

Benchmarking and Improving PDDL Formalization Ability of Large Language Models with Planner-in-the-Loop Feedback

Anonymous ACL submission

Abstract

Planning depends on symbolic specifications that are both executable and verifiable. However, large language models often generate Planning Domain Definition Language (PDDL) problem descriptions that appear syntactically well formed yet fail under strict precondition, effect-consistency, and reachability constraints. Even minor specification errors can render a task unsolvable, motivating benchmarks and learning signals grounded in planner-based verification rather than surface plausibility.

We present NL-PDDL-Bench, a multi-domain benchmark for natural-language-to-PDDL specification construction, with planner-verified executability and object-count difficulty scaling. We also propose a planner-in-the-loop framework that uses validator and planner diagnostics to revise non-executable specifications via localized edits. Building on this toolchain, we present a planner-grounded optimization recipe combining parameter-efficient Low-Rank Adaptation supervised fine-tuning, offline planner-derived preference pairs for Direct Preference Optimization, and inference-time planner-in-the-loop repair, without online planner calls during training. We further provide a unified evaluation suite for parseability, solvability, specification similarity, and outcome-aware plan-level consistency against planner references. Experiments on eight representative model families show higher planner success and plan-level agreement, improved robustness under difficulty scaling and cross-domain transfer. Code and data are available at: <https://anonymous.4open.science/r/NL-PDDL-Bench-BF76>

1 Introduction

Planning is a core capability in artificial intelligence. Given an initial state and a set of goal conditions, a planner could synthesize an executable sequence of actions. This requirement underlies many applications, including robot control, autonomous driving,

and logistics scheduling. Planning is intrinsically difficult because it demands a faithful model of action preconditions and effects, and it must satisfy reachability constraints over long horizons and large action spaces (Bercher et al., 2025; Pozo and Seipp, 2025).

Classical symbolic planners address these requirements through explicit, checkable specifications. PDDL is a widely adopted formalism that represents planning domains and problem instances in a structured, planner-verifiable form (Aeronautiques et al., 1998). When specifications are correct, state transitions and goal satisfaction can be validated by established planners and verification tools, yielding reproducible executability guarantees (Helmert, 2006; Howey et al., 2004). However, constructing such specifications is costly and typically requires expert knowledge, which limits scalability and practical deployment.

Recent large language models (LLMs) offer strong language understanding and structured generation abilities, motivating their use to reduce modeling effort (Valmeekam et al., 2023b). In practice, existing approaches largely follow two paradigms. The first paradigm directly generates action sequences from a task description. The second paradigm produces a symbolic specification that is subsequently solved by a symbolic planner (Huang et al., 2025; Mahdavi et al., 2024). Although both paradigms can produce outputs that appear plausible, reliability under planner verification remains the key obstacle.

For direct plan generation, actions that look reasonable often violate preconditions, break state consistency, or miss necessary constraints, and they fail under established planners (Valmeekam et al., 2023a; Kokel and Katz, 2025). More broadly, current LLMs remain far behind strong symbolic planners in terms of correctness and reliability on classical planning tasks, especially under strict reachability and long-horizon constraints. For symbolic

specification generation, outputs may be syntactically well formed yet still unsolvable due to subtle modeling errors such as missing objects, incorrect predicate usage, or incomplete initial facts (Silver et al., 2024). These failure modes show that planning requires consistency, executability, and verifiability that cannot be ensured by surface plausibility alone.

We address this challenge by releasing NL-PDDL-Bench and a unified planner-grounded pipeline that supports executable specification construction, optimization, and evaluation at scale. It couples large-scale planner-verified NL-PDDL pairs with a planner-in-the-loop feedback mechanism and planner-grounded training signals, enabling reproducible analysis and robust performance comparisons across domains, difficulty levels, and transfer settings.

Contributions:

- We propose a planner-in-the-loop feedback framework that turns validator and planner diagnostics into localized guidance for minimal edits toward executability.
- We open-source NL-PDDL-Bench, a multi-domain benchmark with planner-verified executable instances and controlled scaling by instance size.
- We introduce a planner-grounded optimization recipe for specification construction. It combines LoRA-based supervised fine-tuning, planner-derived preference signals via DPO, and inference-time planner-in-the-loop repair to systematically improve executability and plan-quality alignment.
- We develop a planner-grounded evaluation suite that measures solvability and plan-level agreement against reference solutions. Using this suite, we show consistent gains in planner success and plan-level agreement across model families, with improved robustness under difficulty scaling, cross-domain generalization, and transfer to PlanBench.

2 Definitions and Background

2.1 Planning Formulation

We focus on deterministic classical planning. A planning domain is denoted by $D = (V, A)$, where V is a finite set of predicate symbols (which induce ground atomic facts), and A is a set of action schemas with preconditions and effects. Given D , a problem instance specifies an object set O , an

initial state s_0 (a set of ground atoms), and a goal condition g_{true} (Volkema, 1983):

$$P = (D, O, s_0, g_{\text{true}}) \quad (1)$$

A classical planner searches for a finite plan $\pi = (a_0, \dots, a_n)$ such that each a_t is applicable in s_t and induces a valid transition

$$s_{t+1} = \gamma(s_t, a_t) \quad (2)$$

with the terminal state satisfying $s_n \models g_{\text{true}}$. Under the PDDL convention, the domain file declares predicates, types, and action schemas, while the problem file instantiates O , s_0 , and g_{true} .

2.2 Introduction to PDDL

PDDL is a widely used formalism for classical planning (Gerevini, 2020). A specification is separated into a domain file and a problem file: the domain file defines predicate symbols, object types, and action schemas with explicit preconditions and effects, while the problem file instantiates a concrete task by declaring objects, initial facts, and goal conditions. This decomposition yields a normalized state-action representation that can be parsed and solved by symbolic planners, enabling reproducible verification of executability and solution quality.

In this work, PDDL serves as the formal interface for connecting natural-language task descriptions with symbolic planning toolchains. A large language model is required to produce problem specifications that respect the predicate signatures and action constraints defined by the domain. Deviations such as missing object declarations, inconsistent predicate usage, or incomplete initial facts can render an instance invalid or unsolvable, and these failures are reliably exposed by planner-based validation and solving, providing grounded signals for diagnosis and correction.

3 Method

3.1 Overall Framework

We present a planner-in-the-loop feedback framework for natural language planning that couples large language models with a standard planning toolchain to produce specifications that are parseable, solvable, and verifiable, together with executable plans. The key principle is to place the planner in the generation loop. The model focuses on extracting task entities and constraints from text and expressing them as a symbolic specification,

181 while the validator and planner provide rigorous
 182 checks of well-formedness and solvability. Their
 183 diagnostics are treated as reproducible evidence
 184 that supports subsequent targeted revisions. This
 185 design separates semantic interpretation from state-
 186 space search, enabling evidence-driven minimal
 187 edits that improve executability under fixed sym-
 188 bolic semantics, rather than relying on surface-form
 189 plausibility.

190 To enable systematic evaluation and optimiza-
 191 tion, we build NL-PDDL-Bench and use a unified
 192 symbolic toolchain throughout dataset construc-
 193 tion, verification, training, and evaluation. Each in-
 194 stance includes an aligned natural-language descrip-
 195 tion, its corresponding specification, and planner-
 196 derived reference signals, which support consistent
 197 measurement of validity, solvability, and plan-level
 198 alignment. Figure 1 summarizes the workflow, in-
 199 cluding benchmark construction, diagnosis-guided
 200 refinement, and unified symbolic evaluation.

201 On the learning side, we combine parameter-
 202 efficient supervised fine-tuning(SFT) with planner-
 203 grounded preference optimization to improve ex-
 204 ecutability and solution-quality alignment. At in-
 205 ference time, when parsing or planning fails, the
 206 system converts diagnostic logs into structured feed-
 207 back and prompts the model to perform localized
 208 repairs until verification succeeds or a fixed budget
 209 is reached. Formal definitions, training objectives,
 210 and evaluation metrics are specified in the subse-

quent subsections.

3.2 NL-PDDL-Bench

211 NL-PDDL-Bench is designed to provide a repro-
 212 ducible and planner-verifiable foundation for train-
 213 ing and evaluating natural-language specification
 214 construction. Rather than relying on offline text-
 215 level alignment or manual rewriting, we build an
 216 end-to-end automated pipeline that combines top-
 217 down task design with bottom-up symbolic verifica-
 218 tion, following the general benchmark-construction
 219 spirit of automatically deriving LLM planning eval-
 220 uations from PDDL resources (Stein et al., 2023).
 221 This pipeline enforces executable alignment: only
 222 instances that satisfy symbolic constraints and pass
 223 planner verification are included in the executable
 224 subset used for learning and evaluation, while failed
 225 instances yield structured diagnostic evidence under
 226 a unified toolchain. To make verification repro-
 227 ducible and diagnostically useful, we rely on
 228 planner-backed validation as a systematic check of
 229 semantic executability and inconsistency exposure,
 230 consistent with classical-planner-based verification
 231 and counterexample generation practices (Goldman
 232 et al., 2012), and we adopt a strong IPC-tested plan-
 233 ning system as the underlying solver backbone (Cor-
 234 rêa et al., 2023).

235 Starting from classical IPC domains and their
 236 problem generators, we generate multi-domain
 237 planning instances with controlled scaling. We use
 238
 239

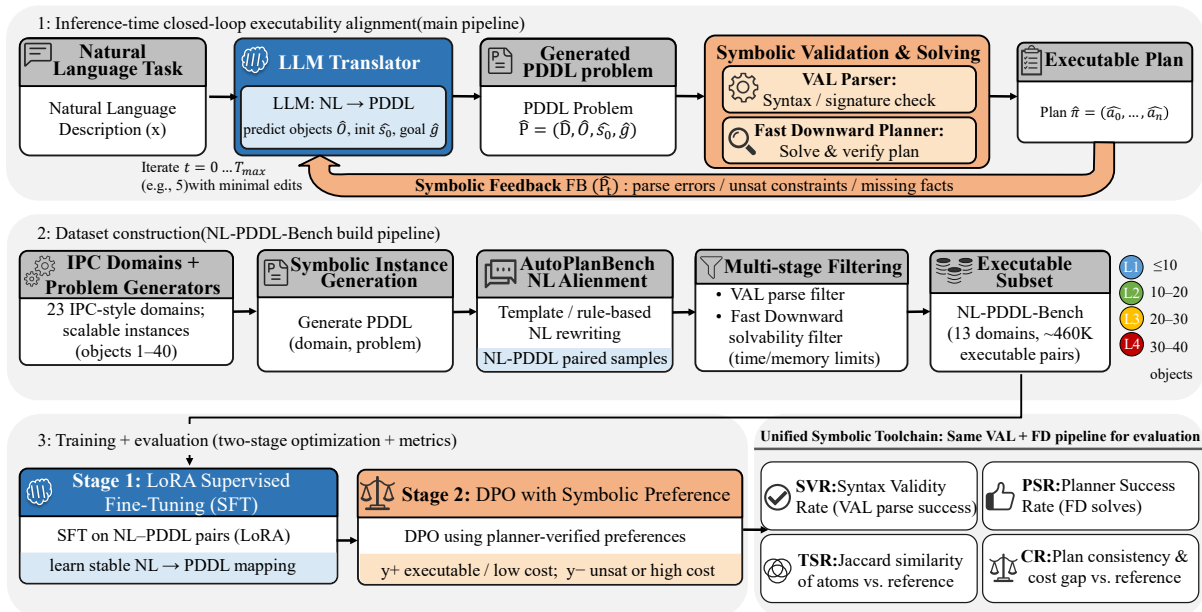


Figure 1: Overview of the planner-in-the-loop feedback framework, including benchmark construction, planner-in-the-loop feedback, and unified symbolic evaluation.

the number of objects as a unified difficulty surrogate and stratify problems into four levels (L1–L4), covering instances from 1 to 40 objects. For each symbolic instance, we then produce an aligned natural-language description via rule- and template-based, structure-preserving rewriting. The resulting descriptions explicitly encode initial conditions, goal conditions, and domain-critical constraints such as connectivity, capacities, and mutual exclusions, ensuring traceability between text and symbolic structure.

