehivl8auwt

Table 2: Accuracy (percentage of fully correct plans with complete, stepwise reasoning) on PlanBench Blocksworld variants, Classic, Mystery, and Random comparing plain-text prompts (best of sampled Zero & One shot) and TMK structured prompts (One shot); tasks and evaluation follow PlanBench definitions (Valmeekam et al., 2023a). Bold values indicate significantly improvements (\*Note: o1Preview has been depreciated and replaced by o1. Results extracted from Valmeekam (2023))

Model	Type	Plain Text (%)	TMK (%)
<b>LLM (Large Language Models)</b>			
GPT4	Classic	34.6	39.7
	Mystery	0	3.8
	Random	0	4.17
GPT4o	Classic	35.5	<b>45.3</b>
	Mystery	0	<b>5.5</b>
	Random	0.83	4.83
<b>LRM (Large Reasoning Models)</b>			
o1mini	Classic	56.7	57
	Mystery	19.1	16.83
	Random	9.33	<b>27.0</b>
o1preview*	Classic	97.8	NA
	Mystery	52.8	NA
	Random	37.3	NA
o1	Classic	95.7	98.5
	Mystery	74.3	<b>83.3</b>
	Random	31.5	<b>97.33</b>
GPT5	Classic	99.3	99.7
	Mystery	98.1	98.3
	Random	92.5	<b>99.0</b>

## 5 DISCUSSION

The experimental results in Table 2 demonstrate that integration of a TMK framework into an LLM prompt improves procedural capabilities in a language model. The improvements are not merely incremental in some cases as in o1, the performance on random blocksworld dataset can transform a competent (31.5%) into a SoTA result (97.3%) and for o1-mini, near failure (9.33%) into a competent result (27%). This discussion explores the underlying reasons. The paper argues that the TMK structure acts as a "cognitive scaffold" that decomposes tasks on behalf of language models. The authors posit that the TMK structure can help shift a model's inference processes towards a more formal symbolic mode of operations, more akin to its code generation training data set than plaintext data set. The existence of a code or textual skew in latent space has only recently been reported by Chen et al. (2024), also referencing PlanBench.

## 5.1 CAN LLMs REASON BETTER GIVEN PROMPTING?

Some of the criticisms (Stechly et al., 2024; Bhambri et al., 2025) of language model reasoning include: (1) Papers claiming that prompting increasing reasoning ability in LLMs often use N-shot examples from which the LLM engages in pattern matching, (2) In cases where LLM show CoT as part of their reasoning, the CoT often contradicts the final answer given, and (3) Even in Zero-Shot prompting cases, LLMs do not show significant ability to plan across different domains or models.

Our research shows gains in LLM planning ability while, for the most part, avoiding the criticism above, through taking the following precautions where applicable, respectively: (1) We tested and demonstrated learning gains in only the one-shot case. In the one-shot case, we provided one constituent example, but unlike many of the papers criticized, our example does not match the problem in problem length or block description. Additionally, we demonstrate that in the PlanBench dataset, the zero-shot case often outperforms the one-shot case for plan text, reinforcing that it is the TMK, not the given example, that is creating the gains in accuracy. (2) We evaluate the entire explana-

432  
 433  
 434  
 435  
 436  
 437  
 438  
 439  
 440  
 441  
 442  
 443  
 444  
 445  
 446  
 447  
 448  
 449  
 450  
 451  
 452  
 453  
 454  
 455  
 456  
 457  
 458  
 459  
 460  
 461  
 462  
 463  
 464  
 465  
 466  
 467  
 468  
 469  
 470  
 471  
 472  
 473  
 474  
 475  
 476  
 477  
 478  
 479  
 480  
 481  
 482  
 483  
 484  
 485  
 tion as part of our results, meaning every step of the planning task must be correct in order to be considered. (3) We show for the blocksworld domain, TMK prompting increases the correctness of LLM reasoning across models and subdomains, although only incrementally. Further research is needed to know if TMK prompting is useful across domains. This shows that in the TMK language of describing tasks, methods, and knowledge, along with their functions, descriptions, input, output, pre- and post-conditions that can be used by LLMs to improve their course of action.

## 5.2 EXPLANATIONS FOR ACCURACY GAINS

The authors of this paper have a hypothesis in two parts about why TMK assists in LLM planning, the first relating to code in the training data, and the second as tied to cognitive and educational science

### 5.2.1 STEERING BETWEEN CODE EXECUTION AND TEXTUAL REASONING

Language models are often trained with both code and plain text data (Langlais et al., 2025; Aryabumi et al., 2024). With the TMK framework with nested parenthesis, keywords, explicit variable assignment is structurally analogous to more formal syntax of programming languages. This structural similarity likely prompts the model to shift its inference strategy from a linguistic mode to a more code-like symbolic mode of internal token manipulation.

It is feasible that the superior performance of TMK structure prompts can be attributed to their ability to activate formal reasoning pathways inherent in a model’s code training data. Unlike plain text that a language model was trained on, code training data will naturally contain variables that, while may not look exactly like “1jpkithdyjmlikck” (pick up), can allow a language model to consider and use “1jpkithdyjmlikck” as an arbitrary identifier. This identifier can substitute the word “pick up” and be used in various logical or non-linguistically sequential context that allows attention tokens within transformer architectures (Vaswani et al., 2017) to look further than if it were linguistic (English) sentences. In a more code analogous structure, attention must track long-range dependencies such as linking a variable’s definition to its uses many lines later. The hypothesis is reinforced by the generally greater improvements in LRM scores compared to LLMs from TMK structured prompts. It shows that LRM with longer test time inference can better capitalize on the TMK hierarchical and teleological structures. Compared to plain text, it is reasonable to deduce that the inference (reasoning) tokens of TMK should contain more code-like structures, given that TMK is in JSON format. This should be tested in models that have transparent reasoning tokens as part of future work. Overall, a TMK-formatted prompt appears to be effective at unlocking latent capabilities of LRM for procedural tasks.

This aligns with recent findings by Chen et al. (2024), who demonstrate that textual reasoning has inherent limitations for logic and optimization tasks, often necessitating a shift to code-based execution for accuracy. The TMK framework effectively acts as a steering mechanism, forcing the model to bypass its default textual reasoning pathways, which are susceptible to the semantic interference observed in the Mystery domain—and engage the code-execution reasoning pathways that Chen et al. (2024) identify as superior for symbolic tasks.

The ‘Performance Inversion’ observed in the o1 model—where symbolic tasks (Random) become significantly easier than semantic ones (Mystery) under TMK, serves as empirical validation of this steering effect. If TMK were simply providing additional context, we would expect uniform gains across domains. Instead, the reversal of domain difficulty indicates a fundamental shift in the underlying reasoning modality: the prompt successfully steers the model out of the semantic interference zone and into a robust symbolic manipulation zone.”

