

TEXT2WORLD: Benchmarking Large Language Models for Symbolic World Model Generation

Anonymous ACL submission

Abstract

001 Recently, there has been growing interest in
002 leveraging large language models (LLMs) to
003 generate symbolic world models from textual
004 descriptions. Although LLMs have been exten-
005 sively explored in the context of world mod-
006 eling, prior studies encountered several chal-
007 lenges, including evaluation randomness, de-
008 pendence on indirect metrics, and a limited
009 domain scope. To address these limitations, we
010 introduce a novel benchmark, TEXT2WORLD,
011 based on planning domain definition language
012 (PDDL), featuring hundreds of diverse domains
013 and employing multi-criteria, execution-based
014 metrics for a more robust evaluation. We bench-
015 mark current LLMs using TEXT2WORLD and
016 find that reasoning models trained with large-
017 scale reinforcement learning outperform oth-
018 ers. However, even the best-performing model
019 still demonstrates limited capabilities in world
020 modeling. Building on these insights, we ex-
021 amine several promising strategies to enhance
022 the world modeling capabilities of LLMs, in-
023 cluding test-time scaling, agent training, and
024 more. We hope that TEXT2WORLD can serve
025 as a crucial resource, laying the groundwork for
026 future research in leveraging LLMs as world
027 models.

028 1 Introduction

029 The significance of world models for intelligent be-
030 havior has been historically acknowledged in early
031 psychological theories, which posited that organ-
032 isms employ internal representations of the exter-
033 nal world for prediction and planning (Craig, 1967).
034 Furthermore, LeCun (2022) extends this concept
035 by highlighting world modeling as a core compo-
036 nent of autonomous machine intelligence. In this
037 paper, we primarily study *symbolic world models*
038 (also known as domain models), which are formal
039 representations of an environment’s dynamics and
040 constraints. In recent years, Large Language Mod-
041 els (LLMs) (OpenAI, 2022, 2023; Meta AI, 2024)

042 have showcased their understanding of common-
043 world knowledge, making them promising candi-
044 dates for generating symbolic world models, which
045 requires inferring action dynamics and constraints
046 from solely natural language description. Some
047 works have already explored this across numerous
048 tasks, including planning (Hu et al., 2024b; Guan
049 et al., 2023), game design (Wang et al., 2023a,
050 2024), reinforcement learning (Tang et al., 2024)
051 among others.

052 Despite extensive exploration, previous work for
053 evaluating symbolic world model generation suf-
054 fers from several key limitations: (i) **Limited Do-
055 main Scope:** These studies are often confined to
056 a narrow set of domains (typically fewer than 20),
057 which limits the generalizability and applicabil-
058 ity of their findings (Oswald et al., 2024; Silver
059 et al., 2024; Wong et al., 2023). (ii) **Evaluation
060 Randomness:** Some works rely on LLM-based
061 evaluation methods, which may introduce addi-
062 tional margins of error (Wang et al., 2023a). Pre-
063 liminary experiments in Section 3.6 demonstrate
064 that the LLM-based evaluation exhibits a low inter-
065 annotator agreement with human annotators (Co-
066 hen’s $\kappa = 0.10$). (iii) **Indirect Evaluation:** Some
067 studies evaluate world models based on end-to-end
068 success rates in model-based planning, making it
069 difficult to identify specific failure modes (Guan
070 et al., 2023; Dainese et al., 2024).

071 Motivated by these issues, this paper introduces
072 a novel benchmark TEXT2WORLD based on the
073 Planning Domain Definition Language (PDDL) as
074 illustrated in Figure 1. Specifically, to address
075 the first issue, we initially gathered a broad set
076 of domains, which were then filtered through an
077 automated pipeline and manually curated to ensure
078 their quality, ultimately resulting in a collection
079 of hundreds of diverse domains. Furthermore, to
080 tackle the second issue, we designed multi-criteria,
081 execution-based metrics to ensure a more robust as-
082 sessment. Specifically, we not only employed struc-

tural similarity for an overall evaluation but also designed component-wise F1 scores to assess finer-grained aspects such as action dynamics. Moreover, to overcome the third issue, we systematically apply these metrics to assess the generated world model directly, eliminating reliance on indirect feedback mechanisms.

We also performed data contamination analysis using n-gram matching (Touvron et al., 2023), revealing a lower contamination rate ($\mu = 0.04$) compared to prior works (Guan et al., 2023; Smirnov et al., 2024), indicating that TEXT2WORLD effectively evaluates LLMs’ world modeling capabilities rather than pattern memorization.

We used TEXT2WORLD to benchmark the world modeling capabilities of 16 different LLMs from 9 model families. Experimental results in Table 1 highlight several key findings: (i) *The most advanced LLMs still struggle with TEXT2WORLD*; (ii) *large reasoning models trained by reinforcement learning show stronger world modeling capabilities*; and (iii) *error correction significantly improves model performance*. To gain a deeper understanding, we performed a manual error analysis and found that the majority were due to the LLMs’ inability to include essential preconditions or effects. We also explored several strategies to enhance the world modeling capabilities of LLMs. Specifically, we initially experimented with scaling the test-time budget and observed consistent improvements as the test-time budget increased. Additionally, methods like fine-tuning and in-context learning contributed positively to model effectiveness. Moreover, we found that supervised fine-tuning on agent trajectory data yielded unexpected gains, underscoring the importance of robust world modeling for developing high-performing agents.

To facilitate further research, benchmark and code are available at this [anonymous URL](#).

2 Preliminary

2.1 World Model

We formally define a symbolic world model as $D = \langle F, A \rangle$, where F represents the set of fluents (state variables represented as predicates) and A is the set of possible actions. Each fluent $f \in F$ is a predicate of the form $p(x_1, \dots, x_n)$, where p is the predicate name and x_1, \dots, x_n are typed variables. Each action $a \in A$ is defined as a tuple $a = \langle \alpha, \mathcal{P}, \varphi, \mathcal{E} \rangle$ where: i) α denotes the action signature (identifier); ii) \mathcal{P} represents a list of typed

parameters (p_1, \dots, p_k) ; iii) φ specifies the preconditions: a logical formula over fluents that must hold for the action to be applicable; and iv) \mathcal{E} defines the effects: a set of fluent literals describing how the action changes the world state.

2.2 Task Definition

The task is formally defined as: $\mathcal{M} : \mathcal{N} \rightarrow D, D \models \mathcal{N}$, where \mathcal{M} is a mapping function (implemented by an LLM) that generates world model D from the natural language description \mathcal{N} . \models denotes semantic satisfaction. Each \mathcal{N} contains the following components: i) A general description describing the overall objective of the domain; ii) A set of predicates $\mathcal{N}_F = \{f_1, \dots, f_n\}$ where each predicate is described with its signature (e.g., “*conn ?x ?y*”) and an explanation (e.g., “*Indicates a connection between two places ?x and ?y*”); iii) A set of actions $\mathcal{N}_A = \{a_1, \dots, a_m\}$ where each action is described with: its signature (e.g., “*move <?curpos> <?nextpos>*”) and an explanation (e.g., “*Allows the robot to move from place <?curpos> to place <?nextpos>*”). Note that to evaluate LLMs’ inherent world modeling capabilities, action descriptions in \mathcal{N}_A are intentionally kept at a high level, without explicit specifications of preconditions φ and effects \mathcal{E} . This design choice allows us to assess how well LLMs can infer the underlying world dynamics and constraints from purely descriptive text. A comparative analysis of model performance conditioned on different description styles is presented in Section 6.5.

2.2.1 Evaluation Metrics

We directly evaluate generated world models, addressing the ambiguity associated with indirect evaluations (Guan et al., 2023; Dainese et al., 2024). In addition, we propose using execution-based metrics, overcoming the randomness of LLM-based evaluation (Wang et al., 2023a). Specifically, we established the following evaluation metrics: (i) **Executability (EXEC.)**: Measures whether the generated PDDL can be successfully parsed and validated by standard PDDL validators. (ii) **Structural Similarity (SIM.)**: Quantifies the textual similarity between the generated and ground truth PDDL using normalized Levenshtein ratio. (iii) **Component-wise F1 Scores**: When generated PDDL achieves executability (EXEC. = 1), we perform fine-grained analysis by calculating the macro-averaged F1 score for each component type (predicates, actions, etc.). More specifically, we compute F1 scores for

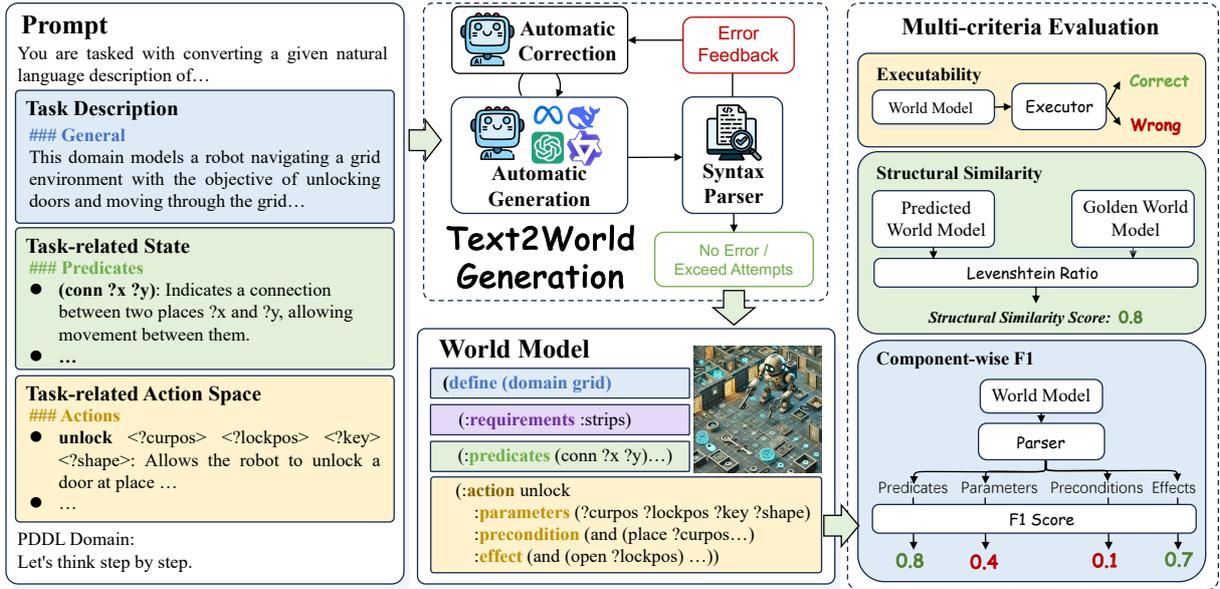


Figure 1: Overview of TEXT2WORLD.

predicates ($F1_{\text{PRED}}$), parameters ($F1_{\text{PARAM}}$), pre-conditions ($F1_{\text{PRECOND}}$), and effects ($F1_{\text{EFF}}$) by parsing both generated and ground truth PDDL into structured representations.