Table 1: Dataset distribution and planning statistics of NL-PDDL-Bench across difficulty levels.

Level	Objects	Solvable pairs	Plan length (median [Q1, Q3])
L1	1–10	154219	15.0 [12.8, 17.5]
L2	11–20	112003	39.0 [21.5, 91.5]
L3	21–30	93909	56.0 [19.0, 144.0]
L4	31–40	96918	86.5 [30.2, 208.5]

To guarantee executability and support planner-in-the-loop use cases, we apply multi-stage verification. We first run a validator to filter specifications that violate syntax, predicate signatures, typing, or arity constraints. We then invoke Fast Downward (FD) under a fixed resource budget and retain only planner-solvable instances as the executable subset, recording planner-derived reference signals including plans and plan cost/length, as well as failure logs for diagnosis. This pipeline starts from approximately 1.707M aligned NL–PDDL pairs and yields about 460K planner-solvable high-quality instances. Table 1 summarizes the resulting dataset distribution across difficulty levels and reports basic planning statistics.

3.3 LoRA-Based SFT and Planner-Grounded Preference Optimization

We adopt a two-stage optimization scheme to improve executability and plan-quality alignment of generated specifications. The first stage performs parameter-efficient SFT to learn a stable mapping from text to problem specifications, emphasizing object declarations, type consistency, predicate signatures, and faithful encoding of initial and goal facts. The second stage introduces planner-grounded preference optimization: we construct reusable preference signals via offline validation and planning, pushing the model toward specifications that are solvable and quality-aligned with reference solutions rather than merely well-formed.

For SFT, we use LoRA (Hu et al., 2022) and learn a low-rank update on top of frozen base weights,

$$W = W_0 + \Delta W = W_0 + \alpha B_{\text{LoRA}} A_{\text{LoRA}} \quad (3)$$

and minimize the standard negative log-likelihood,

$$\mathcal{L}_{\text{NLL}} = -\frac{1}{T} \sum_{t=1}^T \log P_{\theta}(y_t | y_{<t}, x) \quad (4)$$

Following common practice under constrained compute, we adopt QLoRA (Detmeters et al., 2023), which fine-tunes LoRA adapters on a frozen 4-bit quantized base model to reduce memory footprint while maintaining performance.

Supervision alone does not guarantee planner solvability. We therefore build an offline preference dataset: for each input x , the reference specification is treated as a positive sample y^+ ; candidate negatives are generated by programmatic perturbations and filtered by a validator and a planner to obtain y^- (unsolvable, or solvable but substantially worse than the reference in plan quality). We then apply Direct Preference Optimization (DPO) (Rafailov et al., 2023) by optimizing the policy f_{θ} against a frozen reference f_{θ_0} :

$$\begin{aligned} \mathcal{L}_{\text{DPO}} &= \mathbb{E}[-\log \sigma(\beta \Delta)], \\ \Delta &= \left(\log f_{\theta}(y^+ | x) - \log f_{\theta}(y^- | x) \right) \\ &\quad - \left(\log f_{\theta_0}(y^+ | x) - \log f_{\theta_0}(y^- | x) \right) \end{aligned} \quad (5)$$

Since preference labels are produced entirely by offline planner verification, this stage avoids planner calls during backpropagation while injecting stable symbolic signals for executability and quality alignment. Details of perturbations and thresholds are provided in the appendix.

3.4 Planner-in-the-Loop Feedback and Multi-Dimensional Evaluation

We build a verifiable neuro-symbolic loop by using a symbolic planner as a unified interface for parsing, solving, and diagnosis, and by placing it in the generation process at inference time. Due to space constraints, the full feedback workflow is illustrated in Appendix C.2 (Figure 5). The model first produces a candidate problem specification; a validator and a planner then attempt to parse and solve it, returning structured diagnostics. When parsing or planning fails, we convert diagnostics into localized constraints and revision directives, and prompt the

model to apply minimal edits while preserving the intended task semantics. The loop terminates once the specification becomes executable or the budget is exhausted. This design turns non-executable outputs into attributable symbolic defects and provides formally grounded evidence for both evaluation and optimization.

At test time, we create an isolated workspace per instance and record the domain file, the reference specification, intermediate model outputs, and planner logs. We first run a one-shot generation and invoke validation and planning. If parsing fails, we trigger localized repairs guided by validator errors; if parsing succeeds but planning fails, we provide feedback targeting structural omissions and consistency conflicts and continue iterative repairs. The final prediction is the best verified specification produced within the evaluation budget.

Let $\mathcal{D}_{\text{test}} = \{(x_i, y_i)\}_{i=1}^N$, where x_i is the input description, y_i is the reference specification, and \hat{y}_i is the model output (optionally after planner-in-the-loop refinement). We report four complementary metrics that emphasize executable alignment. Syntax Validity Rate (SVR) measures whether \hat{y}_i is parseable:

$$\text{SVR} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[\text{Parse}(\hat{y}_i) = \text{OK}] \quad (6)$$

Planner Success Rate (PSR) measures solvability conditional on parseability. Define $I_{\text{parse}} = \{i \mid \text{Parse}(\hat{y}_i) = \text{OK}\}$ and $I_{\text{solve}} = \{i \mid \text{Plan}(\hat{y}_i) = \text{SUCCESS}\}$, then

$$\text{PSR} = \frac{|I_{\text{parse}} \cap I_{\text{solve}}|}{|I_{\text{parse}}|} \quad (7)$$

Text Similarity Rate (TSR) measures atom-level overlap with Jaccard similarity:

$$\text{TSR} = \frac{1}{N} \sum_{i=1}^N \frac{|\text{Atoms}(\hat{y}_i) \cap \text{Atoms}(y_i)|}{|\text{Atoms}(\hat{y}_i) \cup \text{Atoms}(y_i)|} \quad (8)$$

where $\text{Atoms}(y)$ extracts the set of atomic facts from a specification y .

For plan-level alignment, we solve both the generated and reference problems under a fixed planner configuration and obtain $(\text{status}_i^{\text{gen}}, \pi_i^{\text{gen}}, c_i^{\text{gen}})$ and $(\text{status}_i^{\text{ref}}, \pi_i^{\text{ref}}, c_i^{\text{ref}})$, where c denotes plan cost or length. Consistency Rate (CR) is outcome-aware: it measures whether the planner outcomes agree (SUCCESS vs. FAILURE); when the reference instance

is solvable, it further requires agreement in plan quality and structure:

$$\begin{aligned} \text{CR} &= \frac{1}{N} \sum_{i=1}^N \mathbf{1}[\text{status}_i^{\text{gen}} = \text{status}_i^{\text{ref}} \wedge \Psi_i], \\ \Psi_i &= \left(\text{status}_i^{\text{ref}} = \text{FAILURE} \right) \\ &\quad \vee \left(|c_i^{\text{gen}} - c_i^{\text{ref}}| \leq \epsilon_i \wedge \pi_i^{\text{gen}} \approx \pi_i^{\text{ref}} \right) \end{aligned} \quad (9)$$

We detail the operational definition of $\pi_i^{\text{gen}} \approx \pi_i^{\text{ref}}$ and the choice of ϵ in Appendix C.3.

4 Experiments

4.1 Experimental Setup

Experiments are conducted on NL-PDDL-Bench to validate executability alignment and feedback-driven optimization across multiple model families, thirteen planning domains, and stratified difficulty levels. All evaluations rely on a unified symbolic toolchain: a validator checks well-formedness, and a symbolic planner assesses solvability and produces reproducible plans and failure diagnostics. Fast Downward (FD) is used throughout under a fixed configuration, yielding deterministic behavior and logs for consistent verification and attribution (Helmert, 2006). All reported numbers are computed with the same tool versions and resource limits to ensure comparability.

The test set is drawn from the executable subset of NL-PDDL-Bench. From approximately four hundred and sixty thousand planner-verified instances, we sample evenly across the thirteen domains and difficulty strata to obtain one thousand test instances. The test set is disjoint from the data used during model fine-tuning. For planner-in-the-loop evaluation, the maximum number of feedback iterations is fixed to $T_{\text{max}} = 5$.

Evaluation targets a common failure mode in language-based specification generation: outputs may be syntactically well formed yet violate preconditions, effect consistency, or reachability constraints under planner verification. We therefore report executability metrics and plan-level alignment against planner-derived references, and conduct three analyses: (i) comparison with direct action-sequence generation, (ii) transfer experiments on PlanBench, and (iii) difficulty- and domain-wise analyses for robustness under scaling and cross-domain variation.

4.2 Comparative Results

4.2.1 Comparing Against Mainstream LLMs for Direct Plan Generation

This paper compares mainstream large language models that directly generate complete action sequences, including GPT-4o, Claude 3.5, and DeepSeek-R1, against our final method. Our method uses a low-rank-adapted Llama 3.1-8B model with symbolic feedback to produce a planning specification, and then applies the FD planner to synthesize an executable plan. We report planner success rate on three classic domains. For Logistics and Zenotravel, we evaluate 100 instances per domain sampled from our benchmark; for Blocksworld, we evaluate 600 instances from the PlanBench Mystery variant.

Table 2: Planner success rate for direct plan generation vs. our method.

Domain	GPT-4o	Claude 3.5	DeepSeek-R1	Ours
Logistics	22.4%	22.6%	56.8%	78.0%
Zenotravel	33.7%	57.9%	98.9%	99.0%
Blocksworld	0.0%	0.0%	43.3%	87.3%

We include plan generation as a baseline to test a central claim: when correctness is governed by preconditions, effect consistency, and reachability, fluent action sequences are not reliable evidence of executability. Direct generation must satisfy long chains of implicit constraints, so a single local omission can invalidate downstream transitions. This brittleness is pronounced in Blocksworld, where validity is tightly coupled to long-horizon state changes, whereas our approach leverages a verifiable symbolic interface that delegates search to a planner and uses planner-in-the-loop diagnostics for localized repair. Although direct generation can be strong in some domains (e.g., Zenotravel), it is less stable across domains. The gap is driven less by surface form than by missing constraint closure; planner-in-the-loop feedback turns global plan failures into localized, formally grounded corrections.

4.2.2 Overall Results of Mainstream Models

We evaluate eight representative model families on the NL-PDDL-Bench test set under two settings: Baseline, where each open-source base model follows the same prompt to directly generate a PDDL problem that is then solved under the same Fast Downward configuration, and our full pipeline that combines LoRA-based supervised

fine-tuning, planner-grounded preference optimization, and inference-time planner-in-the-loop repair. We test a central claim: when correctness is jointly governed by preconditions, effect consistency, and reachability, failures are dominated by missing constraint closure rather than surface formatting.

Table 3 shows that this behavior holds broadly across model families. Averaged over all eight models, planner-in-the-loop raises planner success from 17.8% to 54.8% and plan-level agreement from 19.9% to 55.7%, while syntax validity increases more modestly from 66.8% to 74.7%. Except for GPT-OSS-20B, most baselines exhibit low PSR/CR, leading to larger absolute gains once planner feedback is introduced: PSR and CR increase (or remain unchanged) for every model, and TSR improves for seven out of eight models, consistent with planner-grounded signals correcting structural omissions and constraint conflicts.

The per-model results illustrate the same trend under diverse base capabilities. For instance, Llama 3.1-8B and Qwen3-8B both exhibit large gains in PSR/CR (e.g., for Llama 3.1-8B, PSR increases from 9.7% to 80.7% and CR increases from 10.9% to 81.2%; for Qwen3-8B, PSR increases from 24.8% to 78.9% and CR increases from 25.6% to 79.7%), indicating that planner diagnostics effectively convert global failures into localized, correctable revisions. In contrast, GPT-OSS-20B shows essentially flat PSR/CR with decreased TSR, suggesting deeper semantic mismatches that are less amenable to conservative local repair; we therefore conduct a domain-wise analysis of failure patterns. This analysis further localizes the anomaly to specific domains, indicating that the lack of improvement is driven by concentrated domain-level failures rather than uniform degradation across the benchmark. Overall, the results support our thesis that planner-in-the-loop feedback is most effective when failures are diagnosable and local, and that syntax compliance is neither necessary nor sufficient for executable alignment.

