

CRISP: Complex Reasoning with Interpretable Step-based Plans

Anonymous ACL submission

Abstract

Recent advancements in large language models (LLMs) underscore the need for stronger reasoning capabilities to solve complex problems effectively. While Chain-of-Thought (CoT) reasoning has been a step forward, it remains insufficient for many domains. A promising alternative is explicit high-level plan generation, but existing approaches largely assume that LLMs can produce effective plans through few-shot prompting alone, without additional training. In this work, we challenge this assumption and introduce CRISP (Complex Reasoning with Interpretable Step-based Plans), a multi-domain dataset of high-level plans for mathematical reasoning and code generation. The plans in CRISP are automatically generated and rigorously validated—both intrinsically, using an LLM as a judge, and extrinsically, by evaluating their impact on downstream task performance. We demonstrate that fine-tuning a small model on CRISP enables it to generate higher-quality plans than much larger models using few-shot prompting, while significantly outperforming Chain-of-Thought reasoning. Furthermore, our out-of-domain evaluation reveals that fine-tuning on one domain improves plan generation in the other, highlighting the generalizability of learned planning capabilities.

1 Introduction

Large language models (LLMs) abilities advance rapidly in logical reasoning, code generation, and mathematical problem-solving (Plaat et al., 2024; Jiang et al., 2024a; Ahn et al., 2024). A key factor behind recent breakthroughs is the ability of LLMs to break down complex tasks into manageable steps—an approach exemplified by chain-of-thought prompting (Wei et al., 2022a).

While chain-of-thought reasoning has led to notable performance gains, it remains prone to errors such as missing intermediate steps and semantic misunderstandings (Wei et al., 2022a; Jiang et al., 2024b). To address these challenges, recent studies have explored prompting strategies that explicitly break down problems into subtasks (Dua et al., 2022; Zhou et al., 2023; Khot et al., 2023; Prasad et al., 2024; Ding et al., 2024). One prominent approach is self-generating a high-level plan. by the LLM before executing the task. This plan-and-solve approach lead to significant improvements in mathematical, commonsense, and symbolic reasoning, as well as code generation (Wang et al., 2023a; Jiang et al., 2024b). However, the self-generated plans were only partially effective, as they did not match the performance of ground-truth planning across various downstream tasks.

In this work, we argue that generating high-quality high-level plans is a hard challenge for LLMs, as current models often struggle to produce accurate and effective plans across different domains. Part of this difficulty arises from the scarcity of explicit planning data, as people rarely externalize or document their high-level reasoning in a structured way. To address this gap, we introduce CRISP (Complex Reasoning with Interpretable Step-based Plans), a novel dataset designed to enhance high-level planning capabilities. CRISP spans two domains: mathematics and code generation—where solutions naturally decompose into structured, high-level steps. These plans were derived from annotated detailed solutions of Magpie-Reasoning-V1-150K (Xu et al., 2025) and validated both extrinsically, by measuring their impact on the original task performance, and intrinsically, using LLM-based judgment to assess coherence, conciseness, clarity, and

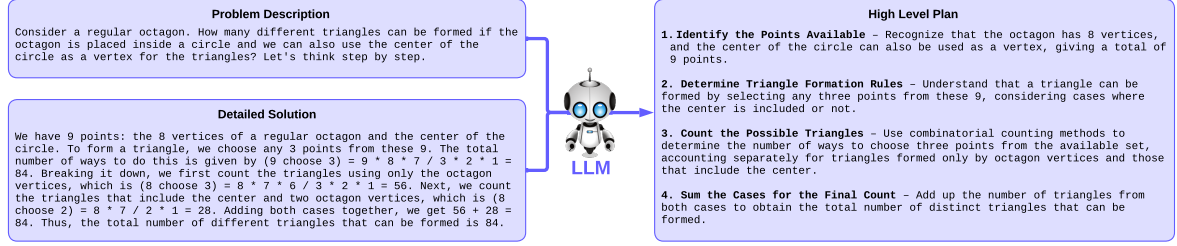


Figure 1: Example from the Math domain showing a problem statement, its detailed solution, and the generated high-level plan. The LLM was prompted to retain the high-level strategy while omitting lower-level details. It was also provided with a few examples and guidelines tailored to the specific domain.

completeness.