3 Benchmark Construction

The overall process of benchmark construction is shown in Figure 2. In this section, we provide a detailed explanation of each stage.

3.1 Data Acquisition

Our benchmark construction process began with collecting PDDL files from various public repositories and planning competitions. Through this initial collection phase, we accumulated 1,801 raw PDDL files. We performed several preprocessing steps to standardize the data format (e.g., convert files with BOM encoding to standard UTF-8). The processed files served as the foundation for our dataset construction.

3.2 Data Filtering and Manual Selection

To ensure the quality and reliability of TEXT2WORLD, we implemented a comprehensive filtering pipeline: (i) **Validation**: We employed a PDDL domain parser to perform syntax validation on each file; (ii) **Similarity Deduplication**: We eliminated duplicate entries by computing pairwise cosine similarity on TF-IDF vectorized PDDL content, removing files with similarity scores exceeding 0.9; (iii) **Complexity Control**: We removed domains with over 40

predicates or 20 actions to balance expressiveness with practical utility. (iv) **Token Length Filtering**: We removed files exceeding 5,000 tokens using GPT-2 (Radford et al., 2019) tokenizer to ensure compatibility with model context windows. Additionally, we conducted manual selection to eliminate domains that were not designed for world modeling (such as blocksworld-mystery) and low-quality cases that were not captured by the automated filtering methods. After this process, we obtained 264 high-quality PDDL domain specifications.

3.3 Data Annotation

After obtaining the high-quality PDDL domains, we manually annotated natural language descriptions for each domain. To ensure the quality of annotations, we recruited 6 computer science graduates as annotators. The annotated description followed the structured format described in Section 2.2, and annotators were required to follow the annotation criteria: (i) **Descriptive Completeness**: Annotations must contain all required components; (ii) **Action Abstraction**: Action descriptions should avoid explicit references to formal preconditions and effects; (iii) **Inference-Enabling**: Descriptions should contain sufficient contextual information to allow models to infer the underlying dynamics; (iv) **Natural Language Priority**: Technical terminology should be minimized in favor of natural language explanations. Examples of TEXT2WORLD can be found in Appendix A.1.

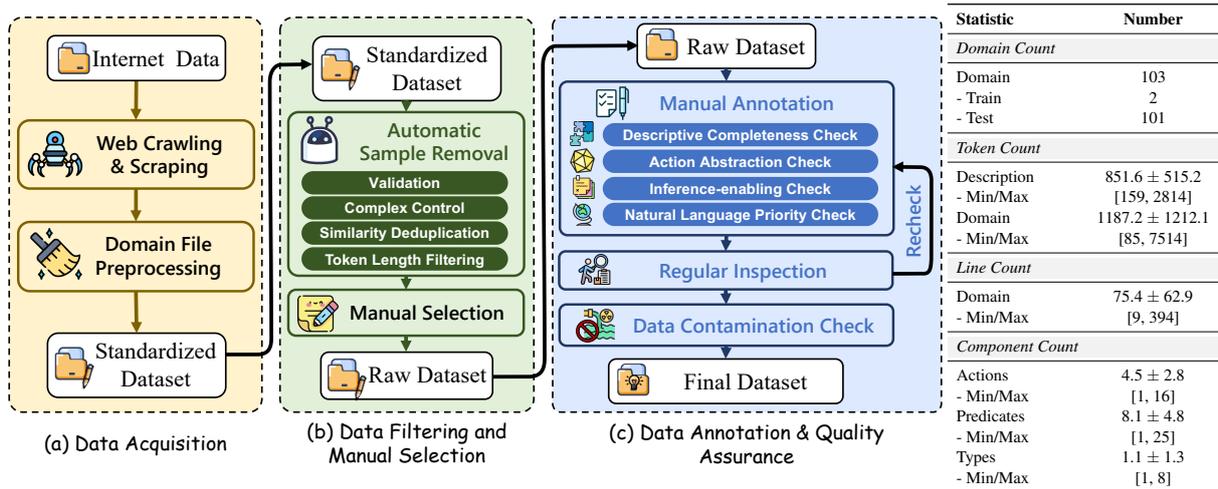


Figure 2: *Left*: Dataset construction process including: (a) *Data Acquisition* (§3.1); (b) *Data Filtering and Manual Selection* (§3.2); (c) *Data Annotation and Quality Assurance* (§3.3 and §3.4). *Right*: Key statistics of Text2World. Tokens are counted by GPT-2 (Radford et al., 2019) tokenizer. The style is referenced from Chen et al. (2024b).

3.4 Quality Assurance

Manual Recheck To maintain rigorous quality standards throughout the annotation process, we established a review system supervised by two senior experts. These experts conducted regular inspections of the annotations, ensuring accuracy and consistency. Inspectors must verify all data twice to determine if the annotated examples meet the specified annotation standards. Examples are accepted only if both inspectors approve them. The verification results showed "almost perfect agreement" with a Fleiss Kappa (Landis and Koch, 1977) score of 0.82. Through this comprehensive quality control process, we compiled a final curated dataset of 103 domains with gold-standard descriptions.

Data Contamination As shown by Carlini et al. (2021), LLMs can memorize training data rather than truly model the world. To assess potential contamination between LLMs' training data and TEXT2WORLD, we generated complete PDDL domains from the first 20 tokens using GPT-4 (OpenAI, 2023) and calculated contamination rates based on tokenized 10-grams with up to 4 mismatches (Touvron et al., 2023), excluding PDDL-specific keywords and variables. We also compared these results with previous studies (Guan et al., 2023; Smirnov et al., 2024). Figure 3 shows that TEXT2WORLD has a lower contamination rate ($\mu = 0.04$ vs. $\mu = 0.47$), suggesting its performance reflects domain understanding rather than memorization. However, the complete elimination of contamination remains challenging due to PDDL's widespread use.

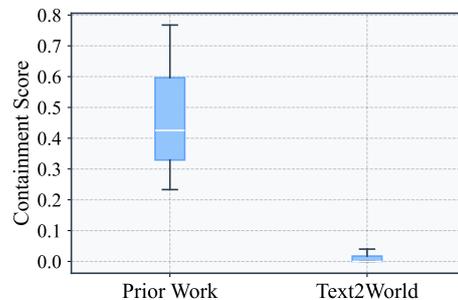


Figure 3: n-gram contamination rate of TEXT2WORLD and prior works.

3.5 Data Analysis

This section provides some detailed data analysis to better understand TEXT2WORLD.

Core Statistics We designated 2 domains as in-context exemplars (train set), with the remaining 101 samples forming our test set.

Semantic Analysis We use LLMs to extract high-level domain characteristics to better understand the conceptual distribution of TEXT2WORLD, As shown in Figure 4 (Bottom), common themes such as *path planning*, *constraint satisfaction*, and *task allocation*, among others, emerge.

Requirements Analysis A PDDL requirement specifies a formal capability needed to express a domain, often reflecting its complexity. For instance, `:typing` stands for allowing the usage of typing for objects. As shown in Figure 4 (Top), there are eight different requirement type in TEXT2WORLD. We also provide an in-depth analysis of requirement type in Appendix A.3.

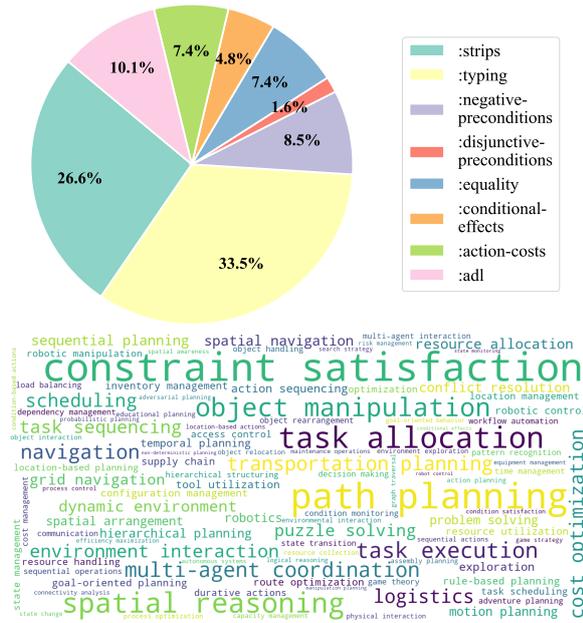


Figure 4: *Top*: The frequency of requirements distribution. *Bottom*: Word cloud of concepts in TEXT2WORLD.

3.6 Preliminary Experiment

In previous works, LLMs have been employed to evaluate the action dynamics of world models generated by LLMs themselves (Wang et al., 2023a). To further assess the ability of LLMs to detect errors in world models, we conducted a preliminary experiment where we first used claude-3.5-sonnet for TEXT2WORLD. Subsequently, human annotators and the LLM independently evaluated the generated action dynamics to identify potential errors. The inter-annotator agreement between human ratings and LLM ratings, measured using Cohen’s κ , was 0.10, indicating a low level of agreement. This suggests that predicting the correctness of PDDL domains using an LLM is particularly challenging, highlighting the need for more discriminative evaluation metrics. Prompting examples and more results can be found in Appendix A.2.