4.2.3 Comparison and Transferability on PlanBench

Table 4 contrasts NL-PDDL-Bench and PlanBench along task focus, scale, and verification. Both benchmarks are grounded in PDDL and rely on planner-based checking, but they target different bottlenecks. PlanBench evaluates planning and state-change reasoning from natural-language prompts, with a limited set of domains and small

Table 3: Overall results on the NL-PDDL-Bench test set. Parentheses report the absolute change of Planner-in-the-loop over Baseline for each metric.

Model	Setting	SVR	PSR	TSR	CR
Gemma3-4B	Baseline	37.2%	0.8%	6.8%	10.1%
	Planner-in-the-loop	60.6% (↑23.4%)	25.5% (↑24.7%)	32.3% (↑25.5%)	29.4% (↑19.3%)
GLM-4.1V-9B	Baseline	75.6%	9.9%	21.1%	11.1%
	Planner-in-the-loop	80.0% (↑4.4%)	63.3% (↑53.4%)	75.9% (↑54.8%)	64.2% (↑53.1%)
Llama3.1-8B	Baseline	76.3%	9.7%	52.9%	10.9%
	Planner-in-the-loop	84.9% (↑8.6%)	80.7% (↑71.0%)	82.8% (↑29.9%)	81.2% (↑70.3%)
Qwen2.5-7B	Baseline	76.8%	21.7%	61.2%	23.3%
	Planner-in-the-loop	81.8% (↑5.0%)	62.5% (↑40.8%)	78.8% (↑17.6%)	61.9% (↑38.6%)
Qwen3-8B	Baseline	82.1%	24.8%	61.5%	25.6%
	Planner-in-the-loop	82.7% (↑0.6%)	78.9% (↑54.1%)	80.8% (↑19.3%)	79.7% (↑54.1%)
GPT-OSS-20B	Baseline	69.9%	37.5%	39.9%	38.3%
	Planner-in-the-loop	79.0% (↑9.1%)	37.6% (↑0.1%)	31.5% (↓8.4%)	38.4% (↑0.1%)
DeepSeek-R1-7B	Baseline	84.6%	32.4%	34.1%	33.7%
	Planner-in-the-loop	90.2% (↑5.6%)	58.3% (↑25.9%)	57.0% (↑22.9%)	58.1% (↑24.4%)
DeepSeek-Coder-V2-16B	Baseline	32.2%	5.3%	13.1%	6.1%
	Planner-in-the-loop	38.5% (↑6.3%)	32.0% (↑26.7%)	32.8% (↑19.7%)	32.8% (↑26.7%)

problem sizes. NL-PDDL-Bench emphasizes verifiable specification construction under a unified symbolic toolchain, enabling systematic measurement of syntactic validity, planner solvability, and plan-level agreement against reference solutions. This design shifts evaluation from surface plausibility to executable correctness, and makes failures reproducible and diagnosable through planner logs.

Table 4: Key differences between Ours and PlanBench.

Aspect	NL-PDDL-Bench	PlanBench
Primary target	Spec construction and executability / plan-level alignment	Plan generation + cost-oriented evaluation
Domain coverage	13 planner-verified domains; cross-domain by design	2 IPC domains: Blocksworld, Logistics.
Scale	1.7M pairs (460K executable)	26.3K prompts
Instance size	Difficulty-stratified; 1–40 objects	No explicit strata; typically 1–5 objects
Verification	Unified parsing/solving and log-level diagnostics	Planner/validator checks; less standardized diagnostics
Training support	Large executable subset enables training + feedback optimization	Limited training-oriented support

Compared with PlanBench, NL-PDDL-Bench provides broader coverage and stronger training support. It spans thirteen planner-verified domains and controls difficulty through object-scale stratifica-

tion, while PlanBench covers two domains with substantially smaller instances. NL-PDDL-Bench also offers a substantially larger collection of aligned pairs with an executable subset, which supports supervised training and feedback-driven optimization. Finally, NL-PDDL-Bench standardizes verification through consistent parsing and solving procedures and retains log-level diagnostics, enabling fine-grained attribution of errors such as signature mismatches, missing initial facts, and constraint conflicts. These properties make NL-PDDL-Bench complementary to PlanBench, and better suited for studying executability alignment and robustness under scaling and cross-domain transfer.

Table 5: Results on PlanBench (Baseline vs. Ours).

Metric	Qwen3-8B	Llama3.1-8B	GLM4.1V-9B
SVR	43.3%	99.8%	96.2%
	99.8% (↑56.5%)	100.0% (↑0.2%)	93.2% (↓3.0%)
PSR	0.0%	7.2%	21.7%
	32.2% (↑32.2%)	49.2% (↑42.0%)	69.0% (↑47.3%)
TSR	0.0%	0.0%	2.8%
	0.0%	0.0%	4.9% (↑2.1%)
CR	14.3%	7.2%	44.7%
	44.3% (↑30.0%)	49.3% (↑42.1%)	64.3% (↑19.6%)

To assess transferability, we evaluate our models on PlanBench using its original verification protocol. Table 5 shows that the full training and planner-guided refinement pipeline substantially improves planner success and plan-level agreement across

model families, even though the benchmark differs in domain scope and instance characteristics. The gains indicate that planner-grounded training signals and refinement generalize beyond the source benchmark, improving reliability under a different evaluation setting rather than merely overfitting to the construction pipeline.

4.3 Analytical Experiments

4.3.1 Performance under Difficulty Scaling

To test our central claim that planning reliability is governed by executable constraint satisfaction rather than surface-form fluency, we conduct a difficulty-scaling analysis. We stratify instances by object scale and evaluate how specification generation degrades as horizons and constraint coupling intensify, and whether planner-in-the-loop feedback can arrest such degradation.

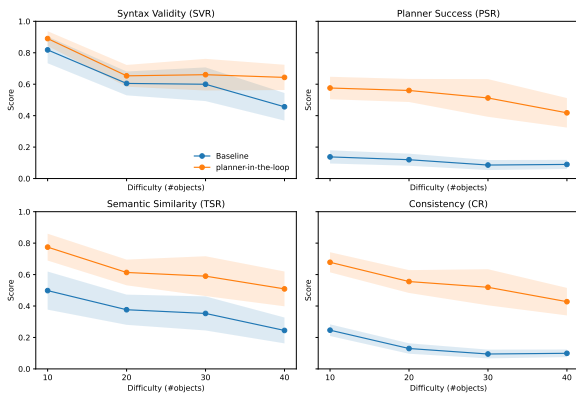


Figure 2: Performance under difficulty scaling.

The results support this diagnosis. As the object scale increases from 10 to 40, the baseline planner success drops from about 13% to about 9%, plan-level agreement drops from about 25% to about 9%, and specification similarity drops from about 50% to about 26%, with the overall score decreasing from about 42% to about 22%. In contrast, with planner-in-the-loop feedback, performance remains substantially more stable at the hardest scale: syntax validity is about 64%, planner success and plan-level agreement are both about 42%, and the overall score stays near 50%, suggesting improved robustness under scaling.

4.3.2 Cross-Domain Generalization

We further study cross-domain generalization to assess whether gains reflect domain-specific memorization or a general improvement in executable alignment. Using the same evaluation protocol,

we compare baseline and planner-in-the-loop performance by domain and interpret metric patterns under different constraint structures. Because domains vary in their dominant constraints, executable alignment should primarily improve planner success and plan-level agreement across domains rather than only syntax.

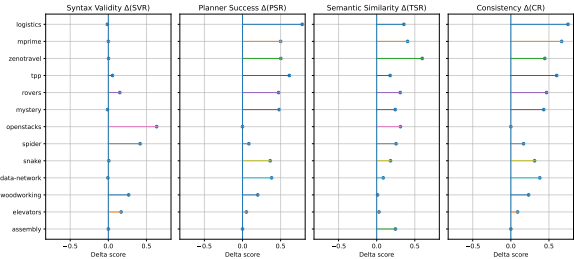


Figure 3: Cross-domain generalization results on NL-PDDL-Bench.

The observed trends support this expectation. In many domains, syntax validity is already relatively high and changes modestly, whereas planner success and plan-level agreement improve substantially, indicating that the main benefit comes from constraint-level alignment. For example, in Logistics, planner success increases from about 13% to about 78%, and plan-level agreement increases from about 20% to about 82%, consistent with repairing parseable but structurally invalid specifications such as missing connectivity facts or capacity conflicts. In contrast, in resource-intensive domains such as Openstacks, syntax improvements do not yield comparable gains in executability, reinforcing that syntax compliance is not a reliable proxy for executable alignment.

5 Conclusion

We introduce NL-PDDL-Bench, a multi-domain benchmark for natural-language-to-PDDL specification generation, with planner-verified executability and object-count scaling. We present a planner-in-the-loop framework that maps tasks to executable PDDL and applies localized repairs from validator and planner diagnostics. We further propose a planner-grounded optimization recipe, together with a unified evaluation suite for parseability, solvability, similarity, and outcome-aware plan consistency. Across eight model families and thirteen domains, we substantially improve planner success and plan-level agreement with only modest syntax changes, and transfer robustly to harder scales, new domains, and PlanBench.

6 Limitations

Although NL-PDDL-Bench and our planner-in-the-loop feedback framework provide a unified, planner-verifiable pipeline for executable alignment, our current study still has limitations along several dimensions:

- **Broader settings.** We evaluate on a fixed set of classical PDDL domains under a standardized validation and planning setup; extending to richer PDDL fragments (e.g., numeric fluents, temporals, derived predicates), alternative planners, and heterogeneous resource conditions is a natural next step.
- **Richer language.** Our benchmark uses structure-preserving, rule/template-based generation to maintain traceability between text and symbols; incorporating more diverse, naturally occurring instructions would broaden linguistic coverage.
- **Computational cost and configuration generality.** Our plan-level metrics and verification loop run under a fixed planner configuration and a bounded iteration budget, which improves reproducibility but incurs additional inference-time overhead and limits systematic evaluation across planners and heuristics. Future work may reduce cost via caching, early stopping, and adaptive budgeting, and further study stability across planner/heuristic configurations to strengthen scalability and deployment practicality.

We are actively exploring these directions in ongoing experiments, aiming to broaden coverage and improve practicality without changing the core executable-alignment protocol.

Acknowledgments

We thank the anonymous reviewers for their valuable comments.

References

Constructions Aeronautiques, Adele Howe, Craig Knoblock, ISI Drew McDermott, Ashwin Ram, Manuela Veloso, Daniel Weld, David Wilkins Sri, Anthony Barrett, Dave Christianson, and 1 others. 1998. Pddl—the planning domain definition language. *Technical Report, Tech. Rep.*

Pascal Bercher, Sarath Sreedharan, and Mauro Vallati. 2025. A survey on model repair in ai planning. In

34th International Joint Conference on Artificial Intelligence, page 26. IJCAI Organization.

Augusto B Corrêa, Guillem Frances, Markus Hecher, Davide Mario Longo, and Jendrik Seipp. 2023. The powerlifted planning system in the ipc 2023. *Tenth International Planning Competition (IPC-10): Planner Abstracts*.

Augusto B Corrêa, André G Pereira, and Jendrik Seipp. 2025. Classical planning with llm-generated heuristics: Challenging the state of the art with python code. *arXiv preprint arXiv:2503.18809*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115.

Alfonso Emilio Gerevini. 2020. An introduction to the planning domain definition language (pddl): Book review. *Artificial Intelligence*, 280:103221.

Robert P Goldman, Ugur Kuter, and A Schneider. 2012. Using classical planners for plan verification and counterexample generation. In *Proceedings of AAAI Workshop on Problem Solving Using Classical Planning*.

Malte Helmert. 2006. The fast downward planning system. *Journal of Artificial Intelligence Research*, 26:191–246.