To assess CRISP’s impact, we perform a series of experiments on four reasoning-related datasets: MBPP (Austin et al., 2021), HumanEval (Chen et al., 2021), GSM8K (Cobbe et al., 2021), and MATH (Hendrycks et al., 2021). Our findings indicate that while larger models generate plans that lead to better performance on the reasoning-related benchmarks, a small model can surpass them with a lightweight fine-tuning on CRISP using LoRA (Hu et al., 2022). Furthermore, with lightweight fine-tuning with LoRA (Hu et al., 2022) on our dataset, LLMs exhibit substantial improvements in plan generation, as reflected in both higher performance on structured reasoning benchmarks such as MBPP (Austin et al., 2021), HumanEval (Chen et al., 2021), and higher quality scores assigned by LLM-based evaluations.

Additionally, we assess the out-of-domain generalization of planning abilities and find that CRISP effectively transfers these capabilities across different tasks and domains. For example, a model fine-tuned on the Math domain achieves a pass@1 score of 84.6 on the HumanEval code dataset—only 0.4 points lower than the same model fine-tuned on the Coding and Debugging domain. This strong transferability highlights CRISP’s potential for seamless integration into existing training pipelines, where it can enhance LLM reasoning abilities and improve performance across diverse downstream tasks.

The significance of this work lies in its emphasis on high-level planning as a beneficial trainable capability that current off-the-shelf models do not excel in. By providing LLMs with explicit fine-tuning on high-level planning,

we enhance their ability to decompose tasks which in turn improves their applicability to real-world scenarios requiring robust, domain-agnostic reasoning. The dataset is publicly available [here](#).

Our contributions are threefold:

- We show that generating high-level plans is a **challenging task** for LLMs, yet highly beneficial.
- We introduce the **CRISP dataset**, a multi-domain dataset of high-level plans derived from annotated detailed solutions.
- We demonstrate that even a short LoRA-based **fine-tuning improves** the quality of generated plans, significantly outperforming larger models.
- We show that high-level **planning abilities transfer well** between tasks and domains through out-of-domain evaluation.

2 Related Work

A growing body of research investigates methods to enhance LLMs’ performance on tasks requiring multi-step reasoning and structured problem-solving. One prominent approach is Chain of Thought (CoT) reasoning (Wei et al., 2022b), which improves LLMs’ ability to handle complex tasks by explicitly breaking down reasoning into intermediate steps. This method has significantly improved performance on arithmetic, commonsense, and symbolic reasoning benchmarks.

Building on CoT, several works have introduced further improvements, such as exploring multiple reasoning trajectories (Wang et al., 2023b), backtracking and applying

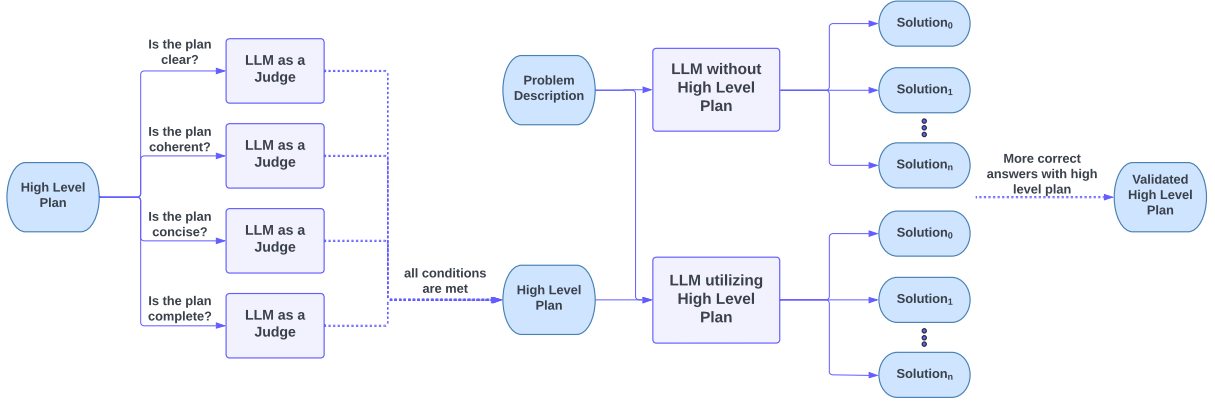


Figure 2: The validation and filtering pipeline of CRISP. Each generated high-level plan undergoes binary validation, where its clarity, coherence, conciseness, and completeness are assessed by the LLM. If all attributes receive a positive judgment, the plan is then externally validated by comparing the solutions generated with and without the high-level plan in the prompt, filtering out those that do not lead to improved performance on the original task.

search algorithms (Yao et al., 2023; Hao et al., 2023; Besta et al., 2024; Zhou et al., 2024), self-evaluation (Xie et al., 2023), self-reflection (Shinn et al., 2023), and self-refinement (Madaan et al., 2023). These methods enable models to explore multiple solution paths, iteratively improving their responses.