### 5.2.2 COGNITIVE SCAFFOLDING

Within the literature review of cognitive and educational science, TMK encourages LLMs to return procedural explanations when compared to unstructured output (Dass et al., 2025). These more procedurally focused explanations are generally preferred by domain experts, but novice learners tend to prefer the more knowledge-focused conversational style of unstructured text (Lum et al., 2025). Under Bloom’s taxonomy of learning, procedural understanding is considered a higher level of learning compared to factual knowledge-focused answers; however, this type of learning is considered more

486 cognitively demanding under cognitive load theory (Krathwohl, 2002). Given this, the preference  
 487 for unstructured knowledge-focused responses in LLMs may be a side effect of the worked example  
 488 effect (Kirschner et al., 2010). The worked example effect is a well-known effect in education  
 489 science where, for novice learners, having worked examples, that is, examples where much of the  
 490 problem is solved, minus key steps that are being taught, has been shown to be immensely effective  
 491 for novice learners. Despite this, novices have been shown to dislike using worked examples in  
 492 their learning. For this reason, TMK prompting may encourage LLMs to produce more procedural  
 493 reasoning explanations, which are more likely to be made by experts, increasing their planning abilities.  
 494 From the perspective of LLM reasoning critics, while TMK prompting does not offer n-shot  
 495 solutions to which the LLM then pattern matches, it instead offers expert knowledge to which the  
 496 LLM pattern matches. That is, prompting with information structured in ways preferred by experts  
 497 in procedural tasks allows the LLM to perform better on those tasks.  
 498

### 499 5.3 LIMITATIONS

500 Our research only investigates improving language model reasoning and planning within the domain  
 501 of Blocksworld and its variant found in the PlanBench benchmark. Future work should expand into  
 502 Logistics problems (within PlanBench) or other extensible domains and other families of models.  
 503 Another key component of planning that introduces distinct challenges is multi-agent coordination  
 504 and path finding. Investigating whether TMK’s efficacy transfers to these distinct planning patterns  
 505 would help validate the generalizability of our findings beyond stacking tasks.  
 506

## 507 6 CONCLUSION

510 Our work indicates a general improvement of plan correctness on PlanBench Blocksworld variants.  
 511 Improvements across (LLM and LRM) models tested were shown for Classic Blocksworld  
 512 and Random Blocksworld, with Mystery Blocksworld showing improvement on all but o1-mini  
 513 model. We also saw our greatest increase in model performance in the flagship o1 model in Random  
 514 Blocksworld with a 65.8% increase in planning accuracy. Crucially, this result represents a funda-  
 515 mental ‘performance inversion’: whereas standard prompting struggles with opaque symbolic tasks  
 516 (Random) compared to semantic ones (Mystery), TMK prompting reversed this dynamic for the o1  
 517 model, enabling it to outperform its own semantic baseline.  
 518

519 This confirms that TMK acts as a symbolic scaffold, effectively steering reasoning models toward  
 520 formal code-like manipulation when semantic cues are absent. Improvements were consistent across  
 521 flagship models, however, the smaller o1-mini model proved to be an outlier, showing regression in  
 522 the Mystery domain, likely due to capacity limitations in resolving semantic interference. Larger  
 523 increases were demonstrated for models that performed worse on Blocksworld in general, but even  
 524 models such as GPT-5 which demonstrated a high starting success rate, still saw non-trivial in-  
 525 creases. For the LRMs, a more robust dataset or a different domain may be needed, as while TMK  
 526 did show increases in accuracy on planning tasks, their high starting success rate leaves little room  
 527 for improvement. We theorize this increase is due to the TMK prompting structure that intrinsically  
 528 allows language models to defer to a more code-adjacent inference. This contrasts with one that is  
 529 prompted with plain text drawing on a more linguistically adjacent inference.

530 There was a singular notable case where TMK prompting caused a decrease in accuracy for o1-mini  
 531 that warrants further investigation, although the authors theorize it is caused by semantic overload.  
 532 Determining what about the TMK framework causes an increase in random and classic blocksworld  
 533 but a decrease in mystery would be an interesting path for future research. Additionally, across  
 534 subdomains, we also noted that the introduction of TMK shifted a pattern in plain text where LLMs  
 535 often did better on mystery blocksworld when compared to random blocksworld. This could also be  
 a side effect of TMK, causing semantic overload.

536 For future work, we hope to investigate other domains, as well as evaluate how well TMK performs  
 537 when compared to other knowledge models such as BDI and HTNs. Our research gives some indi-  
 538 cation that prompting can increase the accuracy of LLMs in reasoning tasks within the PlanBench  
 539 blocks world variant, while avoiding the common criticism of previous research; however, the cause  
 of that increase is left to future work.

540 REFERENCES  
541542 Anonymous. Supplementary Materials. [https://osf.io/gufdm/?view\\_only=c0e90b9254a1432cb9f30e1e06297d23](https://osf.io/gufdm/?view_only=c0e90b9254a1432cb9f30e1e06297d23), 2025. Accessed: 2025-09-17.544 Viraat Aryabumi, Yixuan Su, Raymond Ma, Adrien Morisot, Ivan Zhang, Acyr Locatelli, Marzieh  
545 Fadaee, Ahmet Üstün, and Sara Hooker. To code, or not to code? exploring impact of code in  
546 pre-training. *arXiv preprint arXiv:2408.10914*, 2024.548 Siddhant Bhambri, Mudit Verma, and Subbarao Kambhampati. Do think tags really help llms plan?  
549 a critical evaluation of react-style prompting. *Transactions on Machine Learning Research*, 2025.

550 Emunah Chan. Understanding logical reasoning ability of large language models, 2024.

552 Yongchao Chen, Harsh Jhamtani, Srinagesh Sharma, Chuchu Fan, and Chi Wang. Steering large lan-  
553 guage models between code execution and textual reasoning. *arXiv preprint arXiv:2410.03524*,  
554 2024.555 Carlos G Correa, Mark K Ho, Frederick Callaway, Nathaniel D Daw, and Thomas L Griffiths.  
556 Humans decompose tasks by trading off utility and computational cost. *PLoS computational  
557 biology*, 19(6):e1011087, 2023.559 Rahul K Dass, Rochan H Madhusudhana, Erin C Deye, Shashank Verma, Timothy A Bydlon, Grace  
560 Brazil, and Ashok K Goel. Ivy: a hybrid knowledge-based and generative ai coach for explaining  
561 procedural skills. In *International Conference on Artificial Intelligence in Education*, pp. 233–  
562 246. Springer, 2025.563 Kutluhan Erol, James A Hendler, and Dana S Nau. *Semantics for hierarchical task-network plan-  
564 ning*. University of Maryland College Park, 1994.566 Alan Fern. The strips subset of pddl for the learning track of IPC-08. Technical re-  
567 port, School of Electrical Engineering and Computer Science, Oregon State University, April  
568 2008. URL [https://ipc08.icaps-conference.org/learning/documents/  
569 strips-pddl-subset.pdf](https://ipc08.icaps-conference.org/learning/documents/strips-pddl-subset.pdf). Accessed: 2025-09-24.571 Michael Georgeff, Barney Pell, Martha Pollack, Milind Tambe, and Michael Wooldridge. The  
572 belief-desire-intention model of agency. In *International workshop on agent theories, architec-  
573 tures, and languages*, pp. 1–10. Springer, 1998.574 Ashok K Goel and Spencer Rugaber. Gaia: A cad-like environment for designing game-playing  
575 agents. *IEEE Intelligent Systems*, 32(3):60–67, 2017.577 Malte Helmert. The fast downward planning system. *Journal of Artificial Intelligence Research*, 26:  
578 191–246, 2006.579 Richard Howey, Derek Long, and Maria Fox. Val: Automatic plan validation, continuous effects  
580 and mixed initiative planning using pddl. In *16th IEEE International Conference on Tools with  
581 Artificial Intelligence*, pp. 294–301. IEEE, 2004.583 Hanxu Hu, Hongyuan Lu, Huajian Zhang, Yun-Ze Song, Wai Lam, and Yue Zhang. Chain-of-  
584 symbol prompting elicits planning in large langauge models. *arXiv preprint arXiv:2305.10276*,  
585 2023.586 Jennifer Hu, Kyle Mahowald, Gary Lupyan, Anna Ivanova, and Roger Levy. Language models align  
587 with human judgments on key grammatical constructions. *Proceedings of the National Academy  
588 of Sciences*, 121(36):e2400917121, 2024. doi: 10.1073/pnas.2400917121. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2400917121>.591 Paul A Kirschner, John Sweller, Richard E Clark, PA Kirschner, and RE Clark. Why minimal guid-  
592 ance during instruction does not work: An analysis of the failure of constructivist. *Based Teaching  
593 Work: An Analysis of the Failure of Constructivist, Discovery, Problem-Based, Experiential, and  
Inquiry-Based Teaching*, (November 2014), pp. 37–41, 2010.