Richard Howey, Derek Long, and Maria Fox. 2004. Val: Automatic plan validation, continuous effects and mixed initiative planning using pddl. In *16th IEEE International Conference on Tools with Artificial Intelligence*, pages 294–301. IEEE.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.

Cassie Huang and Li Zhang. 2025. On the limit of language models as planning formalizers. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4880–4904.

Sukai Huang, Nir Lipovetzky, and Trevor Cohn. 2025. Planning in the dark: Llm-symbolic planning pipeline without experts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 26542–26550.

Harsh Kokel and Michael Katz. 2025. *Acpbench: Reasoning about action, change, and planning*. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 26559–26568.

Sadegh Mahdavi, Raquel Aoki, Keyi Tang, and Yanshuai Cao. 2024. Leveraging environment interaction for automated pddl translation and planning with large language models. *Advances in Neural Information Processing Systems*, 37:38960–39008.

706	Martin Pozo and Jendrik Seipp. 2025. Abstraction
707	heuristics for classical planning tasks with conditional
708	effects. In <i>Proceedings of the Thirty-Fourth Inter-</i>
709	<i>national Joint Conference on Artificial Intelligence</i> ,
710	pages 8608–8616.
711	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-
712	pher D Manning, Stefano Ermon, and Chelsea Finn.
713	2023. Direct preference optimization: Your language
714	model is secretly a reward model. <i>Advances in neural</i>
715	<i>information processing systems</i> , 36:53728–53741.
716	Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B
717	Tenenbaum, Leslie Kaelbling, and Michael Katz.
718	2024. Generalized planning in pddl domains with pre-
719	trained large language models. In <i>Proceedings of the</i>
720	<i>AAAI conference on artificial intelligence</i> , volume 38,
721	pages 20256–20264.
722	Pavel Smirnov, Frank Joubin, Antonello Ceravola, and
723	Michael Gienger. 2024. Generating consistent pddl
724	domains with large language models. <i>arXiv preprint</i>
725	<i>arXiv:2404.07751</i> .
726	Katharina Stein, Daniel Fišer, Jörg Hoffmann, and
727	Alexander Koller. 2023. Autoplanbench: Automat-
728	ically generating benchmarks for llm planners from
729	pddl. <i>arXiv preprint arXiv:2311.09830</i> .
730	Karthik Valmeekam, Matthew Marquez, Alberto Olmo,
731	Sarath Sreedharan, and Subbarao Kambhampati.
732	2023a. Planbench: An extensible benchmark for
733	evaluating large language models on planning and
734	reasoning about change. <i>Advances in Neural Infor-</i>
735	<i>mation Processing Systems</i> , 36:38975–38987.
736	Karthik Valmeekam, Sarath Sreedharan, Matthew Mar-
737	quez, Alberto Olmo, and Subbarao Kambhampati.
738	2023b. On the planning abilities of large language
739	models (a critical investigation with a proposed bench-
740	mark). <i>arXiv preprint arXiv:2302.06706</i> .
741	Roger J Volkema. 1983. Problem formulation
742	in planning and design. <i>Management science</i> ,
743	29(6):639–652.
744	Xiaopan Zhang, Hao Qin, Fuquan Wang, Yue Dong,
745	and Jiachen Li. 2025. Lamma-p: Generalizable
746	multi-agent long-horizon task allocation and plan-
747	ning with lm-driven pddl planner. In <i>2025 IEEE In-</i>
748	<i>ternational Conference on Robotics and Automation</i>
749	(<i>ICRA</i>), pages 10221–10221. IEEE.
750	Max Zuo, Francisco Piedrahita Velez, Xiaochen Li,
751	Michael Littman, and Stephen Bach. 2025. Plane-
752	tarium: A rigorous benchmark for translating text
753	to structured planning languages. In <i>Proceedings of</i>
754	<i>the 2025 Conference of the Nations of the Americas</i>
755	<i>Chapter of the Association for Computational Lin-</i>
756	<i>guistics: Human Language Technologies (Volume 1:</i>
757	<i>Long Papers)</i> , pages 11223–11240.

A Appendix: Introduction to PDDL

Planning Domain Definition Language (PDDL) is a standard formalism for representing deterministic classical planning tasks. A specification is decomposed into a *domain* and a *problem*. The domain defines predicate symbols, object types, and action schemas with explicit preconditions and effects. The problem instantiates a concrete instance by declaring objects, initial facts, and goal conditions. This separation yields a normalized state–action model that can be parsed and solved by symbolic planners, enabling reproducible verification of executability.

Figure 4 illustrates the semantics using a simplified Logistics example that contains a minimal but complete closed-world task: one truck must deliver one package from Location A to Location B. The example is split into three parts: the domain file `logistics.pddl` (left), the problem file `task.pddl` (top-right), and a schematic state-transition trace consistent with a valid plan (bottom-right).

Domain file (`logistics.pddl`). The domain declares `:requirements` `:strips` `:typing`, indicating that the task uses STRIPS-style add/delete effects with typed objects. It introduces three types—`truck`, `package`, and `location`—and three predicates: `(at ?x ?l)` encodes the location of either a truck or a package, `(in ?p ?t)` encodes that a package is loaded in a truck, and `(connected ?from ?to)` encodes reachability between locations. These predicates define the state vocabulary, i.e., the set of ground atoms that can be true in any state.

The domain further defines three action schemas. Action `drive` moves a truck from `from` to `to` when the truck is currently at `from` and the connectivity relation holds; its effect deletes `(at t from)` and adds `(at t to)`. Action `load` requires a package and the truck to co-locate, and updates the state by removing the package location fact `(at p l)` and adding `(in p t)`. Action `unload` is the inverse: it requires `(in p t)` and the truck to be at a location, then removes `(in p t)` and adds `(at p l)`. Together, these schemas constrain which transitions are legal and make executability checkable via precondition satisfaction and effect application.

Problem file (`task.pddl`). The problem instantiates a specific task under the above domain. It declares a truck `t1`, a package `p1`, and two locations

808 A and B, each with the correct type. The initial state
809 sets $t1$ and $p1$ at A, and provides connectivity facts
810 (`connected A B`) and (`connected B A`). The
811 goal requires the package to be at B, i.e., (`at p1`
812 `B`). Importantly, this problem illustrates a common
813 source of failure in language-to-PDDL generation:
814 if connectivity is omitted or mis-specified, the in-
815 stance can remain syntactically valid yet become
816 unsolvable due to broken reachability constraints.

817 **State-transition trace and an executable plan.**

818 Given the domain dynamics and the instantiated
819 facts, a planner can synthesize an executable
820 plan that achieves the goal. One minimal plan
821 is: `load(p1, t1, A)`, then `drive(t1, A, B)`, then
822 `unload(p1, t1, B)`. Starting from the initial state,
823 `load` is applicable because both $t1$ and $p1$ are
824 at A; after applying its effects, the state contains
825 (`in p1 t1`) and no longer contains (`at p1 A`).
826 Next, `drive` is applicable because $t1$ is at A and
827 (`connected A B`) holds; its effects move the
828 truck to B. Finally, `unload` is applicable because
829 $p1$ is in $t1$ and $t1$ is at B; after unloading, the state
830 satisfies the goal atom (`at p1 B`). The schematic
831 diagram in Figure 4 visualizes these intermediate
832 states and highlights how PDDL makes the underlying
833 constraints explicit and mechanically verifiable.

834 **Relevance to NL \rightarrow PDDL generation.** When a
835 large language model produces a problem specifi-
836 cation from a natural-language description, it must
837 preserve (i) type-correct object declarations, (ii)
838 predicate signatures and arities, and (iii) domain-
839 critical initial facts that ensure reachability and re-
840 source consistency (e.g., `connected` relations). Vi-
841 olations of (i)–(ii) typically yield parser/type errors,
842 whereas violations of (iii) often produce instances
843 that parse but are unsolvable. This example there-
844 fore captures the key distinction between surface
845 well-formedness and true executability that moti-
846 vates planner-based verification and feedback.

847 **B Appendix: Related Work**

848 This appendix situates our study within three
849 closely related threads: language-model-based
850 planning, symbolic interfaces with planner veri-
851 fication, and feedback-driven refinement. Across
852 these threads, a common theme is the gap between
853 surface plausibility and executable correctness un-
854 der formal planning semantics. Our work addresses
855 this gap by unifying benchmark construction, veri-
856 fication, and optimization under a single symbolic

857 toolchain, while explicitly measuring plan-level
858 alignment in addition to syntactic well-formedness
859 and solvability.

860 **B.1 Large Language Models for Planning**

861 Large language models exhibit strong semantic pars-
862 ing and structured generation capabilities, which
863 has spurred growing interest in applying them to
864 planning problems. Nevertheless, classical plan-
865 ning imposes a stringent notion of correctness: ac-
866 tion applicability must satisfy preconditions, effects
867 must preserve state consistency, and goals must
868 be reachable through a valid transition sequence
869 (Valmeekam et al., 2023a). Empirical benchmarks
870 consistently show that models can generate plans
871 that read as coherent yet fail under symbolic execu-
872 tion or validation, indicating that natural-language
873 fluency is not a reliable proxy for constraint satis-
874 faction or long-horizon state-transition correctness
875 (Valmeekam et al., 2023a; Kokel and Katz, 2025).
876 In contrast, established symbolic planners remain
877 strong on well-specified domains, owing to mature
878 heuristic search, soundness guarantees under the
879 given semantics, and predictable performance pro-
880 files when resources are controlled (Corrêa et al.,
881 2025).

882 Recent analyses further characterize the limits of
883 LLMs as planning formalizers, showing that seem-
884 ingly minor specification errors (e.g., in signatures,
885 initial facts, or reachability-relevant constraints) can
886 systematically break executability under strict sym-
887 bolic semantics (Huang and Zhang, 2025). Our
888 study is consistent with these observations, but goes
889 beyond diagnosing failure modes by providing a uni-
890 fied benchmark and a planner-in-the-loop optimiza-
891 tion pipeline that turns validator/planner outcomes
892 into reproducible supervision signals and localized
893 repairs, enabling systematic improvements in exe-
894 cutable alignment rather than surface plausibility.

895 **B.2 Symbolic Interfaces and Planner** 896 **Verification**

897 A complementary line of work adopts symbolic task
898 specifications as the interface between language un-
899 derstanding and planning, enabling parsing, veri-
900 fication, and diagnosis through standard planning
901 toolchains (Mahdavi et al., 2024). This interface
902 has two practical advantages. First, it separates
903 semantic interpretation from search, since the plan-
904 ner operates on a formal specification whose se-
905 mantics are fixed by the domain model. Second, it
906 makes failure modes more interpretable, because

(1) Domain: logistics.pddl

```
(define (domain logistics)
  (:requirements :strips :typing)
  (:types truck package location)
  (:predicates (at x - (either truck package) l - location)
               (in p - package t - truck)
               (connected from - location to - location))

  (:action drive
   :parameters (t - truck from - location to - location)
   :precondition (and (at t from) (connected from to))
   :effect (and (not (at t from)) (at t to)))

  (:action load
   :parameters (p - package t - truck l - location)
   :precondition (and (at p l) (at t l))
   :effect (and (not (at p l)) (in p t)))

  (:action unload
   :parameters (p - package t - truck l - location)
   :precondition (and (in p t) (at t l))
   :effect (and (not (in p t)) (at p l)))
)
```

(2) Problem: task.pddl

```
(define (problem task)
  (:domain logistics)
  (:objects t1 - truck p1 - package A B - location)
  (:init
   (at t1 A)
   (at p1 A)
   (connected A B)
   (connected B A))
  (:goal (and (at p1 B)))
)
```

(3) Schematic Diagram

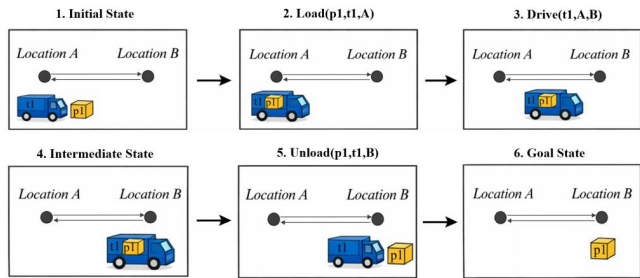


Figure 4: A simple example of the PDDL language.