Another line of work focuses on explicit task decomposition, where problems are broken down into structured subtasks. Some approaches generate successive questions (Dua et al., 2022; Zhou et al., 2023), others decompose problems into lines of code (Chen et al., 2023; Yang et al., 2023; Ding et al., 2024), and some define explicit hierarchical subtasks (Khot et al., 2023). Most similar to our work are "Plan-and-Execute" approaches, where an LLM first generates a structured plan before executing it to solve a problem (Wang et al., 2023a; Jiang et al., 2024b). Additionally, Prasad et al. (2024) proposed iteratively refining plans upon execution failures. However, these methods have largely been evaluated in single-domain settings, and they treat high-quality plan generation as an emergent ability, requiring no additional training. In contrast, our work systematically explores plan generation across multiple domains, demonstrating both the challenges of this task and the tangible benefits of training LLMs to produce structured plans. However, these approaches are primarily benchmarked within a single domain and consider the task of generating a high-quality plan as an emergent capability requiring no additional training. In

contrast, our work systematically explores plan generation across multiple domains, demonstrating both the challenges of this task and the benefits of training LLMs to generate structured plans.

Several datasets incorporate task decomposition, including those designed for web agents (Liu et al., 2018; Yao et al., 2022; Deng et al., 2023; Shi et al., 2017), household activities (Puig et al., 2018; Shridhar et al., 2021, 2020), games (Guss et al., 2019; Prasad et al., 2024), and robotics (Kannan et al., 2024; Zhang et al., 2023; Li et al., 2023). These datasets primarily focus on dynamic decision-making in interactive environments, where planning is contingent on real-time feedback and reinforcement learning. In contrast, CRISP is designed for structured, multi-domain task decomposition, emphasizing myopic problems—tasks that can be solved through a predefined sequence of steps rather than adaptive decision-making.

3 Dataset Collection

In this section we elaborate on the creation of CRISP. In Section 3.1 we describe how we generated the high-level plans based on careful prompt engineering. Then, in Section 3.2 we describe the validation and filtering mechanisms we developed to ensure the quality of the generated high-level plans based on both intrinsic and extrinsic evaluations. Finally, in Section 3.3, we analyze how the filtering and validation affect downstream tasks and the scores assigned to the plans based on LLM-

based judgment.

3.1 High-Level Plan Generation

Our high-level plan dataset is derived from Magpie-Reasoning-V1-150K (Xu et al., 2025), which is licensed under ‘Llama3’ and spans two domains: *Math* and *Coding and Debugging*. Magpie reasoning examples in the domains of math and coding and debugging illustrate how models can decompose complex problems into structured subproblems. In math, this involves breaking down equations or proofs into intermediate steps, while in coding and debugging, it includes identifying error patterns, generating hypotheses about potential bugs, and testing fixes iteratively. Each instance includes a problem statement and a detailed solution, generated by *Qwen2-72B-Instruct* for math and *Llama-3-70B-Instruct* for coding. An example of such a problem statement and its corresponding detailed solution is depicted on the left side of Figure 1. The dataset encompasses a wide range of topics from mathematics and coding such as geometry, algebra, and integrals, differential equations, and probability in mathematics, as well as data structures, algorithm design, syntax and logic error fixes, concurrency, and general software engineering in coding and debugging. In total, we extract 74,225 math-related samples and 66,342 coding-related samples. Since the problems in this dataset are myopic—meaning they can often be solved using a predefined sequence of steps—we believe that generating a high-level plan should be particularly beneficial.

For each problem and its detailed solution, we used *Mixtral-8x22B-Instruct-v0.1* to generate a high-level plan through few-shot prompting. We selected this model after manually evaluating its generated plans and finding them to be of higher quality than those from other LLMs, along with its permissive license. An example of an such a high-level plans is depicted on the right side of Figure 1. In the plan generation prompt, we instructed the model to outline the high-level logical strategy while abstracting away implementation details. This approach ensures that information from the detailed solution, which the model should not have prior knowledge of, remains undisclosed. Simultaneously, it preserves a degree of flexibility, allowing the model to determine the precise

method for executing each step at a later stage. The full generation prompt of Math is provided in Appendix A.3.