594 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
 595 language models are zero-shot reasoners. *Advances in neural information processing systems*,  
 596 35:22199–22213, 2022.

597

598 David R Krathwohl. A revision of bloom’s taxonomy: An overview. *Theory into practice*, 41(4):  
 599 212–218, 2002.

600

601 Pierre-Carl Langlais, Carlos Rosas Hinostroza, Mattia Nee, Catherine Arnett, Pavel Chizhov,  
 602 Eliot Krzysztof Jones, Irène Girard, David Mach, Anastasia Stasenko, and Ivan P Yamshchikov.  
 603 Common corpus: The largest collection of ethical data for llm pre-training. *arXiv preprint*  
 604 *arXiv:2506.01732*, 2025.

605

606 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and  
 607 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the*  
 608 *Association for Computational Linguistics*, 12:157–173, 2024.

609

610 Cherie Lum, Erin Deye, Grace Brazil, Tim Bydlon, Shashank Verma, Rochan Madhusudhana, Rahul  
 611 Dass, and Ashok Goel. Designing an ai coaching system for interactive video-based skill learning.  
 612 In *International Conference on Intelligent Tutoring Systems*, pp. 281–291. Springer, 2025.

613

614 Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The  
 615 sequential learning problem. volume 24 of *Psychology of Learning and Motivation*, pp. 109–165.  
 616 Academic Press, 1989. doi: [https://doi.org/10.1016/S0079-7421\(08\)60536-8](https://doi.org/10.1016/S0079-7421(08)60536-8). URL <https://www.sciencedirect.com/science/article/pii/S0079742108605368>.

617

618 J. W. Murdock. *Self -improvement through self -understanding: Model-based re-*  
 619 *flexion for agent adaptation.* PhD thesis, Georgia Institute of Technology,  
 620 2001. URL <https://www.proquest.com/dissertations-theses/self-improvement-through-understanding-model/docview/252085430/se-2>. Copyright - Database copyright ProQuest LLC; ProQuest does not claim copyright in the  
 621 individual underlying works; Last updated - 2023-02-23.

622

623 J William Murdock and Ashok K Goel. Meta-case-based reasoning: self-improvement through self-  
 624 understanding. *Journal of Experimental & Theoretical Artificial Intelligence*, 20(1):1–36, 2008.

625

626 OpenAI. OpenAI o1-mini: Advancing Cost-Efficient Reasoning. <https://openai.com/index/openai-o1-mini-advancing-cost-efficient-reasoning/>, 2025. Ac-  
 627 cessed: 2025-09-24.

628

629 Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha.  
 630 A systematic survey of prompt engineering in large language models: Techniques and applica-  
 631 tions. *arXiv preprint arXiv:2402.07927*, 2024.

632

633 Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad  
 634 Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning  
 635 models via the lens of problem complexity, 2025. URL <https://ml-site.cdn-apple.com/papers/the-illusion-of-thinking.pdf>.

636

637 Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. Chain of thoughtlessness? an  
 638 analysis of cot in planning. *Advances in Neural Information Processing Systems*, 37:29106–  
 29141, 2024.

639

640 Shalini Sushri, Rahul Dass, Rhea Basappa, Hong Lu, and Ashok Goel. Combining cognitive and  
 641 generative ai for self-explanation in interactive ai agents. *arXiv preprint arXiv:2407.18335*, 2024.

642

643 Karthik Valmeekam. Llms-planning: An extensible benchmark for evaluating large language models  
 644 on planning. <https://github.com/karthikv792/LLMs-Planning>, 2023. Accessed:  
 645 2025-09-24.

646

647 Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kamb-  
 648 hampati. Planbench: An extensible benchmark for evaluating large language models on planning  
 649 and reasoning about change. *Advances in Neural Information Processing Systems*, 36:38975–  
 38987, 2023a.

648 Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kamb-  
649 hampati. Planbench: An extensible benchmark for evaluating large language models on planning  
650 and reasoning about change. Poster presented at the Thirty-seventh Conference on Neural Infor-  
651 mation Processing Systems (NeurIPS 2023), 2023b. URL <https://neurips.cc/media/PosterPDFs/NeurIPS%202023/73553.png>. Accessed: 2025-09-18.  
652

653 Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. On the  
654 planning abilities of large language models-a critical investigation. *Advances in Neural Infor-  
655 mation Processing Systems*, 36:75993–76005, 2023c.  
656

657 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
658 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
659 tion processing systems*, 30, 2017.

660 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
661 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in  
662 neural information processing systems*, 35:24824–24837, 2022.  
663

664 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan  
665 Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international  
666 conference on learning representations*, 2022.

667 Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression  
668 for large language models. *Transactions of the Association for Computational Linguistics*, 12:  
669 1556–1577, 2024.  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

702 **A APPENDICES**  
703704 **Contents**  
705
706     A.1 One shot Classic Blocksworld with Plain Text Example  
707     A.2 Zero shot Classic Blocksworld with Plain Text Example  
708     A.3 One shot Classic Blocksworld with TMK Prompt Example  
709     A.4 One shot Mystery Blocksworld with Plain Text Example  
710     A.5 Zero shot Mystery Blocksworld with Plain Text Example  
711     A.6 One shot Mystery Blocksworld with TMK Prompt Example  
712     A.7 One shot Random Blocksworld with Plain Text Example  
713     A.8 Zero shot Random Blocksworld with Plain Text Example  
714     A.9 One shot Random Blocksworld with TMK Prompt Example  
715     A.10 Declaration  
716

---

717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
**A.1 ONE SHOT CLASSIC BLOCKSWORLD WITH PLAIN TEXT EXAMPLE****One shot Classic Blocksworld with Plain Text**

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.

[STATEMENT]  
 As initial conditions I have that, the red block is clear, the blue block is clear, the  
 ↪ yellow block is clear, the hand is empty, the blue block is on top of the orange  
 ↪ block, the red block is on the table, the orange block is on the table and the yellow  
 ↪ block is on the table.

My goal is to have that the orange block is on top of the blue block.

My plan is as follows:

[PLAN]  
 unstack the blue block from on top of the orange block  
 put down the blue block  
 pick up the orange block  
 stack the orange block on top of the blue block  
 [PLAN END]

[STATEMENT]  
 As initial conditions I have that, the red block is clear, the yellow block is clear, the  
 ↪ hand is empty, the red block is on top of the blue block, the yellow block is on top  
 ↪ of the orange block, the blue block is on the table and the orange block is on the  
 ↪ table.

```

756 My goal is to have that the orange block is on top of the red block.
757
758 My plan is as follows:
759
760 [PLAN]

```

## A.2 ZERO SHOT CLASSIC BLOCKSWORLD WITH PLAIN TEXT EXAMPLE

### Zero shot Classic Blocksworld with Plain Text

```

763 I am playing with a set of blocks where I need to arrange the blocks into stacks. Here
764 → are the actions I can do
765
766 Pick up a block
767 Unstack a block from on top of another block
768 Put down a block
769 Stack a block on top of another block
770
771 I have the following restrictions on my actions:
772 I can only pick up or unstack one block at a time.
773 I can only pick up or unstack a block if my hand is empty.
774 I can only pick up a block if the block is on the table and the block is clear. A block
775 → is clear if the block has no other blocks on top of it and if the block is not picked
776 → up.
777 I can only unstack a block from on top of another block if the block I am unstacking was
778 → really on top of the other block.
779 I can only unstack a block from on top of another block if the block I am unstacking is
780 → clear.
781 Once I pick up or unstack a block, I am holding the block.
782 I can only put down a block that I am holding.
783 I can only stack a block on top of another block if I am holding the block being stacked.
784 I can only stack a block on top of another block if the block onto which I am stacking
785 → the block is clear.
786 Once I put down or stack a block, my hand becomes empty.
787
788 [STATEMENT]
789 As initial conditions I have that, the red block is clear, the blue block is clear, the
790 → yellow block is clear, the hand is empty, the blue block is on top of the orange
791 → block, the red block is on the table, the orange block is on the table and the yellow
792 → block is on the table.
793 My goal is to have that the orange block is on top of the blue block.
794
795 What is the plan to achieve my goal? Just give the actions in the plan.