4 Experiments

4.1 Experimental Setup

We evaluate several state-of-the-art LLMs, including *GPT-4* (OpenAI, 2023), *GPT-3.5* (OpenAI, 2022), *Claude-3.5* (Anthropic), and *LLaMA-3.1* (Meta AI), *DeepSeek-v3* (Liu et al., 2024), *CodeLlama* (Roziere et al., 2023), *LlaMA-2* (Touvron et al., 2023), etc. We also evaluated Large Reasoning Models (LRMs) trained using reinforce-

ment learning, such as *DeepSeek-R1* (DeepSeek-AI et al., 2025), *OpenAI-o1* (OpenAI, 2024) and *OpenAI-o3* (OpenAI, 2025). We set temperature = 0 for each model for all experiments to maintain reproducibility. We employ tarski¹ library to check syntactic correctness and executability. We prompt LLMs to generate symbolic world models under a zero-shot setting with chain-of-thought reasoning (Wei et al., 2022). In error-correction experiments, LLMs refine outputs based on validator-reported syntax errors, denoted as EC₃ for k attempts. Evaluation of open-sourced models were conducted on NVIDIA A100 GPUs with 80GB memory. We access proprietary models through their official API platform. Prompt examples can be found in Appendix B.2.

4.2 Experimental Results

Several conclusions can be drawn from Table 1: (i) **The most advanced LLMs still struggle with TEXT2WORLD.** For example, the best-performing model, *DeepSeek-R1*, achieves F1 scores below 60% for both preconditions (F1_{PRECOND}) and effects (F1_{EFF}) under the without error correction setting. This highlights the limitations of current LLMs in world modeling tasks. (ii) **Large reasoning models trained with reinforcement learning exhibit superior world modeling capabilities.** These models, such as *DeepSeek-R1* (DeepSeek-AI et al., 2025), outperform others in executability, structural similarity, and component-wise performance, indicating that RL-based training enhances the ability of models to generate structured and valid world models. (iii) **The ability of models to benefit from error correction is evident.** For instance, *GPT-4* (gpt-4o-mini) demonstrates a notable improvement in executability, increasing from 48.5% to 72.3% after three correction attempts.

5 Analysis

5.1 Statistical Analysis

We conducted a one-way ANOVA (Girden, 1992) to evaluate the impact of correction attempts on model performance, excluding anomalous zero values. The results showed a significant improvement with three correction attempts ($F = 27.48, p = 0.00012$), indicating that correction attempts lead to a notable enhancement in model performance.

¹<https://github.com/aig-upf/tarski>

Table 1: Performance comparison of different LLMs on TEXT2WORLD. EC_k denotes the setting where models are allowed k correction attempts (EC_0 : zero-shot without correction, EC_3 : with 3 correction attempts).

Model Family	Version	EXEC. \uparrow		SIM. \uparrow		$F1_{\text{PRED}} \uparrow$		$F1_{\text{PARAM}} \uparrow$		$F1_{\text{PRECOND}} \uparrow$		$F1_{\text{EFF}} \uparrow$	
		EC_0	EC_3	EC_0	EC_3	EC_0	EC_3	EC_0	EC_3	EC_0	EC_3	EC_0	EC_3
OPENAI-O1	o1-mini	49.5	69.3	82.5	82.2	48.4	66.3	36.4	49.7	28.9	38.0	31.7	42.1
OPENAI-O3	o3-mini	54.5	84.2	83.0	81.9	53.9	81.1	43.7	63.0	36.8	50.4	39.4	53.8
GPT-4	gpt-4o	60.4	75.2	84.5	84.1	59.6	72.1	56.5	68.1	49.3	56.4	47.8	56.7
	gpt-4o-mini	48.5	72.3	82.6	82.2	48.1	70.1	47.1	67.3	34.9	47.5	38.2	52.7
GPT-3.5	turbo-0125	41.6	56.4	81.9	81.6	41.2	55.8	39.6	53.8	30.2	39.2	27.5	37.7
CLAUDE-3.5	sonnet	45.5	64.4	73.2	66.8	45.5	62.5	41.5	48.8	37.4	44.0	38.4	45.0
LLAMA-2	7b-instruct	0.0	0.0	45.5	33.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	70b-instruct	0.0	0.0	48.7	48.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
LLAMA-3.1	8b-instruct	0.0	0.0	74.3	74.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	70b-instruct	0.0	0.0	83.6	79.2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DEEPSEEK	deepseek-v3	56.4	79.2	84.7	84.2	55.9	75.6	53.7	74.4	45.1	58.6	46.7	61.5
	deepseek-r1	72.3	89.1	84.3	84.0	71.7	86.7	64.0	76.3	57.6	65.0	58.8	67.3
CODELLAMA	7b-instruct	17.8	22.8	60.2	57.6	17.8	18.8	17.2	18.2	11.3	12.2	10.7	11.1
	13b-instruct	7.9	8.9	57.6	55.0	7.9	8.9	7.9	8.9	4.9	5.9	5.2	6.1
	34b-instruct	7.9	8.9	34.2	7.6	7.9	8.6	7.9	8.4	5.0	5.0	5.4	5.4
	70b-instruct	16.8	16.8	54.0	14.0	16.4	16.4	16.8	16.8	10.7	10.7	14.1	14.1

5.2 Error Analysis

The interpretable nature of generating symbolic world models can be utilized for a deeper manual analysis of the failure modes. We select the results from claude-3.5-sonnet under the few-shot setting for manual error analysis. Errors are categorized into syntax and semantic errors, where syntax errors occur when the generated domain cannot be validated ($EXEC. = 0$), and semantic errors arise when the generated world model does not align with action dynamics or fails to follow the natural language description. The distribution for each error type and detailed explanations are presented in Appendix C.

Syntax Errors Figure 5 (Left) shows the distribution of syntax errors during correction. Common errors like UndefinedConstant and IncorrectParentheses decrease over correction steps, indicating improvements in syntax validation, though errors like UndefinedDomainName and UndefinedType persist.

Semantic Error Figure 5 (Right) illustrates the distribution of semantic errors. Semantic errors are categorized into four types: (i) DisobeyDescription involves direct violations of descriptions. (ii) IncompleteModeling, where the world model lacks necessary components. (iii) RedundantSpecifications refers to su-

perfluous preconditions or effects; and (iv) SurfaceDivergence involves surface-level variations that preserve semantic equivalence to gold domain. In addition, since a domain may encompass various action dynamics, different error types can occur simultaneously. For instance, nearly 10% of cases exhibited both IncompleteModeling and RedundantSpecifications concurrently.

6 Exploration

In addition to the zero-shot CoT evaluation in Section 4.2, we further evaluate the models on TEXT2WORLD with five different strategies: (1) *Test-time Scaling*; (2) *In-Context Learning*; (3) *Fine-tuning*; (4) *Agent Training*; (5) *Inference with Concrete Description*.

6.1 Test-time Scaling

Recently, test-time scaling has demonstrated remarkable potential (OpenAI, 2024; DeepSeek-AI et al., 2025). We use the error information from the syntax parser as feedback and assess whether increasing the test-time compute budget can enhance the LLM’s performance. As shown in Figure 6, the model exhibits consistent improvement with increased test-time computation. More advanced test-time scaling strategies may serve as a viable approach to enhancing the model’s world modeling ability (Chen et al., 2025).

Table 2: The experimental results of models under different settings: (1) In-context learning (§6.2); (2) Fine-tuning, and fine-tuning with LoRA (Hu et al., 2021) (§6.3); (3) Agent training (§6.4).

Model Family	EXEC.		SIM.		F1 _{PRED}		F1 _{PARAM}		F1 _{PRECOND}		F1 _{EFF}	
	EC ₀	EC ₃										
<i>In-Context Learning</i>												
CLAUDE-3.5-SONNET <i>w.</i> 2-SHOT	45.5 78.2 _{+32.7}	64.4 88.1 _{+23.7}	73.2 83.9 _{+10.7}	66.8 82.3 _{+15.5}	45.5 77.0 _{+31.5}	62.5 86.1 _{+23.6}	41.5 75.2 _{+33.7}	48.8 82.1 _{+33.3}	37.4 65.6 _{+28.2}	44.0 71.3 _{+27.3}	38.4 67.2 _{+28.8}	45.0 73.4 _{+28.4}
DEEPSEEK-R1 <i>w.</i> 2-SHOT	72.3 69.3 _{-3.0}	89.1 90.1 _{+1.0}	84.3 83.8 _{-0.5}	84.0 83.5 _{-0.5}	71.7 68.4 _{-3.3}	86.7 87.7 _{+1.0}	64.0 64.6 _{+0.6}	76.3 79.1 _{+2.8}	57.6 56.0 _{-1.6}	65.0 66.9 _{+1.9}	58.8 57.6 _{-1.2}	67.3 68.9 _{+1.6}
GPT-4O-MINI <i>w.</i> 2-SHOT	48.5 40.6 _{-7.9}	72.3 69.3 ₋₃	82.6 82.9 _{+0.3}	82.2 82.4 _{+0.2}	48.1 40.3 _{-7.8}	70.1 67.2 _{-2.9}	47.1 40.1 ₋₇	67.3 67.0 _{-0.3}	34.9 31.6 _{-3.3}	47.5 49.3 _{+1.8}	38.2 32.5 _{-5.7}	52.7 54.8 _{+2.1}
<i>Fine-tuning (FT)</i>												
LLAMA-3.1-8B <i>w.</i> FT	0.0 52.5 _{+52.5}	0.0 68.3 _{+68.3}	74.3 80.8 _{+6.5}	74.9 80.6 _{+5.7}	0.0 51.4 _{+51.4}	0.0 65.4 _{+65.4}	0.0 48.5 _{+48.5}	0.0 60.6 _{+60.6}	0.0 31.5 _{+31.5}	0.0 38.1 _{+38.1}	0.0 32.4 _{+32.4}	0.0 40.2 _{+40.2}
LLAMA-3.1-70B <i>w.</i> LORA	0.0 48.5 _{+48.5}	0.0 70.3 _{+70.3}	83.6 83.8 _{+0.2}	79.2 82.3 _{+3.1}	0.0 47.9 _{+47.9}	0.0 68.5 _{+68.5}	0.0 48.5 _{+48.5}	0.0 66.4 _{+66.4}	0.0 39.9 _{+39.9}	0.0 52.8 _{+52.8}	0.0 40.6 _{+40.6}	0.0 52.1 _{+52.1}
<i>Agent Training (AT)</i>												
LLAMA-2-70B <i>w.</i> AT	0.0 7.9 _{+7.9}	0.0 9.9 _{+9.9}	48.7 65.6 _{+16.9}	48.6 47.9 _{-0.7}	0.0 7.3 _{+7.3}	0.0 8.8 _{+8.8}	0.0 7.3 _{+7.3}	0.0 9.1 _{+9.1}	0.0 6.1 _{+6.1}	0.0 6.5 _{+6.5}	0.0 5.7 _{+5.7}	0.0 6.1 _{+6.1}