907 validators and planners can return concrete error
908 messages, unsatisfied conditions, or evidence of
909 unreachable goals.

910 Benchmark efforts for text-to-structured plan-
911 ning languages provide an important foundation for
912 positioning NL→PDDL as a rigorous translation-
913 and-verification problem (Zuo et al., 2025). While
914 such benchmarks emphasize careful evaluation of
915 structured outputs, our focus is on *executable align-*
916 *ment at scale*: NL-PDDL-Bench provides a large
917 planner-verifiable executable subset with controlled
918 difficulty scaling, and our evaluation explicitly mea-
919 sures plan-level agreement under a fixed symbolic
920 configuration. In addition, we treat planner logs and
921 diagnostics as first-class artifacts that are reused
922 consistently for dataset construction, evaluation,
923 and optimization.

924 Related work also studies consistency checking
925 for LLM-generated PDDL artifacts. For example,
926 Smirnov et al. (2024) targets domain-level consis-
927 tency when generating PDDL domains with LLMs,
928 reducing contradictions in action/predicate defini-
929 tions via automated checks. Our setting is comple-
930 mentary: we focus on problem-level specifications
931 and emphasize failures that remain syntactically
932 well-formed yet are unsolvable due to missing reach-
933 ability facts or constraint closure. Accordingly, we
934 rely on a validator-plus-planner pipeline that ex-
935 poses not only static consistency issues but also
936 planning-time reachability and resource conflicts,
937 and we convert these reproducible diagnostics into

localized repair signals.

B.3 Feedback-Driven Refinement

938
939
940 To address brittleness, recent studies incorporate
941 external feedback to iteratively refine model out-
942 puts. One approach leverages environment or simu-
943 lator interaction to produce executability signals
944 and guide progressive correction (Mahdavi et al.,
945 2024). Another approach couples generation with
946 planner-based checking and replanning to recover
947 feasible solutions in long-horizon settings (Huang
948 et al., 2025; Zhang et al., 2025). These methods
949 share the intuition that correctness should be an-
950 chored in externally verifiable signals rather than
951 purely in model likelihood.

952 Despite clear progress, two limitations are recur-
953 rent in the literature. First, evaluations often focus
954 primarily on feasibility, while providing less em-
955 phasis on whether the generated solutions match
956 reference plan structure or achieve comparable so-
957 lution quality under a fixed planning configuration.
958 Second, analyses of robustness under controlled
959 scaling and cross-domain variation are compara-
960 tively limited, making it difficult to characterize
961 when feedback leads to general improvements ver-
962 sus domain-specific gains.

963 Our work addresses these gaps by introducing
964 plan-level consistency and quality-alignment mea-
965 surements grounded in planner-derived reference
966 solutions, and by evaluating behavior under explicit
967 difficulty stratification and multi-domain coverage.

968	This enables a more fine-grained view of what feed-	tion. The released executable subset contains 13	1017
969	back corrects, and under what conditions improve-	fully verified domains: <i>Assembly</i> , <i>Data Network</i> ,	1018
970	ments persist as constraints become more tightly	<i>Elevators</i> , <i>Logistics</i> , <i>MPrime</i> , <i>Mystery</i> , <i>Openstacks</i> ,	1019
971	coupled.	<i>Rovers</i> , <i>Snake</i> , <i>Spider</i> , <i>Traveling Purchaser Prob-</i>	1020
972	B.4 Our Position	<i>lem (TPP)</i> , <i>Woodworking</i> , and <i>ZenoTravel</i> . These	1021
973	Our study lies at the intersection of symbolic in-	domains provide complementary constraint struc-	1022
974	terfaces and feedback-driven refinement, with a	tures, spanning transportation and navigation tasks	1023
975	primary focus on executable alignment and re-	(e.g., <i>Logistics</i> , <i>Rovers</i> , <i>ZenoTravel</i>), resource- and	1024
976	producibility. We contribute NL-PDDL-Bench, a	capacity-intensive settings (e.g., <i>Openstacks</i> , <i>TPP</i> ,	1025
977	multi-domain benchmark with planner-verified ex-	<i>Woodworking</i>), and domains with stronger combi-	1026
978	ecutability and controlled scaling by instance size,	inatorial structure in objects and relations (e.g., <i>As-</i>	1027
979	which supports systematic analysis across domains	<i>sembly</i> , <i>Data Network</i> , <i>Mystery</i>), which supports	1028
980	and difficulty levels. We further provide a uni-	controlled cross-domain comparison and robust-	1029
981	fied symbolic toolchain that is used consistently	ness analyses.	1030
982	for benchmark construction, verification, training-	The reduction from 23 initial domains to the fi-	1031
983	signal derivation, and evaluation, ensuring that cor-	nal 13 prioritizes reproducible executability and	1032
984	rectness claims are grounded in the same planner	consistent NL–symbol alignment. Concretely, we	1033
985	semantics and are directly comparable across ex-	retain domains whose generators produce stable	1034
986	perimental settings.	multi-scale instances, whose domain/problem spec-	1035
987	Building on this foundation, we develop an end-	ifications are reliably parseable and solvable under a	1036
988	to-end pipeline for specification construction that	fixed validation and planning configuration, whose	1037
989	combines parameter-efficient supervised training,	rule/template-based rewriting preserves traceabili-	1038
990	planner-grounded preference optimization based	ty between text and symbolic structure, and whose	1039
991	on offline planner-derived preference pairs, and	verification pipeline yields a sufficient volume of ex-	1040
992	inference-time planner-in-the-loop repair guided by	ecutable instances for training and evaluation. The	1041
993	validator and planner diagnostics. The pipeline is	remaining domains are excluded when they (i) rely	1042
994	evaluated with complementary metrics that separate	on PDDL features outside the fragment supported	1043
995	syntactic well-formedness from planner solvabil-	in our standardized toolchain, (ii) produce too few	1044
996	ity and from outcome-aware plan-level consistency,	solvable instances under the fixed planner configura-	1045
997	thereby measuring executable correctness beyond	tion and resource budget, or (iii) exhibit instability	1046
998	surface plausibility.	in generation-to-rewriting consistency that harms	1047
999	In contrast to benchmarks that mainly frame	reproducibility. For transparency, we record the full	1048
1000	natural-language-to-structured planning as a trans-	initial domain pool and the excluded domains with	1049
1001	lation and evaluation task (Zuo et al., 2025), our	their exclusion reasons in the released benchmark	1050
1002	work unifies verification and optimization under	metadata.	1051
1003	the same symbolic infrastructure. The toolchain	Action costs and cost/length convention. To	1052
1004	that certifies executability also produces reusable	keep metrics comparable across domains, we use	1053
1005	diagnostics that support both preference construc-	a uniform plan-cost convention in all experiments.	1054
1006	tion and localized repair, enabling scalable and re-	By default, we assume unit action costs and treat	1055
1007	producible improvements in executable alignment	plan cost as plan length; correspondingly, plan-level	1056
1008	under controlled difficulty scaling, cross-domain	agreement is computed using length-based cost and	1057
1009	generalization, and transfer to external benchmarks.	structure under a fixed planner configuration. When	1058
1010	C Appendix: Method Details	an original IPC domain encodes action costs, we	1059
1011	C.1 NL-PDDL-Bench	still report results with length as the cost definition	1060
1012	NL-PDDL-Bench is built from an initial pool of 23	under the same standardized configuration, so that	1061
1013	IPC-style planning domains and their problem gen-	cross-domain trends reflect executable alignment	1062
1014	erators collected in AutoPlanBench, and processed	rather than heterogeneous cost modeling choices.	1063
1015	under a unified symbolic toolchain for instance gen-	C.2 Planner-in-the-Loop Feedback Workflow	1064
1016	eration, natural-language rewriting, and verifica-	Figure 5 illustrates the prompting and verification	1065
		loop used in our planner-in-the-loop workflow. The	1066

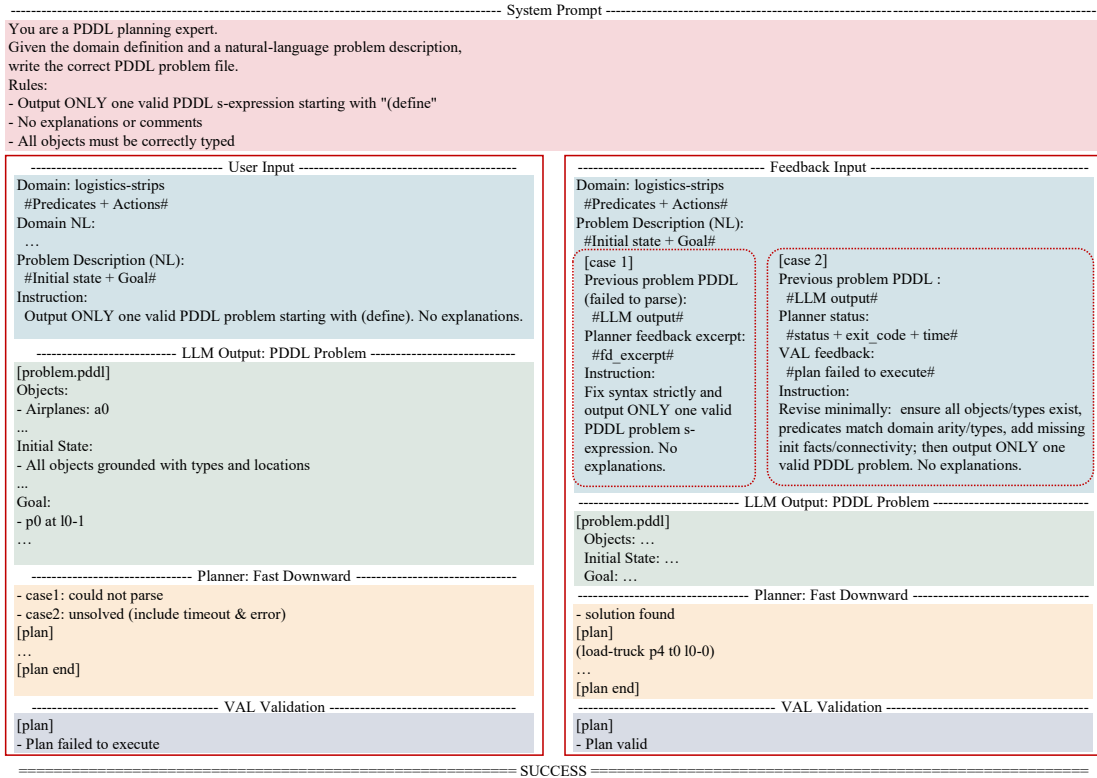


Figure 5: Prompting and planner-in-the-loop feedback process for iterative specification repair and verification.

goal is to convert a natural-language planning request into an executable PDDL problem by iteratively repairing errors revealed by a symbolic toolchain. Concretely, given a domain definition D and a natural-language problem description x , the model is instructed to output exactly one PDDL problem s-expression starting with `(define)`, with no additional commentary. This strict output contract minimizes formatting drift and makes downstream parsing and checking deterministic.

Step 1: One-shot specification generation. We first provide the model with (i) the domain name, (ii) a concise domain summary containing predicates and actions (including arities and types), and (iii) the problem description containing the initial state and goal conditions. The model outputs a candidate PDDL problem instance P . We emphasize type correctness and schema consistency by requiring that all objects be declared with valid domain types, and that each predicate instance matches the declared arity and argument types in D .

Step 2: Symbolic checking and plan synthesis. The generated problem P is fed into a standardized symbolic pipeline. We run a parser/validator to

check well-formedness (e.g., balanced parentheses, well-typed objects, and legal predicate usage). If parsing succeeds, we invoke FD under a fixed configuration to assess solvability and, when possible, to produce a plan π . We then validate π using VAL to ensure that the plan is executable with respect to D and P (i.e., all preconditions hold along the execution trace and the goal is achieved).