```

## A.3 ONE SHOT CLASSIC BLOCKSWORLD WITH TMK PROMPT EXAMPLE

### One shot Classic Blocksworld with TMK

```

793 You must adhere strictly to the JSON below, paying attention to the rules, ensuring to
794 → use only legal moves to achieve the final plan.
795 {
796   "Goals": [
797     {
798       "name": "PickUpBlock",
799       "description": "Pick up a block from the table.",
800       "inputParameters": [
801         "block",
802         "configuration"
803       ],
804       "outputParameters": [
805         "newConfiguration"
806     ],
807       "given": [
808         "On(block, table)",
809         "IsClear(block)",
810         "HandIsEmpty()"
811     ],
812       "makes": [
813         "Holding(block)",
814         "NOT On(block, table)",
815         "NOT HandIsEmpty()"
816     ],
817       "mechanism": "PickUpBlockMechanism"
818     },
819   ],
820 }

```

```

810
811    {
812        "name": "PutDownBlock",
813        "description": "Put down a held block onto the table.",
814        "inputParameters": [
815            "block",
816            "configuration"
817        ],
818        "outputParameters": [
819            "newConfiguration"
820        ],
821        "given": [
822            "Holding(block)"
823        ],
824        "makes": [
825            "On(block, table)",
826            "IsClear(block)",
827            "HandIsEmpty()"
828        ],
829        "mechanism": "PutDownBlockMechanism"
830    },
831    {
832        "name": "StackBlock",
833        "description": "Stack a held block onto another clear block.",
834        "inputParameters": [
835            "blockToStack",
836            "blockTarget",
837            "configuration"
838        ],
839        "outputParameters": [
840            "newConfiguration"
841        ],
842        "given": [
843            "Holding(blockToStack)",
844            "IsClear(blockTarget)"
845        ],
846        "makes": [
847            "On(blockToStack, blockTarget)",
848            "IsClear(blockToStack)",
849            "NOT IsClear(blockTarget)",
850            "HandIsEmpty()"
851        ],
852        "mechanism": "StackBlockMechanism"
853    },
854    {
855        "name": "UnstackBlock",
856        "description": "Unstack a block from on top of another block.",
857        "inputParameters": [
858            "blockToUnstack",
859            "blockFrom",
860            "configuration"
861        ],
862        "outputParameters": [
863            "newConfiguration"
864        ],
865        "given": [
866            "On(blockToUnstack, blockFrom)",
867            "IsClear(blockToUnstack)",
868            "HandIsEmpty()"
869        ],
870        "makes": [
871            "Holding(blockToUnstack)",
872            "IsClear(blockFrom)",
873            "NOT On(blockToUnstack, blockFrom)"
874        ],
875        "mechanism": "UnstackBlockMechanism"
876    },
877    ],
878    "Mechanisms": [
879        {
880            "name": "PickUpBlockMechanism",
881            "description": "Pick up {block} .",
882            "inputParameters": [
883                "block",
884                "configuration"
885            ],
886            "outputParameters": [
887                "newConfiguration"
888            ],
889            "type": "operation",
890            "requires": [
891                "On(block, table)",
892                "IsClear(block)"
893            ]
894        }
895    ]
896 }

```

```

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

```

```

        "HandIsEmpty() "
    ],
    "provides":[
        "Holding(block)",
        "Hand not empty",
        "Block not on table"
    ],
    "process":"Remove On(block, table), add Holding(block), set hand state"
},
{
    "name":"PutDownBlockMechanism",
    "description":"Put down {block}.",
    "inputParameters":[
        "block",
        "configuration"
    ],
    "outputParameters":[
        "newConfiguration"
    ],
    "type":"operation",
    "requires":[
        "Holding(block)"
    ],
    "provides":[
        "On(block, table)",
        "HandIsEmpty()",
        "IsClear(block)"
    ],
    "process":"Remove Holding(block), add On(block, table), clear hand state"
},
{
    "name":"StackBlockMechanism",
    "description":"Stack {blockToStack} on {blockTarget}.",
    "inputParameters":[
        "blockToStack",
        "blockTarget",
        "configuration"
    ],
    "outputParameters":[
        "newConfiguration"
    ],
    "type":"operation",
    "requires":[
        "Holding(blockToStack)",
        "IsClear(blockTarget)"
    ],
    "provides":[
        "On(blockToStack, blockTarget)",
        "HandIsEmpty()",
        "IsClear(blockToStack)"
    ],
    "process":"Remove Holding(blockToStack), add On(blockToStack, blockTarget),
    ↪ update clear states"
},
{
    "name":"UnstackBlockMechanism",
    "description":"Unstack {blockToUnstack} from {blockFrom}.",
    "inputParameters":[
        "blockToUnstack",
        "blockFrom",
        "configuration"
    ],
    "outputParameters":[
        "newConfiguration"
    ],
    "type":"operation",
    "requires":[
        "On(blockToUnstack, blockFrom)",
        "IsClear(blockToUnstack)",
        "HandIsEmpty()"
    ],
    "provides":[
        "Holding(blockToUnstack)",
        "IsClear(blockFrom)"
    ],
    "process":"Remove On(blockToUnstack, blockFrom), add Holding(blockToUnstack),
    ↪ update states"
}
],
"Knowledge":{
    "Concepts":[

```

```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```

```

        {
            "name": "block",
            "description": "A block in the blocks world that can be pick up, put down,
                           ↳ stacked or unstacked"
        },
        {
            "name": "table",
            "description": "The surface where blocks can be pick up, put down or unstacked
                           ↳ onto"
        },
        {
            "name": "hand",
            "description": "The manipulator that can pick up, put down, stacked or
                           ↳ unstacked blocks"
        },
        {
            "name": "IsClear",
            "description": "A block is clear if no other block is on top of it"
        },
        {
            "name": "HandIsEmpty",
            "description": "The hand is not holding any block"
        }

    ],
    "Relations": [
        {
            "name": "On",
            "description": "Relates a block to what it's on top of (another block or
                           ↳ table)"
        },
        {
            "name": "Holding",
            "description": "Relates the hand to the block it's holding"
        }
    ]
}
}

Below, within [Plan] and [Plan End], is the format you will use for the answer. The first
→ one is an example. Focus on only the second plan.

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

My plan is as follows:

[PLAN]
unstack the blue block from on top of the orange block
put down the blue block
pick up the orange block
stack the orange block on top of the blue block
[PLAN END]

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

My plan is as follows:

[PLAN]

```

#### A.4 ONE SHOT MYSTERY BLOCKSWORLD WITH PLAIN TEXT EXAMPLE

##### One shot Mystery BlocksWorld with Plain Text

```

I am playing with a set of objects. Here are the actions I can do

Attack object
Feast object from another object
Succumb object

```

```

972 Overcome object from another object
973
974 I have the following restrictions on my actions:
975 To perform Attack action, the following facts need to be true: Province object, Planet
976 → object, Harmony.
977 Once Attack action is performed the following facts will be true: Pain object.
978 Once Attack action is performed the following facts will be false: Province object,
979 → Planet object, Harmony.
980 To perform Succumb action, the following facts need to be true: Pain object.
981 Once Succumb action is performed the following facts will be true: Province object,
982 → Planet object, Harmony.
983 Once Succumb action is performed the following facts will be false: Pain object.
984 To perform Overcome action, the following needs to be true: Province other object, Pain
985 → object.
986 Once Overcome action is performed the following will be true: Harmony, Province object,
987 → Object Craves other object.
988 Once Overcome action is performed the following will be false: Province other object,
989 → Pain object.
990 To perform Feast action, the following needs to be true: Object Craves other object,
991 → Province object, Harmony.
992 Once Feast action is performed the following will be true: Pain object, Province other
993 → object.
994 Once Feast action is performed the following will be false:, Object Craves other object,
995 → Province object, Harmony.
996 My goal is to have that object c craves object b.