486 detailed experimental results in Appendix D.1.

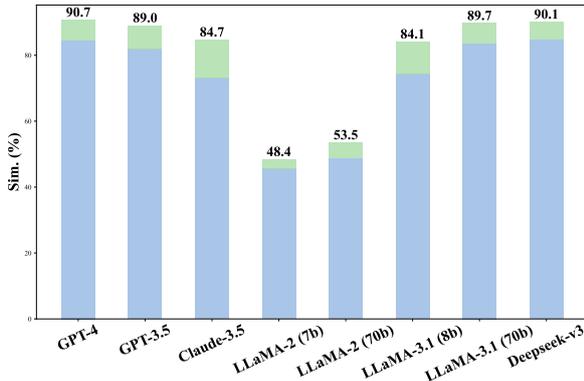


Figure 7: Comparison of model performance on abstract versus concrete domain descriptions, showing the base score for abstract descriptions (blue) and the improvement gained from concrete descriptions (green).

487 7 Related Work

488 Neural world modeling is a long-standing research
489 topic with widespread applications across various
490 fields, including reinforcement learning (Ha and
491 Schmidhuber, 2018b,a), robotics (Wu et al., 2023),
492 and autonomous driving (Guan et al., 2024), among
493 others. In recent years, LLMs trained on massive
494 datasets have demonstrated zero-shot capabilities
495 across a variety of tasks, including planning (Zhao
496 et al., 2023; Qin et al., 2024; Huang et al., 2022;
497 Hu et al., 2024a), robotics (Mu et al., 2024; Chen
498 et al., 2024a), analog design (Lai et al., 2024), and
499 more. Preliminary studies propose directly using
500 LLMs as world models (Hao et al., 2023; Wang
501 et al., 2024, 2023b; Li et al., 2022), by taking the
502 state and action as input and predicting the next

state, but the unreliability and limited interpretability of LLM outputs can lead to accumulating errors. Moreover, some studies have shown that autoregressive models perform poorly in predicting action effects (Banerjee et al., 2020; Luo et al., 2023). Tree-planner (Hu et al., 2023) instead proposes to constructing the possible action space using LLMs before executing. Another line of work focuses on leveraging LLMs to construct symbolic world models (Oswald et al., 2024; Silver et al., 2024; Smirnov et al., 2024; Zhu et al., 2024; Wang et al., 2023a; Wong et al., 2023; Vafa et al., 2024). For example, Guan et al. (2023) uses LLMs to generate a PDDL domain model and relies on human feedback to correct errors. AgentGen (Hu et al., 2024b) synthesizes diverse PDDL domains, aiming to create high-quality planning data. Xie et al. (2024) propose to finetune LLMs for predicting precondition and effect of actions. Despite the growing interest in this research direction, there is currently a lack of a comprehensive benchmark in this area.

524 8 Conclusion

525 We present TEXT2WORLD, a novel benchmark
526 consisting of hundreds of domains designed to
527 evaluate the world modeling capabilities of large
528 language models (LLMs). Developed through a
529 meticulous and thorough process, TEXT2WORLD
530 provides a robust foundation for analysis. Additionally,
531 we conducted an extensive evaluation involving
532 16 different LLMs from 9 model families based
533 on TEXT2WORLD. We hope that TEXT2WORLD
534 will inspire future research in leveraging LLMs as
535 world models.

9 Ethical Considerations

Data Access. We collected the TEXT2WORLD data from open-source repositories and ensured that these repositories are available for academic research in accordance with our commitment to ethical data use.

Participant Recruitment. We recruited graduate students as annotators and required all participants to achieve an IELTS score of 6 or above. To mitigate potential biases stemming from participants' geographical backgrounds, we minimized national differences in the dataset by focusing on human commonsense. All annotators provided informed consent and were compensated above the local minimum wage—\$10 per hour for standard annotators and \$20 per hour for senior annotators.

Potential Risk. After careful examination, we confirmed that our dataset does not contain any personal data (e.g., names, contacting information), and our data collection procedures adhere to ethical guidelines.

10 Limitation

Due to the limited number of available domains online, we did not construct a large-scale training set. Future work should focus on expanding the dataset by incorporating additional data sources, such as synthesized data (Hu et al., 2024b), to cover a broader range of domains. Furthermore, although we conducted regular inspections to minimize the introduction of subjectivity into the dataset, the unavoidable influence of human subjectivity during manual annotation may introduce potential biases.

568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624

References

Anthropic. [Introducing claude 3.5 sonnet](#).

Pratyay Banerjee, Chitta Baral, Man Luo, Arindam Mitra, Kuntal Pal, Tran C Son, and Neeraj Varshney. 2020. Can transformers reason about effects of actions? *arXiv preprint arXiv:2012.09938*.

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.

Junting Chen, Checheng Yu, Xunzhe Zhou, Tianqi Xu, Yao Mu, Mengkang Hu, Wenqi Shao, Yikai Wang, Guohao Li, and Lin Shao. 2024a. [EMOS: Embodiment-aware heterogeneous Multi-robot Operating System with llm agents](#). *Preprint*, arXiv:2410.22662.

Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiaqi Wang, Mengkang Hu, Zhi Chen, Wanxiang Che, and Ting Liu. 2025. Ecm: A unified electronic circuit model for explaining the emergence of in-context learning and chain-of-thought in large language model. *arXiv preprint arXiv:2502.03325*.

Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. 2024b. [M³CoT: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8199–8221, Bangkok, Thailand. Association for Computational Linguistics.

Kenneth James Williams Craik. 1967. *The nature of explanation*, volume 445. CUP Archive.

Nicola Dainese, Matteo Merler, Minttu Alakuijala, and Pekka Marttinen. 2024. Generating code world models with large language models guided by monte carlo tree search. *arXiv preprint arXiv:2405.15383*.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li,

Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanxia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.

Ellen R Girden. 1992. *ANOVA: Repeated measures*. 84. Sage. 659
660

Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. 2023. Leveraging pre-trained large language models to construct and utilize world models for model-based task planning. *Advances in Neural Information Processing Systems*, 36:79081–79094. 661
662
663
664
665
666

Yanchen Guan, Haicheng Liao, Zhenning Li, Jia Hu, Runze Yuan, Yunjian Li, Guohui Zhang, and Chengzhong Xu. 2024. World models for autonomous driving: An initial survey. *IEEE Transactions on Intelligent Vehicles*. 667
668
669
670
671

David Ha and Jürgen Schmidhuber. 2018a. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31. 672
673
674

David Ha and Jürgen Schmidhuber. 2018b. World models. *arXiv preprint arXiv:1803.10122*. 675
676

Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*. 677
678
679
680

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 681
682

683	and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	
684		
685		
686	Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu,	
687	Wenqi Shao, and Ping Luo. 2024a. Hiagent: Hierarchical working memory management for solving	
688	long-horizon agent tasks with large language model.	
689	<i>arXiv preprint arXiv:2408.09559</i> .	
690		
691	Mengkang Hu, Yao Mu, Xinmiao Yu, Mingyu Ding,	
692	Shiguang Wu, Wenqi Shao, Qiguang Chen, Bin	
693	Wang, Yu Qiao, and Ping Luo. 2023. Tree-planner:	
694	Efficient close-loop task planning with large language	
695	models. <i>arXiv preprint arXiv:2310.08582</i> .	
696	Mengkang Hu, Pu Zhao, Can Xu, Qingfeng Sun, Jian-	
697	guang Lou, Qingwei Lin, Ping Luo, Saravan Rajmo-	
698	han, and Dongmei Zhang. 2024b. Agentgen: Enhanc-	
699	ing planning abilities for large language model based	
700	agent via environment and task generation. <i>arXiv</i>	
701	<i>preprint arXiv:2408.00764</i> .	
702	Wenlong Huang, Pieter Abbeel, Deepak Pathak, and	
703	Igor Mordatch. 2022. Language models as zero-shot	
704	planners: Extracting actionable knowledge for em-	
705	bodied agents. In <i>International conference on ma-</i>	
706	<i>chine learning</i> , pages 9118–9147. PMLR.	
707	Yao Lai, Sungyoung Lee, Guojin Chen, Souradip Pod-	
708	dar, Mengkang Hu, David Z Pan, and Ping Luo. 2024.	
709	Analogcoder: Analog circuit design via training-free	
710	code generation. <i>arXiv preprint arXiv:2405.14918</i> .	
711	J. Richard Landis and Gary G. Koch. 1977. The mea-	
712	surement of observer agreement for categorical data.	
713	<i>Biometrics</i> , 33(1):159–174.	
714	Yann LeCun. 2022. A path towards autonomous ma-	
715	chine intelligence version 0.9. 2, 2022-06-27. <i>Open</i>	
716	<i>Review</i> , 62(1):1–62.	
717	Kenneth Li, Aspen K Hopkins, David Bau, Fernanda	
718	Viégas, Hanspeter Pfister, and Martin Wattenberg.	
719	2022. Emergent world representations: Exploring a	
720	sequence model trained on a synthetic task. <i>arXiv</i>	
721	<i>preprint arXiv:2210.13382</i> .	
722	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	
723	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	
724	Deng, Chenyu Zhang, Chong Ruan, et al. 2024.	
725	Deepseek-v3 technical report. <i>arXiv preprint</i>	
726	<i>arXiv:2412.19437</i> .	
727	Man Luo, Shrinidhi Kumbhar, Mihir Parmar, Neeraj	
728	Varshney, Pratyay Banerjee, Somak Aditya, Chitta	
729	Baral, et al. 2023. Towards logigluue: A brief sur-	
730	vey and a benchmark for analyzing logical reason-	
731	ing capabilities of language models. <i>arXiv preprint</i>	
732	<i>arXiv:2310.00836</i> .	
733	Meta AI. Introducing llama 3.1: Our most capable	
734	models to date.	
735	Meta AI. 2024. Introducing meta Llama 3: The most	
736	capable openly available LLM to date. Accessed:	
737	2024-04-18.	
	Yao Mu, Junting Chen, Qinglong Zhang, Shoufa Chen,	738
	Qiaojun Yu, Chongjian Ge, Runjian Chen, Zhix-	739
	uan Liang, Mengkang Hu, Chaofan Tao, et al.	740
	2024. Robocodex: Multimodal code generation	741
	for robotic behavior synthesis. <i>arXiv preprint</i>	742
	<i>arXiv:2402.16117</i> .	743
	OpenAI. 2022. Openai: Introducing chatgpt.	744
	OpenAI. 2023. Gpt-4 technical report. <i>Preprint,</i>	745
	<i>arXiv:2303.08774</i> .	746
	OpenAI. 2024. Learning to reason with llms.	747
	OpenAI. 2025. Openai o3-mini.	748
	James Oswald, Kavitha Srinivas, Harsha Kokel, Junkyu	749
	Lee, Michael Katz, and Shirin Sohrabi. 2024. Large	750
	language models as planning domain generators. In	751
	<i>Proceedings of the International Conference on Au-</i>	752
	<i>tomated Planning and Scheduling</i> , volume 34, pages	753
	423–431.	754
	Libo Qin, Qiguang Chen, Xiachong Feng, Yang Wu,	755
	Yongheng Zhang, Yinghui Li, Min Li, Wanxiang Che,	756
	and Philip S Yu. 2024. Large language models meet	757
	nlp: A survey. <i>arXiv preprint arXiv:2405.12819</i> .	758
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	759
	Dario Amodei, Ilya Sutskever, et al. 2019. Language	760
	models are unsupervised multitask learners. <i>OpenAI</i>	761
	<i>blog</i> , 1(8):9.	762
	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten	763
	Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi,	764
	Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023.	765
	Code llama: Open foundation models for code. <i>arXiv</i>	766
	<i>preprint arXiv:2308.12950</i> .	767
	Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B	768
	Tenenbaum, Leslie Kaelbling, and Michael Katz.	769
	2024. Generalized planning in pddl domains with	770
	pretrained large language models. In <i>Proceedings</i>	771
	<i>of the AAAI Conference on Artificial Intelligence</i> ,	772
	volume 38, pages 20256–20264.	773
	Pavel Smirnov, Frank Joublin, Antonello Ceravola, and	774
	Michael Gienger. 2024. Generating consistent pddl	775
	domains with large language models. <i>arXiv preprint</i>	776
	<i>arXiv:2404.07751</i> .	777
	Hao Tang, Darren Key, and Kevin Ellis. 2024. World-	778
	coder, a model-based llm agent: Building world mod-	779
	els by writing code and interacting with the environ-	780
	ment. <i>arXiv preprint arXiv:2402.12275</i> .	781
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	782
	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	783
	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	784
	Bhosale, et al. 2023. Llama 2: Open founda-	785
	tion and fine-tuned chat models. <i>arXiv preprint</i>	786
	<i>arXiv:2307.09288</i> .	787
	Keyon Vafa, Justin Y Chen, Jon Kleinberg, Sendhil	788
	Mullainathan, and Ashesh Rambachan. 2024. Evalu-	789
	ating the world model implicit in a generative model.	790
	<i>arXiv preprint arXiv:2406.03689</i> .	791

792 Ruoyao Wang, Graham Todd, Ziang Xiao, Xingdi Yuan,
793 Marc-Alexandre Côté, Peter Clark, and Peter Jansen.
794 2024. Can language models serve as text-based
795 world simulators? *arXiv preprint arXiv:2406.06485*.

796 Ruoyao Wang, Graham Todd, Eric Yuan, Ziang Xiao,
797 Marc-Alexandre Côté, and Peter Jansen. 2023a.
798 Bytesized32: A corpus and challenge task for gener-
799 ating task-specific world models expressed as text
800 games. *arXiv preprint arXiv:2305.14879*.

801 Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai,
802 Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P
803 Xing, and Zhiting Hu. 2023b. Promptagent:
804 Strategic planning with language models enables
805 expert-level prompt optimization. *arXiv preprint*
806 *arXiv:2310.16427*.

807 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
808 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
809 et al. 2022. Chain-of-thought prompting elicits reason-
810 ing in large language models. *Advances in Neural*
811 *Information Processing Systems*, 35:24824–24837.

812 Lionel Wong, Gabriel Grand, Alexander K Lew, Noah D
813 Goodman, Vikash K Mansinghka, Jacob Andreas,
814 and Joshua B Tenenbaum. 2023. From word mod-
815 els to world models: Translating from natural lan-
816 guage to the probabilistic language of thought. *arXiv*
817 *preprint arXiv:2306.12672*.

818 Philipp Wu, Alejandro Escontrela, Danijar Hafner,
819 Pieter Abbeel, and Ken Goldberg. 2023. Day-
820 dreamer: World models for physical robot learning.
821 In *Conference on robot learning*, pages 2226–2240.
822 PMLR.

823 Kaige Xie, Ian Yang, John Gunerli, and Mark Riedl.
824 2024. Making large language models into world mod-
825 els with precondition and effect knowledge. *arXiv*
826 *preprint arXiv:2409.12278*.

827 Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao
828 Liu, Yuxiao Dong, and Jie Tang. 2023. Agenttuning:
829 Enabling generalized agent abilities for llms. *arXiv*
830 *preprint arXiv:2310.12823*.

831 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
832 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
833 Zhang, Junjie Zhang, Zican Dong, et al. 2023. A
834 survey of large language models. *arXiv preprint*
835 *arXiv:2303.18223*.

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

A Benchmark Construction

840

A.1 Example

841

A.1.1 Domain Example

842

```
(define (domain grid)
  (:requirements :strips)
  (:predicates (conn ?x ?y) (key-shape ?k ?s) (lock-shape ?x ?s)
    (at ?r ?x) (at-robot ?x) (place ?p) (key ?k) (shape ?s)
    (locked ?x) (holding ?k) (open ?x) (arm-empty ))

  (:action unlock
    :parameters (?curpos ?lockpos ?key ?shape)
    :precondition (and (place ?curpos) (place ?lockpos) (key ?key)
      (shape ?shape) (conn ?curpos ?lockpos)
      (key-shape ?key ?shape) (lock-shape ?lockpos ?shape)
      (at-robot ?curpos) (locked ?lockpos) (holding ?key))
    :effect (and (open ?lockpos) (not (locked ?lockpos))))

  (:action move
    :parameters (?curpos ?nextpos)
    :precondition (and (place ?curpos) (place ?nextpos) (at-robot ?curpos)
      (conn ?curpos ?nextpos) (open ?nextpos))
    :effect (and (at-robot ?nextpos) (not (at-robot ?curpos))))

  (:action pickup
    :parameters (?curpos ?key)
    :precondition (and (place ?curpos) (key ?key) (at-robot ?curpos)
      (at ?key ?curpos) (arm-empty ))
    :effect (and (holding ?key) (not (at ?key ?curpos)) (not (arm-empty ))))

  (:action pickup-and-loose
    :parameters (?curpos ?newkey ?oldkey)
    :precondition (and (place ?curpos) (key ?newkey) (key ?oldkey)
      (at-robot ?curpos) (holding ?oldkey)
      (at ?newkey ?curpos))
    :effect (and (holding ?newkey) (at ?oldkey ?curpos)
      (not (holding ?oldkey)) (not (at ?newkey ?curpos))))

  (:action putdown
    :parameters (?curpos ?key)
    :precondition (and (place ?curpos) (key ?key) (at-robot ?curpos)
      (holding ?key))
    :effect (and (arm-empty ) (at ?key ?curpos) (not (holding ?key))))
)
```

843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882

Listing 1: Grid PDDL

A.1.2 Abstract Description

883

General. This domain models a robot navigating a grid environment with the objective of unlocking doors and moving through the grid. The robot can carry keys that match the shape of locks to unlock doors. The environment includes places, keys with specific shapes, and doors (locks) with corresponding shapes that need to be unlocked.

884
885
886
887

Predicates. The following predicates are used in the domain:

888

- (conn ?x ?y): Indicates a connection between two places ?x and ?y, allowing movement between them.
- (key-shape ?k ?s): Indicates that key ?k has shape ?s.
- (lock-shape ?x ?s): Indicates that lock (or door) at place ?x has shape ?s.
- (at ?r ?x): Indicates that key ?r is at place ?x.
- (at-robot ?x): Indicates that the robot is at place ?x.