Step 3: Feedback construction (two failure cases). When the pipeline fails, we construct feedback that is concise yet actionable and feed it back to the model for minimal repair.

- **Case 1 (Parsing/format failure).** If P cannot be parsed or fails basic well-formedness checks, we return a compact excerpt of the parsing error (e.g., unexpected tokens, missing parentheses, unknown symbols, or type declaration issues). The instruction is to fix syntax strictly and output only one valid PDDL problem. 1103-1110
- **Case 2 (Executable failure).** If P parses but FD reports UNSOLVABLE, times out, or produces a plan that fails VAL validation, we return (a) FD status metadata (e.g., status, exit 1111-1114

code, runtime) and a short log excerpt (e.g., unsatisfied preconditions, unreachable goals), and (b) VAL failure messages when available (e.g., the first violated precondition). The instruction is to revise minimally: ensure all referenced objects/types exist, predicates match domain arities/types, and add missing initial facts (e.g., connectivity) required by reachability constraints.

Step 4: Iterative repair and termination. The model revises the previous problem instance conditioned on the feedback and outputs a new candidate $P^{(t+1)}$. We repeat the symbolic checking and feedback steps until either (i) FD returns a plan and VAL confirms it as valid, or (ii) a maximum number of iterations T_{\max} is reached. In all iterations, we preserve the same strict output contract (single (define) expression; no explanations) to keep the loop machine-checkable and to avoid introducing non-PDDL artifacts.

Rationale. This workflow operationalizes executability alignment by converting opaque generation failures into localized, verifiable repair signals. Parsing feedback targets syntactic correctness, while FD/VAL feedback targets semantic executability (preconditions, effects, reachability, and goal satisfaction). By delegating long-horizon search to a symbolic planner and using diagnostic logs for targeted revisions, the loop turns global plan failure into incremental, formally grounded corrections.

C.3 Consistency Rate Details

This appendix specifies the operational definition of Consistency Rate (CR), including the planner configuration, the cost/length convention, and the plan-structure agreement criterion. CR measures planner-grounded consistency between a generated specification and its reference under a fixed Fast Downward (FD) configuration. It first checks whether the planner outcomes agree (SUCCESS vs. FAILURE); when both sides are solvable, it further checks agreement in plan quality and action-level structure.

C.3.1 Planner configuration and cost definition

For each instance, we solve both the reference specification y_i and the generated specification \hat{y}_i using the same Fast Downward configuration as in the main experiments. The planner outputs are denoted

as $(\text{status}_i^{\text{ref}}, \pi_i^{\text{ref}}, c_i^{\text{ref}})$ and $(\text{status}_i^{\text{gen}}, \pi_i^{\text{gen}}, c_i^{\text{gen}})$. Unless stated otherwise, c denotes the plan length under unit action cost. Thus, cost agreement reduces to length agreement under an identical planner configuration.

C.3.2 Outcome-aware Consistency Rate

Let $\mathcal{D}_{\text{test}} = \{(x_i, y_i)\}_{i=1}^N$ be the test set, and let \hat{y}_i be the model-generated specification. We define Consistency Rate (CR) as

$$\text{CR} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[\text{status}_i^{\text{gen}} = \text{status}_i^{\text{ref}} \wedge \Psi_i]. \quad (10)$$

where Ψ_i enforces additional requirements only when the reference instance is solvable.

$$\Psi_i = \begin{cases} \text{True}, & \text{if } \text{status}_i^{\text{ref}} = \text{FAILURE}, \\ \Phi_i, & \text{if } \text{status}_i^{\text{ref}} = \text{SUCCESS}. \end{cases} \quad (11)$$

The constraint Φ_i requires agreement in both plan cost and plan structure:

$$\Phi_i = \begin{aligned} & |c_i^{\text{gen}} - c_i^{\text{ref}}| \leq \epsilon \\ & \wedge \pi_i^{\text{gen}} \approx \pi_i^{\text{ref}} \end{aligned} \quad (12)$$

Thus, failed instances contribute to CR if and only if the planner outcomes agree, while plan-level comparisons are applied only when both specifications are solvable.

C.3.3 Plan-structure agreement $\pi^{\text{gen}} \approx \pi^{\text{ref}}$

Plans are compared at the level of grounded actions. Let $\text{act}(\pi)$ denote the sequence obtained by serializing each grounded action (e.g., `drive(truck1, locA, locB)`) in execution order.

Sequence similarity. We compute a normalized edit similarity between action-token sequences:

$$\text{Sim}_{\text{edit}} = 1 - \frac{\text{ED}(\text{act}(\pi^{\text{gen}}), \text{act}(\pi^{\text{ref}}))}{\max(|\pi^{\text{gen}}|, |\pi^{\text{ref}}|)} \quad (13)$$

where $\text{ED}(\cdot, \cdot)$ is the Levenshtein edit distance.

Bag-of-actions similarity. To mitigate non-uniqueness arising from commutative or weakly ordered actions, we additionally compute an order-insensitive multiset similarity. Let $\mathcal{M}(\pi)$ be the multiset of grounded actions in π . The multiset Jaccard similarity is defined as

$$\text{Sim}_{\text{bag}} = \frac{\sum_a \min(\mathcal{M}^{\text{gen}}(a), \mathcal{M}^{\text{ref}}(a))}{\sum_a \max(\mathcal{M}^{\text{gen}}(a), \mathcal{M}^{\text{ref}}(a))} \quad (14)$$

Combined criterion (Scheme A). We define a commutativity-robust plan-structure similarity by

$$\text{Sim} = \max(\text{Sim}_{\text{edit}}, \text{Sim}_{\text{bag}}). \quad (15)$$

We consider $\pi^{\text{gen}} \approx \pi^{\text{ref}}$ if

$$\text{Sim} \geq \tau \quad (16)$$

In all experiments, we set $\tau = 0.8$.

C.3.4 Tolerance ϵ for cost/length agreement

We define cost (length) agreement by an absolute tolerance ϵ applied uniformly across domains. Plan costs are considered consistent if

$$|c_i^{\text{gen}} - c_i^{\text{ref}}| \leq \epsilon_i \quad (17)$$

The tolerance ϵ is defined relative to the reference plan length as

$$\epsilon_i = \max(1, \lceil \rho c_i^{\text{ref}} \rceil) \quad (18)$$

where $\rho = 0.05$ in all experiments.

This design allows minor plan-length deviations while still penalizing substantial departures from the reference plan quality under the same planner configuration.

D Appendix: Experiments Result

D.1 Training Environment

We ran our main experiments on a server equipped with six NVIDIA A100 80GB PCIe GPUs. Table 6 summarizes the compute environment.

Table 6: Compute environment for our main experiments.

Component	Specification
GPU	6 × NVIDIA A100 80GB PCIe
CPU	2 × AMD EPYC 9684X (96 cores per socket; 384 logical CPUs in total)
Memory	1.0 TiB RAM
OS	Ubuntu 24.04.1 LTS
Kernel	Linux 6.8.0-88-generic
CUDA	12.8

D.2 Training Dynamics of LoRA Fine-Tuning

As shown in Figure 6, we track the training dynamics of LoRA-based supervised fine-tuning for Qwen3-8B on NL-PDDL-Bench, including the loss curve, learning-rate schedule, and gradient-norm trajectory. The loss exhibits a clear fast-convergence pattern: it starts at approximately

0.18, drops rapidly within a small number of iterations, and falls below 0.01 at around 0.2 epochs. It then enters a stable plateau, fluctuating mildly within roughly 0.004–0.006. This suggests that the model acquires the core structural mapping for NL→PDDL early in training, while later updates primarily refine local constraint details and improve generation stability, yielding diminishing returns from additional epochs.

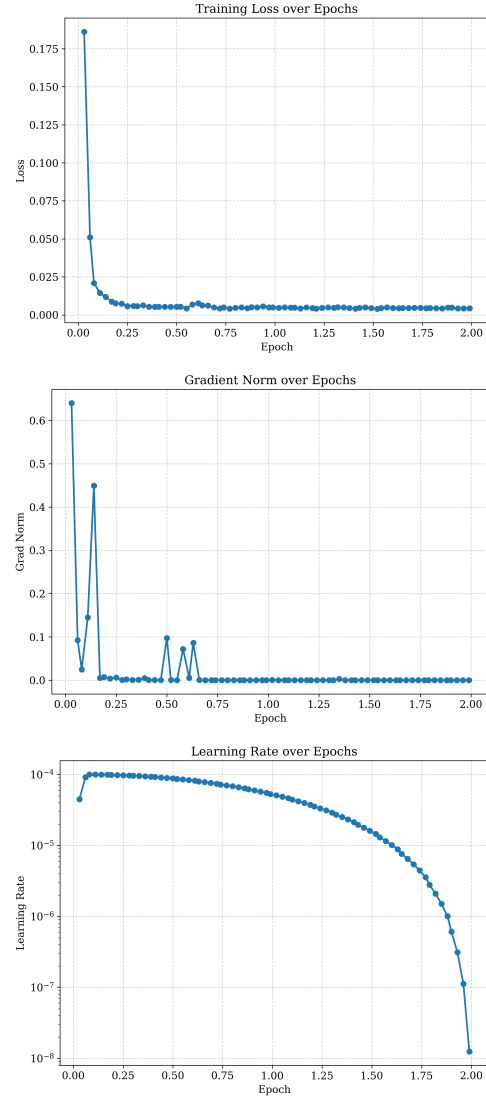


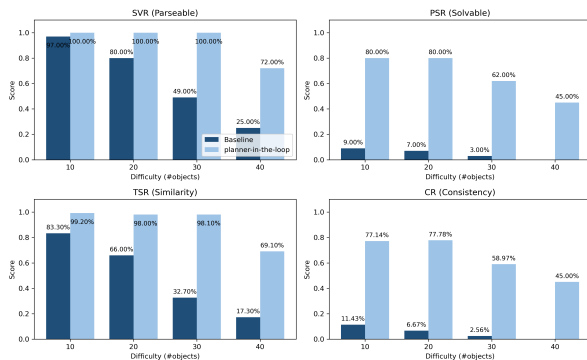
Figure 6: Training dynamics of LoRA fine-tuning for Qwen3-8B on NL-PDDL-Bench.

The learning-rate schedule and gradient-norm evolution further support this interpretation. The learning rate quickly warms up to around 1×10^{-4} , stays relatively high in the early-to-mid stage to encourage effective exploration, and then gradually decays to the order of 1×10^{-7} toward the end to stabilize convergence. Correspondingly, the gra-

1247 dient norm peaks near the beginning and rapidly
 1248 decays toward zero, with only occasional small im-
 1249 pulses and no sustained growth or strong oscilla-
 1250 tions. Overall, these trends indicate stable and well-
 1251 controlled optimization, providing a solid founda-
 1252 tion for subsequently incorporating planner feed-
 1253 back for verification and refinement. Accordingly,
 1254 our experimental focus moves beyond learning sur-
 1255 face syntax templates to evaluating whether sym-
 1256 bolic validation and feedback loops can reliably
 1257 improve executability and plan consistency.

1258 D.3 Case Study: Qwen2.5-7B

1259 Figure 7 presents a case study in the Rovers domain,
 1260 comparing Qwen2.5-7B under Baseline and our
 1261 full framework across four object-count difficulty
 1262 levels (10, 20, 30, and 40 objects) and four eval-
 1263 uation metrics (SVR, PSR, TSR, and CR). Rovers
 1264 is representative of domains where correctness is
 1265 determined by the closure of interacting constraints,
 1266 including typed object declarations, consistent pred-
 1267 icate grounding, and reachability-supporting facts
 1268 that enable long-horizon state transitions. As the
 1269 instance size grows, small omissions in the speci-
 1270 fication tend to propagate into global failures, be-
 1271 cause later actions require precise causal support
 1272 and resource feasibility that cannot be recovered by
 1273 local plausibility alone.