997 [STATEMENT]
998 As initial conditions I have that, object b craves object c, harmony, planet object a,
999 → planet object c, planet object d, province object a, province object b and province
1000 → object d.
1001 My goal is to have that object c craves object b.

1002 [PLAN]
1003
1004 My plan is as follows:
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

```

## A.5 ZERO SHOT MYSTERY BLOCKSWORLD WITH PLAIN TEXT EXAMPLE

### Zero shot Mystery Blocksworld with Plain Text

```

1009 I am playing with a set of objects. Here are the actions I can do
1010
1011 Attack object
1012 Feast object from another object
1013 Succumb object
1014 Overcome object from another object

1015 I have the following restrictions on my actions:
1016 To perform Attack action, the following facts need to be true: Province object, Planet
1017 → object, Harmony.
1018 Once Attack action is performed the following facts will be true: Pain object.
1019 Once Attack action is performed the following facts will be false: Province object,
1020 → Planet object, Harmony.
1021 To perform Succumb action, the following facts need to be true: Pain object.
1022 Once Succumb action is performed the following facts will be true: Province object,
1023 → Planet object, Harmony.
1024 Once Succumb action is performed the following facts will be false: Pain object.
1025 To perform Overcome action, the following needs to be true: Province other object, Pain
1026 → object.
1027 Once Overcome action is performed the following will be true: Harmony, Province object,
1028 → Object Craves other object.
1029 Once Overcome action is performed the following will be false: Province other object,
1030 → Pain object.

```

```

1026 To perform Feast action, the following needs to be true: Object Craves other object,
1027   ↪ Province object, Harmony.
1028 Once Feast action is performed the following will be true: Pain object, Province other
1029   ↪ object.
1030 Once Feast action is performed the following will be false:, Object Craves other object,
1031   ↪ Province object, Harmony.

1032 [STATEMENT]
1033 As initial conditions I have that, object b craves object c, harmony, planet object a,
1034   ↪ planet object c, planet object d, province object a, province object b and province
1035   ↪ object d.
1036 My goal is to have that object c craves object b.

1037 What is the plan to achieve my goal? Just give the actions in the plan.
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

```

## A.6 ONE SHOT MYSTERY BLOCKSWORLD WITH TMK PROMPT EXAMPLE

### One shot Classic Blocksworld with TMK Prompt Example

You must adhere strictly to the JSON below, paying attention to the rules, ensuring to  
 ↪ use only moves spelt out in the JSON to achieve the final plan.

```

{
  "Goals": [
    {
      "name": "AttackObject",
      "description": "Attack an object from the planet.",
      "inputParameters": [
        "object",
        "configuration"
      ],
      "outputParameters": [
        "newConfiguration"
      ],
      "given": {
        "Planet(object)": true,
        "Province(object)": true,
        "Harmony": true
      },
      "makes": {
        "Pain(object)": true,
        "Province(object)": false,
        "Planet(object)": false,
        "Harmony": false
      },
      "mechanism": "AttackObjectMechanism"
    },
    {
      "name": "SuccumbObject",
      "description": "Succumb a Pain object onto the planet.",
      "inputParameters": [
        "object",
        "configuration"
      ],
      "outputParameters": [
        "newConfiguration"
      ],
      "given": {
        "Pain(object)": true
      },
      "makes": {
        "Planet(object)": true,
        "Province(object)": true,
        "Harmony": true,
        "Pain(object)": false
      },
      "mechanism": "SuccumbObjectMechanism"
    },
    {
      "name": "OvercomeObject",
      "description": "Overcome a Pain object onto another Province object.",
      "inputParameters": [
        "objectToOvercome",
        "objectTarget",
        "configuration"
      ],
      "outputParameters": [
        "newConfiguration"
      ]
    }
  ]
}

```

```

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
  ],
  "given": {
    "Pain(objectToOvercome)": true,
    "Province(objectTarget)": true
  },
  "makes": {
    "ObjectCraves(objectToOvercome, objectTarget)": true,
    "Province(objectToOvercome)": true,
    "Province(objectTarget)": false,
    "Harmony": true,
    "Pain(objectToOvercome)": false
  },
  "mechanism": "OvercomeObjectMechanism"
},
{
  "name": "FeastObject",
  "description": "Feast an object from on top of another object (objectFrom).",
  "inputParameters": [
    "objectToFeast",
    "objectFrom",
    "configuration"
  ],
  "outputParameters": [
    "newConfiguration"
  ],
  "given": {
    "ObjectCraves(objectToFeast, objectFrom)": true,
    "Province(objectToFeast)": true,
    "Harmony": true
  },
  "makes": {
    "Pain(objectToFeast)": true,
    "Province(objectFrom)": true,
    "ObjectCraves(objectToFeast, objectFrom)": false,
    "Province(objectToFeast)": false,
    "Harmony": false
  },
  "mechanism": "FeastObjectMechanism"
}
],
"Mechanisms": [
  {
    "name": "AttackObjectMechanism",
    "description": "Attack {object}.",
    "inputParameters": [
      "object",
      "configuration"
    ],
    "outputParameters": [
      "newConfiguration"
    ],
    "type": "operation",
    "requires": {
      "Planet(object)": true,
      "Province(object)": true,
      "Harmony": true
    },
    "provides": {
      "Pain(object)": true,
      "Harmony": false,
      "Planet(object)": false,
      "Province(object)": false
    },
    "process": "Remove Planet(object), add Pain(object), remove Province(object),\n    ↳ set NOT Harmony"
  },
  {
    "name": "SuccumbObjectMechanism",
    "description": "Succumb {object}.",
    "inputParameters": [
      "object",
      "configuration"
    ],
    "outputParameters": [
      "newConfiguration"
    ],
    "type": "operation",
    "requires": {
      "Pain(object)": true
    },
    "provides": {
      "Planet(object)": true,
      "Province(object)": false
    }
  }
]

```

```

1134
1135     "Harmony": true,
1136     "Province(object)": true,
1137     "Pain(object)": false
1138   },
1139   {
1140     "name": "OvercomeObjectMechanism",
1141     "description": "Overcome {objectToOvercome} on {objectTarget}.",
1142     "inputParameters": [
1143       "objectToOvercome",
1144       "objectTarget",
1145       "configuration"
1146     ],
1147     "outputParameters": [
1148       "newConfiguration"
1149     ],
1150     "type": "operation",
1151     "requires": {
1152       "Pain(objectToOvercome)": true,
1153       "Province(objectTarget)": true
1154     },
1155     "provides": {
1156       "ObjectCraves(objectToOvercome, objectTarget)": true,
1157       "Harmony": true,
1158       "Province(objectToOvercome)": true,
1159       "Province(objectTarget)": false,
1160       "Pain(objectToOvercome)": false
1161     },
1162     "process": "Remove Pain(objectToOvercome), add ObjectCraves(objectToOvercome,
1163     ↳ objectTarget), add Province(objectToOvercome), remove
1164     ↳ Province(objectTarget), set Harmony"
1165   },
1166   {
1167     "name": "FeastObjectMechanism",
1168     "description": "Feast {objectToFeast} from {objectFrom}.",
1169     "inputParameters": [
1170       "objectToFeast",
1171       "objectFrom",
1172       "configuration"
1173     ],
1174     "outputParameters": [
1175       "newConfiguration"
1176     ],
1177     "type": "operation",
1178     "requires": {
1179       "ObjectCraves(objectToFeast, objectFrom)": true,
1180       "Province(objectToFeast)": true,
1181       "Harmony": true
1182     },
1183     "provides": {
1184       "Pain(objectToFeast)": true,
1185       "Province(objectFrom)": true,
1186       "ObjectCraves(objectToFeast, objectFrom)": false,
1187       "Province(objectToFeast)": false,
1188       "Harmony": false
1189     },
1190     "process": "Remove ObjectCraves(objectToFeast, objectFrom), add
1191     ↳ Pain(objectToFeast), add Province(objectFrom), remove
1192     ↳ Province(objectToFeast), set NOT Harmony"
1193   },
1194   ],
1195   "Knowledge": {
1196     "Concepts": [
1197       {
1198         "name": "object",
1199         "description": "An object in this domain that when is Province can be
1200         ↳ Attack from the Planet, Succumb onto the Planet, Overcome onto
1201         ↳ another object, or Feast from another object."
1202       },
1203       {
1204         "name": "hand",
1205         "description": "The manipulator that can Attack a object on the Planet,
1206         ↳ Succumb an object onto the Planet, Overcome onto another object, or
1207         ↳ Feast a Province object from another object. When hand is Harmony it
1208         ↳ can Attack, or Feast an object. When hand is Pain object, the same
1209         ↳ object can Succumb an object onto the Planet or Overcome another
1210         ↳ Province object."
1211       },
1212     ]
1213   }
1214 }
```