889
890
891
892
893
894

- 895 • (place ?p): Indicates that ?p is a place in the grid.
- 896 • (key ?k): Indicates that ?k is a key.
- 897 • (shape ?s): Indicates that ?s is a shape.
- 898 • (locked ?x): Indicates that the place ?x is locked.
- 899 • (holding ?k): Indicates that the robot is holding key ?k.
- 900 • (open ?x): Indicates that the place ?x is open.
- 901 • (arm-empty): Indicates that the robot's arm is empty.

902 **Actions.** The following actions are available in the domain:

- 903 • unlock <?curpos> <?lockpos> <?key> <?shape>: Allows the robot to unlock a door at place
- 904 <?lockpos> using a key of a specific shape.
- 905 • move <?curpos> <?nextpos>: Allows the robot to move from place <?curpos> to place
- 906 <?nextpos>.
- 907 • pickup <?curpos> <?key>: Allows the robot to pick up a key at its current location.
- 908 • pickup-and-loose <?curpos> <?newkey> <?oldkey>: Allows the robot to pick up a new key
- 909 while dropping the one it was holding.
- 910 • putdown <?curpos> <?key>: Allows the robot to put down a key it is holding.

911 A.1.3 Concrete Description

912 **General.** This domain models a robot navigating a grid environment with the objective of unlocking

913 doors and moving through the grid. The robot can carry keys that match the shape of locks to unlock

914 doors. The environment includes places, keys with specific shapes, and doors (locks) with corresponding

915 shapes that need to be unlocked.

916 **Predicates.** The following predicates are used in the domain:

- 917 • (conn ?x ?y): Indicates a connection between two places ?x and ?y, allowing movement between
- 918 them.
- 919 • (key-shape ?k ?s): Indicates that key ?k has shape ?s.
- 920 • (lock-shape ?x ?s): Indicates that lock (or door) at place ?x has shape ?s.
- 921 • (at ?r ?x): Indicates that key ?r is at place ?x.
- 922 • (at-robot ?x): Indicates that the robot is at place ?x.
- 923 • (place ?p): Indicates that ?p is a place in the grid.
- 924 • (key ?k): Indicates that ?k is a key.
- 925 • (shape ?s): Indicates that ?s is a shape.
- 926 • (locked ?x): Indicates that the place ?x is locked.
- 927 • (holding ?k): Indicates that the robot is holding key ?k.
- 928 • (open ?x): Indicates that the place ?x is open.
- 929 • (arm-empty): Indicates that the robot's arm is empty.

Actions. The following actions are available in the domain: 930

- unlock <?curpos> <?lockpos> <?key> <?shape>: Allows the robot to unlock a door at place <?lockpos> using a key of a specific shape if the robot is at place <?curpos>, the key matches the lock's shape, the robot is holding the key, there is a connection between <?curpos> and <?lockpos>, and the destination is locked. After the action, the lock is no longer locked. 931-934
- move <?curpos> <?nextpos>: Allows the robot to move from place <?curpos> to place <?nextpos> if there is a connection between them and the destination is open. After the move, the robot is no longer at the original place. 935-937
- pickup <?curpos> <?key>: Allows the robot to pick up a key at its current location if the robot's arm is empty and it is at the same place as the new key. After the action, the robot is holding the key, and the key is no longer at that location. 938-940
- pickup-and-loose <?curpos> <?newkey> <?oldkey>: Allows the robot to pick up a new key while dropping the one it was holding if it is at the same place as the new key. After the action, the robot is holding the new key, and the old key is at the robot's current location. 941-943
- putdown <?curpos> <?key>: Allows the robot to put down a key it is holding if it is at a specific place. After the action, the robot's arm is empty, and the key is at that location. 944-945

A.2 Preliminary Experiment 946

The experimental results show that LLM's effectiveness in detecting PDDL semantic errors is limited, with an accuracy of 55.0%, a precision of 56.2%, a recall rate of 45.0%, an F1 score of 50.0%, and a ROC AUC of 55.0. ROC AUC indicates that the model is close to random performance, making it difficult to reliably distinguish between correct and incorrect PDDL domains. Below is the prompt used for LLMs to detect semantic errors in generated PDDL domains: 947-951

```
You are an expert in automated planning systems and PDDL semantics. Your task is to evaluate whether the LLM are physically accurate models of the world or whether they don't make sense by detecting semantic errors in generated PDDL domain. You need carefully analyze the following PDDL domain by comparing it to the pddl domain description, evaluate whether the generated pddl domain contains SEMANTIC ERRORS in these key aspects: 952-958
1. Predicates consistency. 959
2. Action parameters validity. 960
3. Action preconditions completeness. 961
4. Action effects logical consistency. 962
5. Consistency with the description. 963-964

An example of semantic error would be: 965
1. Missing precondition constraints (e.g. executing "unlock-door" without holding a key). 966-967
2. Contradictory effects (e.g. both adding and deleting the same predicate). 968
3. Incorrect predicate arguments (e.g. reversed parameter order). 969-970

Output Format: 971
{ 972
  "evaluation": "yes/no", 973
  "error_type": "[MissingPrecond|IncorrectEffect|MissingPredicate|...]", 974
  "confidence": "high/medium/low", 975
  "evidence": "<specific code segment>", 976
  "justification": "<short justification>" 977
} 978-979

PDDL Description: 980
{PDDL_DESCRIPTION} 981-982

Generated PDDL: 983
{PDDL_DOMAIN} 984
```

985
986
987

A.3 More Details on Data Analysis

Figure 8 shows the co-occurrence of PDDL requirements across domains, highlighting that `:typing` and `:strips` are the most prevalent features.

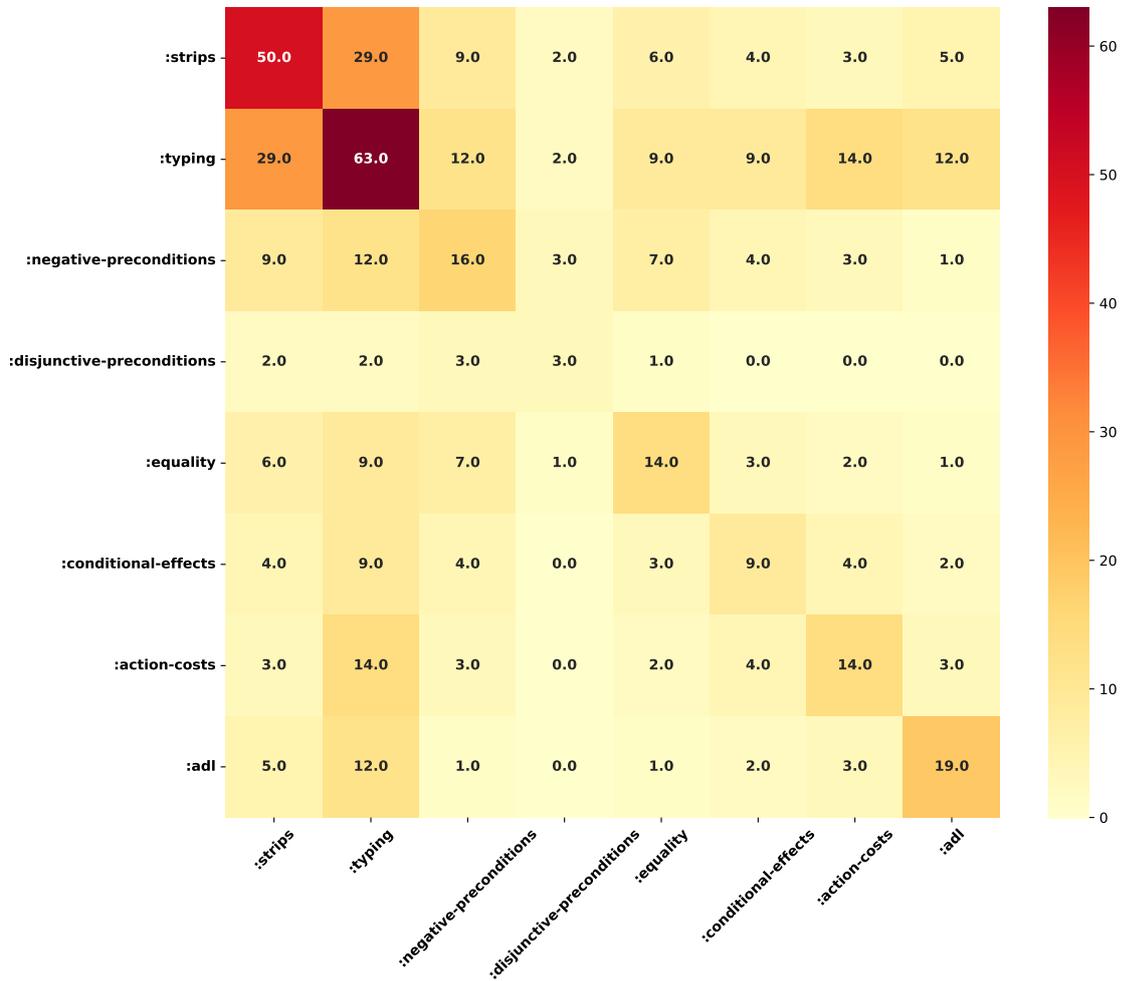


Figure 8: The co-occurrence matrix of requirements of TEXT2WORLD.

988
989
990
991
992
993
994
995
996
997
998
999
1000
1001

B More Details on Experiments

B.1 Evaluation Metrics

Levenshtein Ratio. The Levenshtein Ratio is a value between 0 and 1 that quantifies the similarity between two strings, such as a predicted PDDL domain and a golden PDDL domain. It is derived from the Levenshtein distance, which calculates the minimum number of character-level operations—insertions, deletions, or substitutions—needed to convert one string into the other. The ratio is then computed by dividing the Levenshtein distance by the length of the longer string, providing a measure of how closely the two strings match, where a value closer to 1 indicates high similarity and a value closer to 0 indicates significant differences.