1274 Figure 7: Case study: Qwen2.5-7B on Rovers across
 1275 difficulty levels.

1276 Under the Baseline setting, performance deteri-
 1277 orates sharply with scaling. While many outputs
 1278 remain syntactically parseable at smaller scales, the
 1279 gap between surface well-formedness and planner-
 1280 grounded executability widens as constraints accu-
 1281 mulate. This pattern suggests that baseline er-
 1282 rors increasingly concentrate on semantic structure
 1283 rather than formatting, such as missing enabling
 facts required for navigation and communication,
 inconsistent use of predicate arguments and types,

1284 or incomplete initialization of resources and loca-
 1285 tions. Such defects typically do not prevent parsing,
 1286 yet they break reachability or precondition satis-
 1287 faction, leading the planner to return UNSOLVABLE
 1288 instances and reducing plan-level agreement with
 1289 the reference specification. At the largest scale, the
 1290 baseline often fails to maintain global constraint
 1291 closure, which manifests as a collapse in solvabil-
 1292 ity and consistency even when fragments of the
 1293 specification appear locally reasonable.

1294 In contrast, our framework mitigates this scaling-
 1295 induced brittleness by coupling generation with ver-
 1296 ifiable symbolic feedback. The validator and plan-
 1297 ner expose concrete failure evidence, such as type or
 1298 signature mismatches and unsatisfied preconditions,
 1299 which we convert into targeted revision directives.
 1300 This forces the model to repair missing causal links
 1301 and reconcile constraint conflicts, rather than iterat-
 1302 ing on surface templates. As a result, the framework
 1303 maintains substantially stronger executable align-
 1304 ment as difficulty increases, preserving solvability
 1305 and plan-level consistency under tighter constraint
 1306 coupling. The remaining degradation at the largest
 1307 scale is consistent with the increased density of
 1308 interacting constraints, where a single missed sup-
 1309 port fact can invalidate multiple downstream ac-
 1310 tion chains. Overall, the case study indicates that
 1311 the gains arise from improved semantic closure un-
 1312 der fixed domain dynamics, not merely from better
 1313 formatting, and that planner-grounded diagnostics
 1314 provide an effective mechanism for turning global
 1315 planning failures into localized, correctable specifi-
 1316 cation edits.

1317 D.4 Cross-Domain Generalization

1318 We further study cross-domain generalization to de-
 1319 termine whether the observed gains reflect domain-
 1320 specific memorization or a transferable improve-
 1321 ment in executable alignment. Using the same
 1322 evaluation protocol as in the main experiments, we
 1323 compare the Baseline and the planner-in-the-loop
 1324 setting on each domain and examine how improve-
 1325 ments distribute across the four complementary
 1326 metrics. Because domains differ in their dominant
 1327 constraint structures, a robust method should pri-
 1328 marily improve planner-grounded criteria, namely
 1329 planner success and outcome-aware plan-level con-
 1330 sistency, rather than only increasing surface well-
 1331 formedness.

1332 For space reasons, the main paper presents the
 1333 domain-wise delta plot in a single-column figure
 1334 to fit the layout constraints. Here we include an en-

larged cross-column version (Figure 8) to improve readability and enable more careful inspection of domain-level patterns.

The results show that improvements are concentrated on planner-grounded correctness. Across most domains, changes in syntax validity are comparatively modest, while planner success and plan-level consistency increase substantially, indicating that the main benefit comes from repairing semantic defects that remain invisible to purely syntactic checks. This pattern is consistent with the role of planner-in-the-loop feedback, which targets constraint closure, reachability prerequisites, and resource consistency, and therefore reduces failures caused by missing initial facts, incompatible predicate usage, or latent precondition violations.

We also observe meaningful domain variation that helps interpret when planner feedback is most effective. Transportation-style domains, where failures often arise from omitted connectivity facts or incomplete state constraints, typically exhibit large gains in planner success and consistency once diagnostics are converted into localized revisions. In contrast, in domains with tighter resource coupling or more brittle modeling requirements, improvements in syntax validity do not necessarily translate into comparable gains in solvability, which suggests that the remaining errors are less local and may require richer semantic modeling or stronger priors to resolve. Overall, the domain-wise analysis supports the conclusion that planner-in-the-loop feedback improves executable alignment in a transferable manner, while also revealing constraint regimes where conservative local repair is insufficient.

D.5 Results of Mainstream Models on NL-PDDL-Bench

Tables 7 and 8 summarize domain-wise and difficulty-wise SVR/PSR/TSR/CR for mainstream models under the Baseline setting and our planner-in-the-loop framework. The results highlight a persistent gap between surface well-formedness and planner-grounded correctness: many baselines achieve non-trivial SVR yet frequently yield unsolvable instances, especially in domains with tight structural constraints and long-range dependencies, and the gap widens as object scale increases. Importantly, our CR is outcome-aware: it first measures whether the planner outcomes of the generated and reference instances match (SUCCESS vs. FAILURE); plan-cost and plan-structure agreement are only evaluated when the reference instance is solvable

(and thus both sides succeed). Therefore, it is possible to observe high CR together with low PSR in subsets where the references are predominantly unsolvable, because consistent FAILURE outcomes contribute to CR while PSR only counts SUCCESS of the generated instances. By combining multi-stage symbolic validation (parsing, planning, and execution verification) with evidence-driven minimal repairs, planner-in-the-loop converts previously silent semantic inconsistencies into diagnosable error types and yields broad gains in executability-aligned metrics, with particularly large improvements in PSR and in CR on solvable references, indicating better satisfaction of domain-critical structures such as preconditions, reachability/connectivity facts, and resource constraints. Moreover, the smoother degradation of planner-in-the-loop under scaling suggests that planner-in-the-loop feedback provides an explicit correctness signal that helps preserve global consistency as constraints accumulate, improving reliability beyond what can be obtained by syntax compliance or template memorization alone.

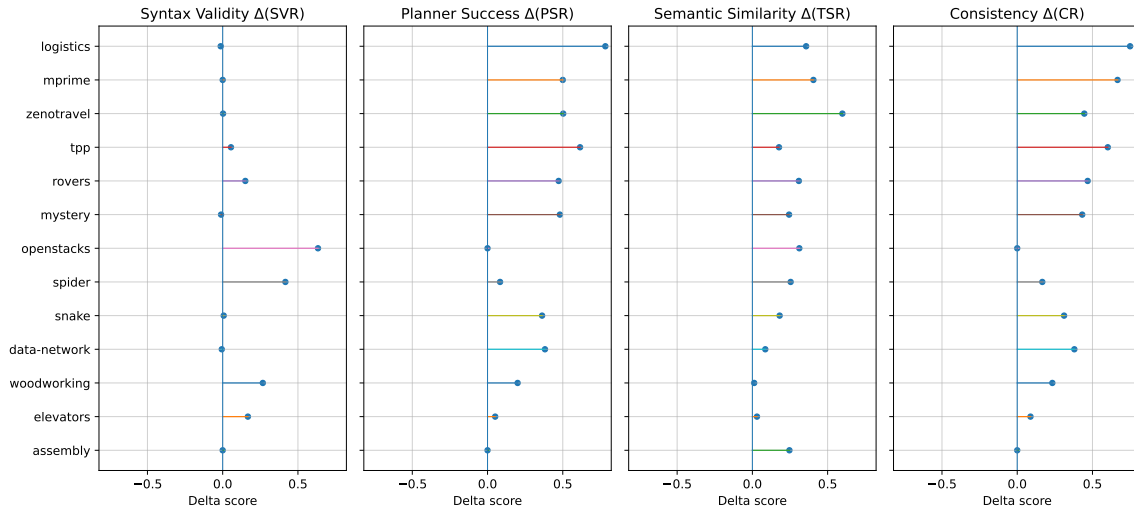


Figure 8: Cross-domain generalization results on NL-PDDL-Bench. This figure provides an enlarged cross-column view to facilitate readability and detailed comparison across planning domains.

Table 7: Domain-wise results on NL-PDDL-Bench. We report SVR/PSR/TSR/CR for each domain under Baseline and planner-in-the-loop.

Domain	Baseline				planner-in-the-loop			
	SVR	PSR	TSR	CR	SVR	PSR	TSR	CR
DeepSeek-Coder-V2-16B								
Assembly	100.0%	0.0%	23.5%	100.0%	100.0%	0.0%	37.3%	100.0%
Data Network	0.0%	0.0%	0.0%	0.0%	10.0%	0.0%	0.0%	0.0%
Elevators	2.0%	0.0%	1.9%	0.0%	12.0%	0.0%	2.3%	0.0%
Logistics	40.0%	6.0%	34.4%	5.6%	38.0%	38.0%	37.6%	37.6%
MPrime	100.0%	0.0%	56.7%	0.0%	100.0%	100.0%	100.0%	100.0%
Mystery	100.0%	84.0%	93.6%	84.2%	100.0%	100.0%	100.0%	100.0%
Openstacks	0.0%	0.0%	0.0%	0.0%	60.0%	0.0%	0.0%	0.0%
Rovers	1.0%	0.0%	0.4%	0.0%	14.0%	11.0%	11.8%	11.3%
Snake	41.0%	0.0%	35.3%	0.0%	46.0%	44.0%	43.6%	43.6%
Spider	0.0%	0.0%	0.0%	0.0%	100.0%	0.0%	0.0%	0.0%
Traveling Purchaser Problem	13.0%	0.0%	9.6%	0.0%	18.0%	11.0%	10.7%	10.7%
Woodworking	0.0%	0.0%	0.0%	0.0%	80.0%	0.0%	0.0%	0.0%
ZenoTravel	100.0%	15.0%	15.5%	15.2%	100.0%	100.0%	100.0%	100.0%
GPT-OSS-20B								
Assembly	100.0%	0.0%	38.0%	100.0%	100.0%	0.0%	10.3%	100.0%
Data Network	88.0%	0.0%	57.5%	0.0%	69.0%	5.0%	17.3%	5.1%
Elevators	10.0%	2.0%	3.6%	1.6%	66.0%	2.0%	3.2%	2.2%
Logistics	90.0%	60.0%	72.7%	60.0%	89.0%	58.0%	63.2%	59.2%
MPrime	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%	0.0%	0.0%
Mystery	95.0%	89.0%	88.9%	89.5%	95.0%	74.0%	72.2%	73.7%
Openstacks	0.0%	0.0%	0.0%	0.0%	100.0%	0.0%	39.9%	0.0%
Rovers	52.0%	18.0%	19.9%	17.6%	59.0%	22.0%	14.4%	22.0%
Snake	46.0%	5.0%	28.0%	5.1%	54.0%	8.0%	12.6%	7.7%
Spider	100.0%	0.0%	46.3%	0.0%	0.0%	0.0%	0.0%	0.0%
Traveling Purchaser Problem	97.0%	30.0%	53.0%	30.0%	86.0%	38.0%	43.6%	37.6%
Woodworking	80.0%	0.0%	50.7%	0.0%	60.0%	0.0%	25.7%	0.0%
ZenoTravel	97.0%	96.0%	48.8%	96.5%	100.0%	86.0%	44.8%	86.3%

Table 7: Domain-wise results on NL-PDDL-Bench (continued).