3.2 Filtering and Validation

After gathering a substantial collection of high-level plans, we implemented a two-step filtering process to validate the plans intrinsically and extrinsically, as illustrated in Figure 2. First, we apply ‘LLM as a Judge’ with *Llama-3.1-70B-Instruct* to determine whether the generated plans are **concise, clear, coherent, and complete**. These four attributes are essential for ensuring that a plan is described efficiently without redundancy or repetition (concise), is easily understandable without ambiguity or vague language (clear), follows a logical sequence without missing critical transitions (coherent), and includes all the essential steps to fully address the problem and derive the solution (complete). We believe that ensuring these four attributes is crucial for generating high-quality plans that are both interpretable and actionable, facilitating their usefulness in various downstream tasks. Each attribute is assessed using a binary judgment, and any plan that fails to meet one or more criteria is discarded. In this step 7,412 math plans and 5,592 code plans were filtered out, accounting for approximately 9% of the original dataset. For the full prompt used in this evaluation, refer to Appendix A.4.

After the intrinsic validation, we assessed the generated plans by testing their impact on the model’s ability to successfully complete the original tasks. To that end, we generated 10 solutions both with and without the plan using *Llama-3.1-70B-Instruct* and discarded cases where the number of correct final answers was higher without the plan. This step removes an additional 1,089 math plans and 4,612 code plans that—while clear, concise, coherent, and complete—failed to produce more correct answers than when the LLM was not provided with a high-level plan.

After filtering, we retain 65,800 math plans and 56,200 coding plans. Table 1 summarizes the final dataset.

3.3 Dataset Analysis

Empirical validation shows that despite reducing the dataset size, each filtering step

DOMAIN	# EXAMPLES	# STEPS
MATH	65,800	3.8
CODE&DEBUGGING	56,200	4.4

Table 1: Statistics by domain in CRISP after applying filtering. We report the total number of generated instances, and the average number of steps per plan.

improves quality enough to warrant the deletion. Specifically, it improves performance on external benchmarks: The intrinsic filtering stage increases accuracy by 0.72 points on GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) and 0.44 points on MBPP (Austin et al., 2021) and HumanEval (Chen et al., 2021). Similarly, the extrinsic filtering, which validates the impact on the original task, further improves accuracy by 0.28 points for mathematics and 0.32 points for coding. These improvements confirm that our filtering pipeline successfully distills the datasets into high-quality high-level plans that are concise, clear, coherent, and complete, as well as beneficial on downstream tasks.

4 Experiments

4.1 Plan Generator Training

To demonstrate the practical benefits of CRISP, we applied LoRA parameter-efficient fine-tuning (Hu et al., 2022) on a small model to show that effective high-level plan generation is a learned capability; consequently, even a relatively small model, when efficiently fine-tuned on CRISP, can outperform a vanilla larger model on this task. Specifically, we used Granite-3.1-70B-Instruct model (Mishra et al., 2024). We trained the model for five epochs with a learning rate of $1e-5$ on four A100-80GB GPUs, using hyperparameters optimized through an extensive sweep. Additional technical details of the training procedure are provided in Appendix A.1.

4.2 Experimental Setup

To systematically assess the impact of high-level plan generation on downstream tasks, we examine various scenarios using both a small model (Granite-3.1-8B-Instruct Mishra et al., 2024) and a large model (Llama-3.1-70B-Instruct). We also experimented with another

small model, Llama-3.1-8B-Instruct, and found comparable results to Granite-8B as described in Appendix A.2. We conducted evaluations on four well-established benchmarks: the HumanEval benchmark (Chen et al., 2021), which assesses code synthesis and problem-solving capabilities, MBPP (Austin et al., 2021) which consists of around 1,000 crowd-sourced Python programming problems, GSM8K (Cobbe et al., 2021)—a grade school math word problems created by human problem writers, and the MATH benchmark (Hendrycks et al., 2021), which measures performance on complex mathematical problem-solving tasks. First, we established a baseline where the plan is based on classic Chain-of-Thought (CoT) prompting that relies solely on the problem description, which could be considered as having an emergent plan. Next, we incorporate high-level plans with plan-and-solve approach (Wang et al., 2023a)—these plans are generated by both the small and large models via few-shot prompting without any additional training, as well as by a fine-tuned version of the small model on CRISP detailed in Section 4.1. Once the plans are generated, another model takes them as input, along with the task description, and attempts to solve the task. We refer to the plan generation model as the ‘planner’ and to the subsequent model as the ‘solver’. We evaluated the generated plans using both the small and large models as solvers in a zero-shot setting without additional fine-tuning. This design enables us to directly compare the contribution of high-level plans produced by different models on various model sizes and domains.