```

1188
1189         "name": "configuration",
1190         "description": "Complete state of this domain world."
1191     },
1192     "Relations": [
1193         {
1194             "name": "ObjectCraves",
1195             "description": "Binary Predicate: Relates an object to what it is on top
1196             ↪ of (another object), represented as ObjectCraves(object,
1197             ↪ anotherObject)."
1198         },
1199         {
1200             "name": "Pain",
1201             "description": "Unary Predicate: Relates the hand to the Pain object by
1202             ↪ setting Pain(object)."
1203         },
1204         {
1205             "name": "Planet",
1206             "description": "Unary Predicate: The surface where objects can be Attack
1207             ↪ from or Succumb onto, represented as Planet(object)."
1208         },
1209         {
1210             "name": "Province",
1211             "description": "Unary Predicate: An object is Province if no other object
1212             ↪ is on top of it, represented as Province(object)."
1213         },
1214         {
1215             "name": "Harmony",
1216             "description": "Predicate: The hand is Harmony, not Pain any object,
1217             ↪ represented as Harmony."
1218         }
1219     ]
1220 }
1221 Below, within [Plan] and [Plan End], is the format you will use for the answer. The first
1222 ↪ one is an example. Focus on only the second plan.
1223
1224 [STATEMENT]
1225 As initial conditions I have that, object b craves object c, harmony, planet object a,
1226 ↪ planet object c, planet object d, province object a, province object b and province
1227 ↪ object d.
1228 My goal is to have that object c craves object b.
1229
1230 My plan is as follows:
1231
1232 [PLAN]
1233 feast object b from object c
1234 succumb object b
1235 attack object c
1236 overcome object c from object b
1237 [PLAN END]
1238
1239 [STATEMENT]
1240 As initial conditions I have that: object a craves object b, object d craves object c,
1241 ↪ harmony, planet object b, planet object c, province object a and province object d.
1242 My goal is for the following to be true: object c craves object a.
1243
1244 My plan is as follows:
1245
1246 [PLAN]

```

## 1231 A.7 ONE SHOT RANDOM BLOCKSWORLD WITH PLAIN TEXT EXAMPLE

### 1233 One shot Random Blocksworld with Plain Text

```

1235 I am playing with a set of objects. Here are the actions I can do
1236 1jpkithdyjmlikck object
1237 xptxjrdkbi3pqsgqr object from another object
1238 9big8ruzarkkqyuu object
1239 21jg9q8swj2shjel object from another object
1240
1241 I have the following restrictions on my actions:
1242 To perform 1jpkithdyjmlikck action, the following facts need to be true: aqcjuuehiv18auwt
1243 ↪ object, 51nbwlachmfarjn object, 3covmuy4yrjthjd.

```

```

1242 Once 1jpkithdyjmlikck action is performed the following facts will be true:
1243   ↪ gk5asm3f7ulfeckpj object.
1244 Once 1jpkithdyjmlikck action is performed the following facts will be false:
1245   ↪ aqcjuuehivl8auwt object, 51nbwlachmfartjn object, 3covmuy4yrjthijd.
1246 To perform 9big8ruzarkkqy whole action, the following facts need to be true: gk5asm3f7ulfeckpj
1247   ↪ object.
1248 Once 9big8ruzarkkqy whole action is performed the following facts will be true:
1249   ↪ aqcjuuehivl8auwt object, 51nbwlachmfartjn object, 3covmuy4yrjthijd.
1250 Once 9big8ruzarkkqy whole action is performed the following facts will be false:
1251   ↪ gk5asm3f7ulfeckpj object.
1252 To perform 2ijg9q8swj2shjel action, the following needs to be true: aqcjuuehivl8auwt
1253   ↪ other object, gk5asm3f7ulfeckpj object.
1254 Once 2ijg9q8swj2shjel action is performed the following will be true: 3covmuy4yrjthijd,
1255   ↪ aqcjuuehivl8auwt object, Object 4DMF1cMTYXGSP94G other object.
1256 Once 2ijg9q8swj2shjel action is performed the following will be false: aqcjuuehivl8auwt
1257   ↪ other object, gk5asm3f7ulfeckpj object.
1258 To perform xptxjrdkbi3pqsrc action, the following needs to be true: Object
1259   ↪ 4DMF1cMTYXGSP94G other object, aqcjuuehivl8auwt object, 3covmuy4yrjthijd.
1260 Once xptxjrdkbi3pqsrc action is performed the following will be true: gk5asm3f7ulfeckpj
1261   ↪ object, aqcjuuehivl8auwt other object.
1262 Once xptxjrdkbi3pqsrc action is performed the following will be false:, Object
1263   ↪ 4DMF1cMTYXGSP94G other object, aqcjuuehivl8auwt object, 3covmuy4yrjthijd.

[STATEMENT]
1264 As initial conditions I have that: aqcjuuehivl8auwt object a, aqcjuuehivl8auwt object b,
1265   ↪ aqcjuuehivl8auwt object d, 3covmuy4yrjthijd, object b 4dmf1cmtxygsp94g object c,
1266   ↪ 51nbwlachmfartjn object a, 51nbwlachmfartjn object c and 51nbwlachmfartjn object d.
1267 My goal is for the following to be true: object c 4dmf1cmtxygsp94g object b.

1268 My plan is as follows:
1269 [PLAN]
1270 xptxjrdkbi3pqsrc object b from object c
1271 9big8ruzarkkqy object b
1272 1jpkithdyjmlikck object c
1273 2ijg9q8swj2shjel object c from object b
1274 [PLAN END]

[STATEMENT]
1275 As initial conditions I have that: aqcjuuehivl8auwt object a, aqcjuuehivl8auwt object d,
1276   ↪ 3covmuy4yrjthijd, object a 4dmf1cmtxygsp94g object b, object d 4dmf1cmtxygsp94g
1277   ↪ object c, 51nbwlachmfartjn object b and 51nbwlachmfartjn object c.
1278 My goal is for the following to be true: object c 4dmf1cmtxygsp94g object a.