Domain	Baseline				planner-in-the-loop			
	SVR	PSR	TSR	CR	SVR	PSR	TSR	CR
GLM-4.1V-9B								
Assembly	100.0%	0.0%	9.5%	100.0%	100.0%	0.0%	62.4%	100.0%
Data Network	92.0%	0.0%	25.2%	0.0%	100.0%	0.0%	94.9%	0.0%
Elevators	10.0%	0.0%	1.4%	0.0%	21.0%	0.0%	7.8%	0.0%
Logistics	93.0%	1.0%	24.1%	0.8%	100.0%	100.0%	99.9%	100.0%
MPrime	100.0%	0.0%	0.0%	0.0%	100.0%	100.0%	100.0%	100.0%
Mystery	100.0%	21.0%	26.3%	21.1%	100.0%	100.0%	100.0%	100.0%
Openstacks	0.0%	0.0%	0.0%	0.0%	100.0%	0.0%	50.1%	0.0%
Rovers	74.0%	8.0%	18.5%	8.2%	75.0%	57.0%	72.4%	57.2%
Snake	90.0%	0.0%	12.2%	10.3%	97.0%	92.0%	97.3%	92.3%
Spider	100.0%	0.0%	0.0%	0.0%	100.0%	50.0%	86.5%	50.0%
Traveling Purchaser Problem	93.0%	3.0%	15.4%	2.5%	94.0%	86.0%	92.1%	85.8%
Woodworking	80.0%	0.0%	59.5%	0.0%	80.0%	0.0%	54.2%	0.0%
ZenoTravel	99.0%	39.0%	45.8%	38.6%	100.0%	96.0%	100.0%	97.5%
DeepSeek-R1-7B								
Assembly	100.0%	0.0%	28.4%	100.0%	100.0%	0.0%	41.7%	100.0%
Data Network	90.0%	8.0%	32.5%	7.8%	94.0%	38.0%	63.1%	37.6%
Elevators	18.0%	0.0%	6.3%	0.0%	24.0%	6.0%	12.1%	5.8%
Logistics	92.0%	35.0%	58.9%	36.4%	100.0%	92.0%	97.6%	91.4%
MPrime	100.0%	12.0%	41.8%	11.6%	100.0%	68.0%	88.9%	66.7%
Mystery	100.0%	54.0%	62.3%	53.2%	100.0%	96.0%	97.1%	95.8%
Openstacks	0.0%	0.0%	0.0%	0.0%	72.0%	0.0%	36.4%	0.0%
Rovers	78.0%	42.0%	47.9%	41.6%	96.0%	83.0%	88.7%	82.4%
Snake	84.0%	21.0%	35.4%	20.8%	93.0%	64.0%	72.8%	63.7%
Spider	100.0%	0.0%	19.6%	0.0%	100.0%	22.0%	51.2%	21.5%
Traveling Purchaser Problem	86.0%	19.0%	39.1%	18.7%	92.0%	61.0%	75.4%	60.3%
Woodworking	72.0%	4.0%	28.7%	3.9%	88.0%	32.0%	55.6%	31.8%
ZenoTravel	98.0%	79.0%	61.2%	78.4%	100.0%	97.0%	98.3%	96.8%
Llama3.1-8B								
Assembly	100.0%	0.0%	57.7%	100.0%	100.0%	0.0%	93.1%	100.0%
Data Network	95.0%	0.0%	89.2%	0.0%	100.0%	100.0%	100.0%	100.0%
Elevators	14.0%	0.0%	5.4%	0.0%	18.0%	8.0%	8.5%	8.2%
Logistics	88.0%	5.0%	47.0%	5.6%	100.0%	100.0%	100.0%	99.2%
MPrime	100.0%	0.0%	100.0%	0.0%	100.0%	100.0%	100.0%	100.0%
Mystery	100.0%	0.0%	91.1%	0.0%	100.0%	100.0%	100.0%	100.0%
Openstacks	20.0%	0.0%	0.0%	0.0%	100.0%	0.0%	57.4%	0.0%
Rovers	69.0%	21.0%	54.9%	21.4%	100.0%	99.0%	99.5%	98.7%
Snake	97.0%	0.0%	49.3%	5.1%	97.0%	92.0%	97.4%	92.3%
Spider	100.0%	0.0%	59.7%	0.0%	100.0%	0.0%	91.5%	0.0%
Traveling Purchaser Problem	97.0%	20.0%	90.2%	19.8%	100.0%	100.0%	100.0%	99.0%
Woodworking	100.0%	0.0%	84.6%	20.0%	100.0%	80.0%	98.9%	80.0%
ZenoTravel	100.0%	9.0%	48.2%	9.1%	100.0%	98.0%	100.0%	98.5%

Table 7: Domain-wise results on NL-PDDL-Bench (continued).

Domain	Baseline				planner-in-the-loop			
	SVR	PSR	TSR	CR	SVR	PSR	TSR	CR
Qwen2.5-7B								
Assembly	100.0%	0.0%	54.7%	100.0%	100.0%	0.0%	60.1%	100.0%
Data Network	100.0%	0.0%	86.8%	0.0%	93.0%	42.0%	92.0%	40.7%
Elevators	15.0%	0.0%	6.8%	0.5%	12.0%	0.0%	6.3%	2.7%
Logistics	94.0%	0.0%	76.4%	0.0%	100.0%	92.0%	99.9%	87.2%
MPrime	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%
Mystery	100.0%	79.0%	94.8%	78.9%	100.0%	95.0%	99.5%	94.7%
Openstacks	0.0%	0.0%	0.0%	0.0%	40.0%	0.0%	33.0%	0.0%
Rovers	62.0%	4.0%	49.4%	5.0%	93.0%	67.0%	91.0%	64.8%
Snake	97.0%	0.0%	76.7%	12.8%	97.0%	69.0%	89.7%	71.8%
Spider	100.0%	0.0%	80.8%	0.0%	100.0%	0.0%	70.4%	50.0%
Traveling Purchaser Problem	100.0%	3.0%	96.6%	3.0%	99.0%	88.0%	99.2%	83.8%
Woodworking	100.0%	0.0%	83.3%	0.0%	100.0%	0.0%	75.8%	20.0%
ZenoTravel	99.0%	96.0%	63.7%	95.9%	100.0%	81.0%	96.5%	79.2%
Qwen3-8B								
Assembly	100.0%	0.0%	61.6%	100.0%	100.0%	0.0%	75.0%	100.0%
Data Network	100.0%	0.0%	100.0%	0.0%	100.0%	100.0%	100.0%	100.0%
Elevators	11.0%	0.0%	9.6%	0.0%	16.0%	6.0%	8.2%	6.0%
Logistics	100.0%	0.0%	62.7%	0.0%	100.0%	100.0%	100.0%	100.0%
MPrime	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%
Mystery	100.0%	100.0%	98.4%	100.0%	89.0%	89.0%	89.5%	89.5%
Openstacks	0.0%	0.0%	0.0%	0.0%	80.0%	0.0%	46.7%	0.0%
Rovers	97.0%	20.0%	75.0%	20.1%	99.0%	97.0%	98.4%	97.5%
Snake	90.0%	0.0%	83.7%	0.0%	97.0%	97.0%	97.4%	97.4%
Spider	50.0%	0.0%	47.2%	0.0%	100.0%	0.0%	91.9%	0.0%
Traveling Purchaser Problem	99.0%	0.0%	86.1%	0.0%	93.0%	93.0%	93.4%	93.4%
Woodworking	100.0%	0.0%	94.3%	0.0%	100.0%	40.0%	100.0%	40.0%
ZenoTravel	100.0%	99.0%	54.9%	99.5%	100.0%	100.0%	100.0%	100.0%
Gemma3-4B								
Assembly	100.0%	0.0%	10.7%	100.0%	100.0%	0.0%	37.5%	100.0%
Data Network	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Elevators	4.0%	0.0%	0.0%	0.5%	96.0%	1.0%	1.1%	21.2%
Logistics	71.0%	2.0%	19.0%	8.0%	38.0%	34.0%	37.6%	30.4%
MPrime	100.0%	0.0%	0.0%	0.0%	100.0%	0.0%	100.0%	100.0%
Mystery	100.0%	32.0%	58.5%	42.1%	100.0%	84.0%	100.0%	63.2%
Openstacks	0.0%	0.0%	0.0%	0.0%	20.0%	0.0%	0.0%	0.0%
Rovers	0.0%	0.0%	0.0%	0.0%	11.0%	7.0%	9.7%	6.3%
Snake	56.0%	0.0%	0.0%	30.8%	44.0%	28.0%	43.6%	25.6%
Spider	0.0%	0.0%	0.0%	0.0%	100.0%	0.0%	0.0%	0.0%
Traveling Purchaser Problem	16.0%	0.0%	5.8%	6.1%	58.0%	8.0%	11.0%	16.8%
Woodworking	0.0%	0.0%	0.0%	0.0%	80.0%	0.0%	0.0%	20.0%
ZenoTravel	99.0%	0.0%	10.6%	25.4%	100.0%	80.0%	100.0%	72.1%

Table 8: Difficulty-wise results on NL-PDDL-Bench. We report SVR/PSR/TSR/CR under Baseline and planner-in-the-loop for each difficulty level (10/20/30/40).

Difficulty Level	Baseline				planner-in-the-loop			
	SVR	PSR	TSR	CR	SVR	PSR	TSR	CR
Gemma3-4B								
Level 1 (10)	58.0%	3.0%	17.3%	22.3%	65.0%	41.0%	54.6%	40.2%
Level 2 (20)	28.0%	0.0%	3.8%	3.6%	49.0%	26.0%	26.8%	26.8%
Level 3 (30)	39.0%	0.0%	2.4%	8.0%	64.0%	22.0%	28.8%	28.4%
Level 4 (40)	23.0%	0.0%	3.6%	6.4%	64.0%	14.0%	18.9%	22.1%
GLM-4.1V-9B								
Level 1 (10)	92.0%	12.0%	24.8%	16.3%	97.0%	73.0%	93.5%	75.7%
Level 2 (20)	74.0%	10.0%	25.6%	9.6%	80.0%	65.0%	77.0%	65.2%
Level 3 (30)	69.0%	12.0%	22.2%	12.0%	77.0%	65.0%	72.2%	65.2%
Level 4 (40)	67.0%	6.0%	11.6%	6.4%	67.0%	51.0%	61.0%	50.6%
Llama3.1-8B								
Level 1 (10)	94.0%	13.0%	65.0%	17.1%	97.0%	91.0%	96.5%	93.6%
Level 2 (20)	70.0%	10.0%	52.2%	10.4%	81.0%	77.0%	78.0%	76.4%
Level 3 (30)	71.0%	7.0%	45.1%	6.8%	80.0%	78.0%	78.2%	77.6%
Level 4 (40)	70.0%	9.0%	49.2%	9.2%	81.0%	77.0%	78.6%	77.1%
Qwen2.5-7B								
Level 1 (10)	98.0%	24.0%	80.2%	30.7%	97.0%	69.0%	91.5%	72.5%
Level 2 (20)	76.0%	19.0%	61.4%	18.8%	77.0%	58.0%	75.6%	55.2%
Level 3 (30)	69.0%	26.0%	52.8%	26.0%	78.0%	62.0%	76.2%	60.8%
Level 4 (40)	63.0%	18.0%	50.4%	17.7%	75.0%	60.0%	71.9%	59.0%
Qwen3-8B								
Level 1 (10)	97.0%	28.0%	78.9%	31.1%	95.0%	88.0%	93.9%	91.2%
Level 2 (20)	77.0%	20.0%	59.1%	20.4%	82.0%	78.0%	77.7%	77.6%
Level 3 (30)	78.0%	30.0%	53.5%	29.6%	78.0%	77.0%	77.1%	76.8%
Level 4 (40)	77.0%	21.0%	54.4%	21.3%	76.0%	73.0%	74.4%	73.1%
DeepSeek-R1-7B								
Level 1 (10)	94.0%	45.0%	48.6%	46.3%	96.0%	78.0%	81.4%	79.6%
Level 2 (20)	86.0%	34.0%	38.9%	35.2%	91.0%	63.0%	65.7%	62.4%
Level 3 (30)	82.0%	27.0%	31.4%	26.8%	89.0%	51.0%	54.1%	50.8%
Level 4 (40)	76.0%	18.0%	24.6%	18.9%	85.0%	41.0%	45.3%	40.7%
DeepSeek-Coder-V2-16B								
Level 1 (10)	53.0%	8.0%	29.1%	10.8%	62.0%	53.0%	56.6%	56.6%
Level 2 (20)	27.0%	3.0%	10.3%	2.8%	31.0%	27.0%	26.8%	26.8%
Level 3 (30)	30.0%	4.0%	6.6%	4.4%	36.0%	29.0%	28.8%	28.8%
Level 4 (40)	19.0%	6.0%	6.3%	6.4%	24.0%	19.0%	18.9%	18.9%