Component-wise F1 Scores. The F1 score is mainly used to measure the similarity between the predicted PDDL domain and the golden PDDL domain, specifically including predicate F1 and action F1. The range of this score is from 0 to 1, which is the harmonic mean of precision and recall.

B.2 Prompt Examples

B.2.1 Error Correction

I would like you to serve as an expert in PDDL, assisting me in correcting erroneous PDDL code. I will provide you with the incorrect PDDL along with the error messages returned by the system. You should output your thought process firstly. You MUST enclose the COMPLETE corrected PDDL within ``pddl``.

Here are some hints to help you debug the pddl domain file:

1. You should start by checking if all the essential domain constructs or domain definition constructs are present. Commonly included domains comprise:
 - a. :domain declaration to name the domain.
 - b. :requirements to specify the PDDL features used in the domain.
 - c. :types to define any object types for categorizing entities in the planning problem.
 - d. :constants (if necessary) to declare any objects that remain static throughout the planning problems.
 - e. :predicates to define the properties and relations between objects that can change over time.
 - f. :functions (if necessary) to define numeric functions for more complex domains.
 - g. :action definitions for each action that agents can perform, including parameters, preconditions, and effects.
2. You need to check the number of parameters of each actions.
3. Having :typing in the domain indicates that it uses a hierarchy to organize objects. Therefore, it's crucial to clearly list all object types related to the planning task in a :types section.
4. '-' should not appear in :types.

Round 0

Incorrect PDDL:

```
(:action clean-up
  :parameters (?robot - robot ?robotTile - tile ?tileToBeCleaned - tile)
  :precondition (and
    (robot-at ?robot ?robotTile)
    (up ?tileToBeCleaned ?robotTile)
    (clear ?tileToBeCleaned)
    (not (cleaned ?tileToBeCleaned)))
  )
  :effect (and
    (cleaned ?tileToBeCleaned)
  )
)

(:action clean-down
  :parameters (?robot - robot ?robotTile - tile ?tileToBeCleaned - tile)
  :precondition (and
    (robot-at ?robot ?robotTile)
    (down ?tileToBeCleaned ?robotTile)
    (clear ?tileToBeCleaned)
    (not (cleaned ?tileToBeCleaned)))
  )
  :effect (and
    (cleaned ?tileToBeCleaned)
  )
)

(:action up
  :parameters (?robot - robot ?robotTile - tile ?moveToNextTile - tile)
  :precondition (and
    (robot-at ?robot ?robotTile)
    (up ?moveToNextTile ?robotTile)
    (clear ?moveToNextTile)
  )
  :effect (and
    (not (robot-at ?robot ?robotTile))
    (robot-at ?robot ?moveToNextTile)
  )
)

(:action down
  :parameters (?robot - robot ?robotTile - tile ?moveToNextTile - tile)
```

```

1072     :precondition (and
1073         (robot-at ?robot ?robotTile)
1074         (down ?moveToNextTile ?robotTile)
1075         (clear ?moveToNextTile)
1076     )
1077     :effect (and
1078         (not (robot-at ?robot ?robotTile))
1079         (robot-at ?robot ?moveToNextTile)
1080     )
1081 )
1082
1083 (:action right
1084     :parameters (?robot - robot ?robotTile - tile ?moveToNextTile - tile)
1085     :precondition (and
1086         (robot-at ?robot ?robotTile)
1087         (right ?moveToNextTile ?robotTile)
1088         (clear ?moveToNextTile)
1089     )
1090     :effect (and
1091         (not (robot-at ?robot ?robotTile))
1092         (robot-at ?robot ?moveToNextTile)
1093     )
1094 )
1095
1096 (:action left
1097     :parameters (?robot - robot ?robotTile - tile ?moveToNextTile - tile)
1098     :precondition (and
1099         (robot-at ?robot ?robotTile)
1100         (left ?moveToNextTile ?robotTile)
1101         (clear ?moveToNextTile)
1102     )
1103     :effect (and
1104         (not (robot-at ?robot ?robotTile))
1105         (robot-at ?robot ?moveToNextTile)
1106     )
1107 )
1108 Error Information:
1109 ParsingError: line 1:1 mismatched input ':action' expecting 'define'
1110 Corrected PDDL:

```

1111 B.2.2 Zero-Shot Prompt

1112 You are tasked with converting a given Planning Domain Definition Language (PDDL)
1113 domain description into its corresponding formal PDDL domain. The description
1114 will outline the essential components of the domains.
1115 Your output should be a well-structured PDDL domain that accurately represents the
1116 given description, adhering to the syntax and semantics of PDDL.
1117 Your output pddl domain must be enclosed in ``pddl``.

1118
1119 You need to generate the corresponding domain pddl for the following description.

1120
1121 PDDL Domain Description:

1122 ### General

1123 This domain is designed for a robot tasked with cleaning floor tiles. The robot can
1124 move in four directions (up, down, right, left) relative to its current position
1125 on a grid of tiles. The goal is to clean all the specified tiles by moving to
1126 them and performing a cleaning action.

1127
1128 ### Types

- 1129 - **robot**: Represents the robot that performs the cleaning.
- 1130 - **tile**: Represents the individual tiles on the floor that may need to be cleaned
- 1131 .

1132
1133 ### Predicates

- 1134 - **(robot-at ?robot - robot ?robotTile - tile)**: Indicates that the robot is
1135 currently at a specific tile.
- 1136 - **(up ?tileAbove - tile ?tileBelow - tile)**: Indicates that one tile is directly
1137 above another.
- 1138 - **(down ?tileBelow - tile ?tileAbove - tile)**: Indicates that one tile is
1139 directly below another.

```

- **right ?tileOnRight - tile ?tileOnLeft - tile)**: Indicates that one tile is 1140
  directly to the right of another. 1141
- **left ?tileOnLeft - tile ?tileOnRight - tile)**: Indicates that one tile is 1142
  directly to the left of another. 1143
- **clear ?clearedTile - tile)**: Indicates that a tile is clear and robot can move 1144
  there. 1145
- **cleaned ?cleanedTile - tile)**: Indicates that a tile has been cleaned. 1146
  1147
### Actions 1148
- **clean-up <?robot> <?robotTile> <?tileToBeCleaned>**: Allows the robot (?robot) 1149
  to clean a tile (?tileToBeCleaned) that is directly above its current position 1150
  (?robotTile). 1151
  1152
- **clean-down <?robot> <?robotTile> <?tileToBeCleaned>**: Allows the robot (?robot) 1153
  to clean a tile (?tileToBeCleaned) that is directly below its current position 1154
  (?robotTile). 1155
  1156
- **up <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a tile 1157
  (?moveToNextTile) directly above its current position (?robotTile). 1158
  1159
- **down <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a 1160
  tile (?moveToNextTile) directly below its current position (?robotTile). 1161
  1162
- **right <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a 1163
  tile (?moveToNextTile) directly to the right of its current position (?robotTile 1164
  ). 1165
  1166
- **left <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a 1167
  tile (?moveToNextTile) directly to the left of its current position (?robotTile) 1168
  . 1169
PDDL Domain: 1170
Let's think step by step. 1171

```

B.2.3 Few-Shot Prompt

```

You are tasked with converting a given Planning Domain Definition Language (PDDL) 1173
domain description into its corresponding formal PDDL domain. The description 1174
will outline the essential components of the domains. Your output should be a 1175
well-structured PDDL domain that accurately represents the given description, 1176
adhering to the syntax and semantics of PDDL. 1177
Your output must strictly adhere to the format exemplified below. 1178
Here are some examples: 1179
  1180
Example 0: 1181
## PDDL Domain Description 1182
### General 1183
You are a robot equipped with a gripper mechanism, designed to move and manipulate 1184
balls between different rooms. The domain focuses on the robot's ability to 1185
navigate rooms, pick up balls, and drop them in designated locations. 1186
### Types 1187
- **room**: Represents the different rooms within the environment. 1188
- **ball**: Represents the objects that the robot can pick up and move. 1189
- **gripper**: Represents the robot's mechanism for holding balls. 1190
### Predicates 1191
- **at-robby ?r - room**: Indicates that Robby, the robot, is currently in room ?r 1192
  . 1193
- **at ?b - ball ?r - room**: Indicates that ball ?b is located in room ?r. 1194
- **free ?g - gripper**: Indicates that the gripper ?g is not currently holding 1195
  any ball. 1196
- **carry ?o - ball ?g - gripper**: Indicates that the gripper ?g is carrying ball 1197
  ?o. 1198
### Actions 1199
- **move <?from> <?to>**: Allows Robby to move from one room to another. 1200
- **pick <?obj> <?room> <?gripper>**: Enables Robby to pick up a ball in a room 1201
  using its gripper. 1202
- **drop <?obj> <?room> <?gripper>**: Allows Robby to drop a ball it is carrying 1203
  into a room. 1204
  1205
## PDDL Domain 1206
```pddl 1207

```

```

1208 (define (domain gripper-strips)
1209 (:types
1210 room - object
1211 ball - object
1212 gripper - object
1213)
1214 (:predicates
1215 (at-robby ?r - room)
1216 (at ?b - ball ?r - room)
1217 (free ?g - gripper)
1218 (carry ?o - ball ?g - gripper))
1219 (:action move
1220 :parameters (?from - room ?to - room)
1221 :precondition (and (at-robby ?from))
1222 :effect (and (at-robby ?to)
1223 (not (at-robby ?from))))
1224 (:action pick
1225 :parameters (?obj - ball ?room - room ?gripper - gripper)
1226 :precondition (and
1227 (at ?obj ?room) (at-robby ?room) (free ?gripper))
1228 :effect (and (carry ?obj ?gripper)
1229 (not (at ?obj ?room))
1230 (not (free ?gripper))))
1231 (:action drop
1232 :parameters (?obj - ball ?room - room ?gripper - gripper)
1233 :precondition (and
1234 (carry ?obj ?gripper) (at-robby ?room))
1235 :effect (and (at ?obj ?room)
1236 (free ?gripper)
1237 (not (carry ?obj ?gripper))))))
1238
1239
1240

```

Example 1:

```

1241 ## PDDL Domain Description
1242 ### General
1243 This domain represents a simplified version of the classic "blocks world" problem,
1244 where a robot arm can stack and unstack blocks. The domain includes actions for
1245 picking up blocks from the table, putting down blocks onto the table, stacking
1246 blocks on top of each other, and unstacking them. The goal is to manipulate the
1247 blocks to achieve a specified configuration.
1248 ### Predicates
1249 (clear ?x): Indicates that there is no block on top of block ?x, making it
1250 accessible for stacking or picking up.
1251 (on-table ?x): Indicates that block ?x is directly on the table.
1252 (arm-empty): Indicates that the robot's arm is not holding any block.
1253 (holding ?x): Indicates that the robot's arm is currently holding block ?x.
1254 (on ?x ?y): Indicates that block ?x is directly on top of block ?y.
1255 ### Actions
1256 - **pickup <?ob>**: Picks up an object (?ob) from the table.
1257 - **putdown <?ob>**: Puts down an object (?ob) onto the table.
1258 - **stack <?ob> <?underob>**: Stacks an object (?ob) on top of another object (?
1259 underob), making the robot arm empty.
1260 - **unstack <?ob> <?underob>**: Unstacks an object (?ob) from another object (?
1261 underob), making the robot arm no longer empty.
1262
1263 ## PDDL Domain:
1264 ```pddl
1265 (define (domain blocksworld)
1266 (:requirements :strips)
1267 (:predicates (clear ?x)
1268 (on-table ?x)
1269 (arm-empty)
1270 (holding ?x)
1271 (on ?x ?y))
1272
1273 (:action pickup
1274 :parameters (?ob)
1275 :precondition (and (clear ?ob) (on-table ?ob) (arm-empty))
1276 :effect (and (holding ?ob) (not (clear ?ob)) (not (on-table ?ob))
1277

```

```

 (not (arm-empty))))))
1278
1279
(:action putdown
1280
 :parameters (?ob)
1281
 :precondition (holding ?ob)
1282
 :effect (and (clear ?ob) (arm-empty) (on-table ?ob)
1283
 (not (holding ?ob))))
1284
1285
(:action stack
1286
 :parameters (?ob ?underob)
1287
 :precondition (and (clear ?underob) (holding ?ob))
1288
 :effect (and (arm-empty) (clear ?ob) (on ?ob ?underob)
1289
 (not (clear ?underob)) (not (holding ?ob))))
1290
1291
(:action unstack
1292
 :parameters (?ob ?underob)
1293
 :precondition (and (on ?ob ?underob) (clear ?ob) (arm-empty))
1294
 :effect (and (holding ?ob) (clear ?underob)
1295
 (not (on ?ob ?underob)) (not (clear ?ob)) (not (arm-empty))))
1296
...
1297
1298
You need to generate the corresponding domain pddl for the following description.
1299
1300
PDDL Domain Description
1301
General
1302
This domain is designed for a robot tasked with cleaning floor tiles. The robot can
1303
move in four directions (up, down, right, left) relative to its current position
1304
on a grid of tiles. The goal is to clean all the specified tiles by moving to
1305
them and performing a cleaning action.
1306
1307
Types
1308
- **robot**: Represents the robot that performs the cleaning.
1309
- **tile**: Represents the individual tiles on the floor that may need to be cleaned
1310
.
1311
1312
Predicates
1313
- **(robot-at ?robot - robot ?robotTile - tile)**: Indicates that the robot is
1314
currently at a specific tile.
1315
- **(up ?tileAbove - tile ?tileBelow - tile)**: Indicates that one tile is directly
1316
above another.
1317
- **(down ?tileBelow - tile ?tileAbove - tile)**: Indicates that one tile is
1318
directly below another.
1319
- **(right ?tileOnRight - tile ?tileOnLeft - tile)**: Indicates that one tile is
1320
directly to the right of another.
1321
- **(left ?tileOnLeft - tile ?tileOnRight - tile)**: Indicates that one tile is
1322
directly to the left of another.
1323
- **(clear ?clearedTile - tile)**: Indicates that a tile is clear and robot can move
1324
there.
1325
- **(cleaned ?cleanedTile - tile)**: Indicates that a tile has been cleaned.
1326
1327
Actions
1328
- **clean-up <?robot> <?robotTile> <?tileToBeCleaned>**: Allows the robot (?robot)
1329
to clean a tile (?tileToBeCleaned) that is directly above its current position
1330
(?robotTile).
1331
1332
- **clean-down <?robot> <?robotTile> <?tileToBeCleaned>**: Allows the robot (?robot)
1333
to clean a tile (?tileToBeCleaned) that is directly below its current position
1334
(?robotTile).
1335
1336
- **up <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a tile
1337
(?moveToNextTile) directly above its current position (?robotTile).
1338
1339
- **down <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a
1340
tile (?moveToNextTile) directly below its current position (?robotTile).
1341
1342
- **right <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a
1343
tile (?moveToNextTile) directly to the right of its current position (?robotTile
1344
).
1345
1346
- **left <?robot> <?robotTile> <?moveToNextTile>**: Moves the robot (?robot) to a
1347

```

1348           tile (?moveToNextTile) directly to the left of its current position (?robotTile)  
 1349           .  
 1350 ## PDDL Domain

## 1351 C More Details on Analysis

### 1352 C.1 Overall

Table 3: Distribution of error types of claude-3.5-sonnet on TEXT2WORLD under few-shot setting.

	Proportion (%)	Number
<b>Correct</b>	23.76	24
<b>Syntax Error</b>	11.88	12
<b>Semantic Error</b>	64.36	65
<b>All</b>	100.00	101

1353 The overall distribution for syntax errors and semantic errors is presented in Table 3.

### 1354 C.2 Syntax Error

Table 4: Distribution of Syntax Errors in PDDL Generation (Total Samples: 66, a task may have 1 to 4 samples.)

Syntax Error	Explanation	Proportion (%)
UndefinedDomainName	Missing mandatory (define (domain ...)) declaration in PDDL header	33.33
IncorrectParentheses	Invalid empty/mismatched parentheses	3.03
UndefinedConstant	Reference to undeclared constants in predicates or actions	13.64
MissingRequirements	Absence of required PDDL extension declarations (e.g., :action-costs)	22.73
UndefinedType	Undeclared parent type in hierarchical type definitions	18.18
UnsupportedFeature	Use of parser-incompatible language features (e.g., either types)	3.03
TypeMismatch	Parameter type conflict with declared type constraints	1.52
UndefinedVariable	Undeclared variables in action preconditions/effects	1.52
DuplicateDefinition	Multiple declarations of identical domain elements	3.03

1355 The distribution and detailed explanation of each syntax error type are presented in Table 4.

### 1356 C.3 Semantic Error

Table 5: Distribution of Semantic Errors in PDDL Generation (Total Samples: 91, a task may have multiple semantic errors.)

Semantic Error	Explanation	Proportion (%)
<b>DisobeyDescription</b>	Direct violation of semantic requirements explicitly stated in the task description.	<b>14.29</b>
IncorrectPredicate	Incorrect or missing the declaration of predicates.	6.59
IncorrectAction	Incorrect or missing the declaration of actions.	7.69
<b>IncompleteModeling</b>	Incomplete world modeling compared to basic requirements.	<b>58.24</b>
IncorrectPrecondition	The precondition of the action is deficient or incorrect.	29.67
IncorrectEffect	The effect of the action is deficient or incorrect.	28.57
<b>RedundantSpecifications</b>	Predicted domain includes superfluous preconditions or effects.	<b>17.58</b>
RedundantPrecondition	Predicted domain includes superfluous preconditions.	10.99
RedundantEffect	Predicted domain includes superfluous effects.	6.59
<b>SurfaceDivergence</b>	Surface variations preserving semantic equivalence with ground truth.	<b>9.89</b>

1357 The distribution and detailed explanation of each semantic error type are presented in Table 5.

## D More Experimental Results

1358

### D.1 Experimental Results with Concrete Description

1359

Table 6: Performance comparison of different LLMs on TEXT2WORLD using concrete domain description.  $EC_k$  denotes the setting where models are allowed k correction attempts ( $EC_0$ : zero-shot without correction,  $EC_3$ : with 3 correction attempts).

Model Family	Version	EXEC. $\uparrow$		SIM. $\uparrow$		$F1_{\text{PRED}} \uparrow$		$F1_{\text{PARAM}} \uparrow$		$F1_{\text{PRECOND}} \uparrow$		$F1_{\text{EFF}} \uparrow$	
		$EC_0$	$EC_3$	$EC_0$	$EC_3$	$EC_0$	$EC_3$	$EC_0$	$EC_3$	$EC_0$	$EC_3$	$EC_0$	$EC_3$
GPT-4	gpt-4o	60.4	75.2	90.7	90.3	59.4	71.8	57.1	69.1	55.3	65.1	54.1	65.2
GPT-3.5	turbo-0125	53.5	68.3	89.0	88.7	52.9	66.7	50.3	64.6	45.1	58.0	46.5	59.9
CLAUDE-3.5	sonnet	64.4	84.2	84.7	77.6	64.4	80.7	55.0	67.5	53.3	65.0	53.3	64.8
LLAMA-2	7b-instruct	0.0	0.0	48.4	32.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	70b-instruct	0.0	0.0	53.5	52.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
LLAMA-3.1	8b-instruct	0.0	1.0	84.1	83.2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	70b-instruct	1.0	1.0	89.7	85.4	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
DEEPSEEK	deepseek-v3	58.4	80.2	90.1	89.3	58.1	76.4	56.2	73.5	53.4	66.0	53.5	67.6