1279 My plan is as follows:
1280 [PLAN]

```

## A.8 ZERO SHOT RANDOM BLOCKSWORLD WITH PLAIN TEXT EXAMPLE

1277 **Zero shot Random Blocksworld with Plain Text**

1278 I am playing with a set of objects. Here are the actions I can do

1279 1jpkithdyjmlikck object  
xptxjrdkbi3pqsrc object from another object  
9big8ruzarkkqy object  
2ijg9q8swj2shjel object from another object

1280 I have the following restrictions on my actions:  
1281 To perform 1jpkithdyjmlikck action, the following facts need to be true: aqcjuuehivl8auwt  
1282 ↪ object, 51nbwlachmfartjn object, 3covmuy4yrjthijd.  
1283 Once 1jpkithdyjmlikck action is performed the following facts will be true:  
1284 ↪ gk5asm3f7ulfeckpj object.  
1285 Once 1jpkithdyjmlikck action is performed the following facts will be false:  
1286 ↪ aqcjuuehivl8auwt object, 51nbwlachmfartjn object, 3covmuy4yrjthijd.  
1287 To perform 9big8ruzarkkqy whole action, the following facts need to be true: gk5asm3f7ulfeckpj  
1288 ↪ object.  
1289 Once 9big8ruzarkkqy whole action is performed the following facts will be true:  
1290 ↪ aqcjuuehivl8auwt object, 51nbwlachmfartjn object, 3covmuy4yrjthijd.  
1291 Once 9big8ruzarkkqy whole action is performed the following facts will be false:  
1292 ↪ gk5asm3f7ulfeckpj object.  
1293 To perform 2ijg9q8swj2shjel action, the following needs to be true: aqcjuuehivl8auwt  
1294 ↪ other object, gk5asm3f7ulfeckpj object.  
1295 Once 2ijg9q8swj2shjel action is performed the following will be true: 3covmuy4yrjthijd,  
1296 ↪ aqcjuuehivl8auwt object, Object 4DMF1cMTYXGSP94G other object.  
Once 2ijg9q8swj2shjel action is performed the following will be false: aqcjuuehivl8auwt  
1297 ↪ other object, gk5asm3f7ulfeckpj object.

```

1296 To perform xptxjrdkbi3pqsrc action, the following needs to be true: Object
1297 ↳ 4DMFlcMTYXGSP94G other object, aqcjuuehivl8auwt object, 3covmuy4yrjthijd.
1298 Once xptxjrdkbi3pqsrc action is performed the following will be true: gk5asm3f7ulfekpj
1299 ↳ object, aqcjuuehivl8auwt other object.
1300 Once xptxjrdkbi3pqsrc action is performed the following will be false:, Object
1301 ↳ 4DMFlcMTYXGSP94G other object, aqcjuuehivl8auwt object, 3covmuy4yrjthijd.

1302 [STATEMENT]
1303 As initial conditions I have that: aqcjuuehivl8auwt object a, aqcjuuehivl8auwt object b,
1304 ↳ aqcjuuehivl8auwt object d, 3covmuy4yrjthijd, object b 4dmflcmtyxgsp94g object c,
1305 ↳ 51nbwlachmfartjn object a, 51nbwlachmfartjn object c and 51nbwlachmfartjn object d.
1306 My goal is for the following to be true: object c 4dmflcmtyxgsp94g object b.

1307 To solve the problem, you will have to provide which actions to take from the initial
1308 conditions and in which order in order to achieve the goal conditions. Provide the
1309 ↳ plan by giving the action names along with the objects \"ACTION_NAME OBJECTS\".
1310 ↳ Provide the plan between these two tags [PLAN] and [PLAN END].

```

## A.9 ONE SHOT RANDOM BLOCKSWORLD WITH TMK PROMPT EXAMPLE

### One shot Random Blocksworld with TMK Prompt

You must adhere strictly to the JSON below, paying attention to the rules, ensuring to  
 ↳ use only moves spelt out in the JSON to achieve the final plan.

```

{
  "Goals": [
    {
      "name": "1jpkithdyjmlikck",
      "description": "1jpkithdyjmlikck an object from the 51nbwlachmfartjn.",
      "inputParameters": [
        "object",
        "configuration"
      ],
      "outputParameters": [
        "newConfiguration"
      ],
      "given": {
        "51nbwlachmfartjn(object)": true,
        "aqcjuuehivl8auwt(object)": true,
        "3covmuy4yrjthijd": true
      },
      "makes": {
        "gk5asm3f7ulfekpj(object)": true,
        "aqcjuuehivl8auwt(object)": false,
        "51nbwlachmfartjn(object)": false,
        "3covmuy4yrjthijd": false
      },
      "mechanism": "1jpkithdyjmlikckMechanism"
    },
    {
      "name": "9big8ruzarkkqyuu",
      "description": "9big8ruzarkkqyuu a gk5asm3f7ulfekpj object onto the
      ↳ 51nbwlachmfartjn.",
      "inputParameters": [
        "object",
        "configuration"
      ],
      "outputParameters": [
        "newConfiguration"
      ],
      "given": {
        "gk5asm3f7ulfekpj(object)": true
      },
      "makes": {
        "51nbwlachmfartjn(object)": true,
        "aqcjuuehivl8auwt(object)": true,
        "3covmuy4yrjthijd": true,
        "gk5asm3f7ulfekpj(object)": false
      },
      "mechanism": "9big8ruzarkkqyuuMechanism"
    },
    {
      "name": "2ijg9q8swj2shjel",
      "description": "2ijg9q8swj2shjel a gk5asm3f7ulfekpj object (objectToOvercome)
      ↳ onto another aqcjuuehivl8auwt object (objectTarget).",
      "inputParameters": [
        "objectToOvercome",

```

```

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
    "objectTarget",
    "configuration"
],
"outputParameters": [
    "newConfiguration"
],
"given": {
    "gk5asm3f7ulfekpj(objectToOvercome)": true,
    "aqcjuuehivl8auwt(objectTarget)": true
},
"makes": {
    "4dmf1cmtyxgsp94g(objectToOvercome, objectTarget)": true,
    "aqcjuuehivl8auwt(objectToOvercome)": true,
    "aqcjuuehivl8auwt(objectTarget)": false,
    "3covmuy4yrjthijd": true,
    "gk5asm3f7ulfekpj(objectToOvercome)": false
},
"mechanism": "2ijg9q8swj2shjelMechanism"
},
{
    "name": "xptxjrdkbi3pqsrc",
    "description": "xptxjrdkbi3pqsrc an object from on top of another object
    → (objectFrom).",
    "inputParameters": [
        "objectToFeast",
        "objectFrom",
        "configuration"
    ],
    "outputParameters": [
        "newConfiguration"
    ],
    "given": {
        "4dmf1cmtyxgsp94g(objectToFeast, objectFrom)": true,
        "aqcjuuehivl8auwt(objectToFeast)": true,
        "3covmuy4yrjthijd": true
    },
    "makes": {
        "gk5asm3f7ulfekpj(objectToFeast)": true,
        "aqcjuuehivl8auwt(objectFrom)": true,
        "4dmf1cmtyxgsp94g(objectToFeast, objectFrom)": false,
        "aqcjuuehivl8auwt(objectToFeast)": false,
        "3covmuy4yrjthijd": false
    },
    "mechanism": "xptxjrdkbi3pqsrcMechanism"
}
],
"Mechanisms": [
    {
        "name": "1jpkithdyjmlikckMechanism",
        "description": "1jpkithdyjmlikck {object}.",
        "inputParameters": [
            "object",
            "configuration"
        ],
        "outputParameters": [
            "newConfiguration"
        ],
        "type": "operation",
        "requires": {
            "51nbwlachmfartjn(object)": true,
            "aqcjuuehivl8auwt(object)": true,
            "3covmuy4yrjthijd": true
        },
        "provides": {
            "gk5asm3f7ulfekpj(object)": true,
            "3covmuy4yrjthijd": false,
            "51nbwlachmfartjn(object)": false,
            "aqcjuuehivl8auwt(object)": false
        },
        "process": "Remove 51nbwlachmfartjn(object), add gk5asm3f7ulfekpj(object),
        → remove aqcjuuehivl8auwt(object), set NOT 3covmuy4yrjthijd"
    },
    {
        "name": "9big8ruzarkkquyMechanism",
        "description": "9big8ruzarkkquy {object}.",
        "inputParameters": [
            "object",
            "configuration"
        ],
        "outputParameters": [

```

```

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
    "newConfiguration"
  ],
  "type": "operation",
  "requires": {
    "gk5asm3f7ulfekpj(object)": true
  },
  "provides": {
    "51nbwlachmfartjn(object)": true,
    "3covmuy4yrjthijd": true,
    "aqcjuuehivl8auwt(object)": true,
    "gk5asm3f7ulfekpj(object)": false
  },
  "process": "Remove gk5asm3f7ulfekpj(object), add 51nbwlachmfartjn(object),
  ↪ add aqcjuuehivl8auwt(object), set 3covmuy4yrjthijd"
},
{
  "name": "2ijg9q8swj2shjelMechanism",
  "description": "2ijg9q8swj2shjel {objectToOvercome} on {objectTarget}.",
  "inputParameters": [
    "objectToOvercome",
    "objectTarget",
    "configuration"
  ],
  "outputParameters": [
    "newConfiguration"
  ],
  "type": "operation",
  "requires": {
    "gk5asm3f7ulfekpj(objectToOvercome)": true,
    "aqcjuuehivl8auwt(objectTarget)": true
  },
  "provides": {
    "4dmf1cmtyxgsp94g(objectToOvercome, objectTarget)": true,
    "3covmuy4yrjthijd": true,
    "aqcjuuehivl8auwt(objectToOvercome)": true,
    "aqcjuuehivl8auwt(objectTarget)": false,
    "gk5asm3f7ulfekpj(objectToOvercome)": false
  },
  "process": "Remove gk5asm3f7ulfekpj(objectToOvercome), add
  ↪ 4dmf1cmtyxgsp94g(objectToOvercome, objectTarget), add
  ↪ aqcjuuehivl8auwt(objectToOvercome), remove
  ↪ aqcjuuehivl8auwt(objectTarget), set 3covmuy4yrjthijd"
},
{
  "name": "xptxjrdkbi3pqsrcMechanism",
  "description": "xptxjrdkbi3pqsrc {objectToFeast} from {objectFrom}.",
  "inputParameters": [
    "objectToFeast",
    "objectFrom",
    "configuration"
  ],
  "outputParameters": [
    "newConfiguration"
  ],
  "type": "operation",
  "requires": {
    "4dmf1cmtyxgsp94g(objectToFeast, objectFrom)": true,
    "aqcjuuehivl8auwt(objectToFeast)": true,
    "3covmuy4yrjthijd": true
  },
  "provides": {
    "gk5asm3f7ulfekpj(objectToFeast)": true,
    "aqcjuuehivl8auwt(objectFrom)": true,
    "4dmf1cmtyxgsp94g(objectToFeast, objectFrom)": false,
    "aqcjuuehivl8auwt(objectToFeast)": false,
    "3covmuy4yrjthijd": false
  },
  "process": "Remove 4dmf1cmtyxgsp94g(objectToFeast, objectFrom), add
  ↪ gk5asm3f7ulfekpj(objectToFeast), add aqcjuuehivl8auwt(objectFrom), remove
  ↪ aqcjuuehivl8auwt(objectToFeast), set NOT 3covmuy4yrjthijd"
},
],
"Knowledge": {
  "Concepts": [
    {
      "name": "object",
      "description": "An object in this domain that when is aqcjuuehivl8auwt
      ↪ can be 1jpkithdyjmlikck from the 51nbwlachmfartjn, 9big8ruzarkkqyu
      ↪ onto the 51nbwlachmfartjn, 2ijg9q8swj2shjel onto another object, or
      ↪ xptxjrdkbi3pqsrc from another object."
    }
  ]
}

```

```

1458
1459    {
1460        "name": "hand",
1461        "description": "The manipulator that can 1jkpithdyjmlikck an object from
1462        the 5lnbwachmfartjn, 9big8ruzarkkqy whole object onto the
1463        5lnbwachmfartjn, 2ijg9q8swj2shjel onto another object, or
1464        xptxjrdkbi3pqsrc an aqcjuuehivl8auwt object from another object. When
1465        hand is 3covmuy4yrjthijd it can 1jkpithdyjmlikck or xptxjrdkbi3pqsrc
1466        an object. When hand is gk5asm3f7ulfekpj object, the same object can
1467        9big8ruzarkkqy whole object onto the 5lnbwachmfartjn or
1468        2ijg9q8swj2shjel another aqcjuuehivl8auwt object."
1469    },
1470    {
1471        "name": "configuration",
1472        "description": "Complete state of this domain world."
1473    },
1474    {
1475        "name": "isWellFormed",
1476        "description": "Configuration follows all domain rules."
1477    },
1478    {
1479        "name": "matches",
1480        "description": "Two configurations are identical."
1481    }
1482 },
1483 "Relations": [
1484     {
1485         "name": "4dmflcmtyxgsp94g",
1486         "description": "Binary Predicate: Relates an object to what it is on top
1487         of (another object), represented as 4dmflcmtyxgsp94g(object,
1488         anotherObject)."
1489     },
1490     {
1491         "name": "gk5asm3f7ulfekpj",
1492         "description": "Unary Predicate: Relates the hand to the held object by
1493         setting gk5asm3f7ulfekpj(object)."
1494     },
1495     {
1496         "name": "5lnbwachmfartjn",
1497         "description": "Unary Predicate: The surface where objects can be picked
1498         up from or put down onto, represented as 5lnbwachmfartjn(object)."
1499     },
1500     {
1501         "name": "aqcjuuehivl8auwt",
1502         "description": "Unary Predicate: An object is clear if no other object is
1503         on top of it, represented as aqcjuuehivl8auwt(object)."
1504     },
1505     {
1506         "name": "3covmuy4yrjthijd",
1507         "description": "Zero-arity Predicate: The hand is empty, represented as
1508         3covmuy4yrjthijd."
1509     }
1510 ],
1511 }
1512 }
1513 Below, within [Plan] and [Plan End], is the format you will use for the answer. The first
1514 one is an example. Focus on only the second plan.
1515
1516 [STATEMENT]
1517 As initial conditions I have that: aqcjuuehivl8auwt object a, aqcjuuehivl8auwt object b,
1518 aqcjuuehivl8auwt object d, 3covmuy4yrjthijd, object b 4dmflcmtyxgsp94g object c,
1519 5lnbwachmfartjn object a, 5lnbwachmfartjn object c and 5lnbwachmfartjn object d.
1520 My goal is for the following to be true: object c 4dmflcmtyxgsp94g object b.
1521
1522 My plan is as follows:
1523
1524 [PLAN]
1525 xptxjrdkbi3pqsrc object b from object c
1526 9big8ruzarkkqy whole object b
1527 1jkpithdyjmlikck object c
1528 2ijg9q8swj2shjel object c from object b
1529 [PLAN END]
1530
1531 [STATEMENT]
1532 As initial conditions I have that: aqcjuuehivl8auwt object a, aqcjuuehivl8auwt object d,
1533 3covmuy4yrjthijd, object a 4dmflcmtyxgsp94g object b, object d 4dmflcmtyxgsp94g
1534 object c, 5lnbwachmfartjn object b and 5lnbwachmfartjn object c.
1535 My goal is for the following to be true: object c 4dmflcmtyxgsp94g object a.
1536
1537 My plan is as follows:
1538
1539 [PLAN]

```

1512 A.10 DECLARATION  
15131514 Declaration of LLM usage: LLMs were used to discover and trace planbench code issues, formatting  
1515 of figures and tables, retrieval and discovery of existing research papers (mixed with google search).  
1516 In the rebuttal phase, language models were used suggest edits of existing content.  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565