4.3 Extrinsic Evaluation

Table 2 compares the results on different planners and solvers as well as the class Chain-of-Thought (CoT), and provides valuable insights into the impact of high-level plan generation across both coding and mathematical problem-solving domains.

Fine-tuned planning model achieves best results. The plans generated by the fine-tuned small model vastly outperforms the plans generated by the vanilla models and CoT across solvers and datasets. For example, improvements reach up to 28% error reduction against CoT in GSM8K. This demonstrates that fine-

Planner	Solver	MBPP Pass@1(Err ↓)	HumanEval Pass@1(Err ↓)	GSM8K Acc.(Err ↓)	MATH Acc.(Err ↓)
CoT (No Plan)	SMALL	60.1	71.2	84.6	49.3
VANILLA SMALL	SMALL	60.6 (1.3%↓)	71.9 (2.4%↓)	84.9 (1.9%↓)	49.6 (0.6%↓)
VANILLA LARGE	SMALL	62.0 (4.98%↓)	73.1 (6.6%↓)	85.7 (7.1%↓)	53.2 (7.7%↓)
FINE-TUNED SMALL	SMALL	64.8 (11.8%↓)	76.4 (18.1%↓)	87.1 (16.2%↓)	60.6 (22.3%↓)
CoT (No Plan)	LARGE	73.7	80.8	94.3	67.2
VANILLA SMALL	LARGE	73.4 (1.1%↑)	80.3 (2.6%↑)	93.6 (1.2%↑)	65.8 (4.3%↑)
VANILLA LARGE	LARGE	74.1 (1.5%↓)	82.2 (7.3%↓)	94.8 (8.8%↓)	70.2 (9.1%↓)
FINE-TUNED SMALL	LARGE	76.2 (9.5%↓)	85.3 (23.4%↓)	95.9 (28.1%↓)	73.1 (18.0%↓)

Table 2: Evaluation on code generation and math benchmarks across plan generator models and solution generator models. ‘Err↓’ represents the relative reduction in error compared to the baseline of Chain-of-Thought prompting (‘CoT (No Plan)’) with the same solver. ‘Small’ refers to Granite-3.1-8B-Instruct model and ‘Large’ refers to Llama-3.1-70B-Instruct model. Notably, the fine-tuned Granite gains the largest improvement in results across the four domains and solver models.

tuning a model for plan generation significantly benefits various myopic downstream tasks, such as code generation and mathematical problem-solving. It also shows that the plan generation capabilities of vanilla models, including larger ones like Llama-3.1-70B-Instruct, can be substantially improved, resulting in enhanced reasoning abilities.

Planning is better than CoT, yet the quality of plans matters. Plans generated by the vanilla models mostly outperformed CoT, with improvement of up to 9.1%. Yet, while the plans generated by the large model achieved significant improvements, the plan generated by the small model achieved only minor improvement and even degradation of up to 4.3% when using the large model as solver. This shows that explicitly generating the high-level plans before the solution is a better approach than CoT in myopic tasks, although the quality of the plan plays a critical role.

The fine-tuned planner equally improves both solvers. Incorporating plans generated by the small fine-tuned model into the prompts of both small and large solvers results in an average error reduction of 17.1% and 19.65%, respectively. This indicates that both models experience similar and significant improvements from receiving a high-quality plan before generating a solution.

Both the planner and the solver impact performance . Improving either the plan-

ner or the solver leads to performance gains. However, the solver has a greater influence on overall performance. This is evident when comparing the results of a fine-tuned small planner paired with a small solver to those of a vanilla small planner paired with a large solver, where the latter configuration yields significantly better results.

4.4 Intrinsic Evaluation

Following the extrinsic evaluation in Section 4.3, we conducted a direct comparison of plan quality using LLM-based judgment. Specifically, we evaluated the coherence, clarity, conciseness, and completeness of plans generated by fine-tuned Granite-3.1-8B-Instruct (described in Section 4.1) and Llama-3.1-8B-Instruct. We chose these two models as we would to further explore how much the fine-tuning helped compared to the best baseline across datasets.

The results are depicted in Figure 3. Surprisingly, although we used Llama-3.1-70B-Instruct as both a judge and a competitor, which should create a bias toward its own generations (Bitton et al., 2023; Koo et al., 2023), it preferred the plans generated by the small fine-tuned model across all datasets in 73.3% of the cases on average. Interestingly, fine-tuned Granite achieved the highest scores in the two harder datasets—MATH and HumanEval. This may indicate that the fine-tuning especially helped with plans that require more complex reasoning. The impact of our fine-tuned model’s high-level plans could be speculated to stem

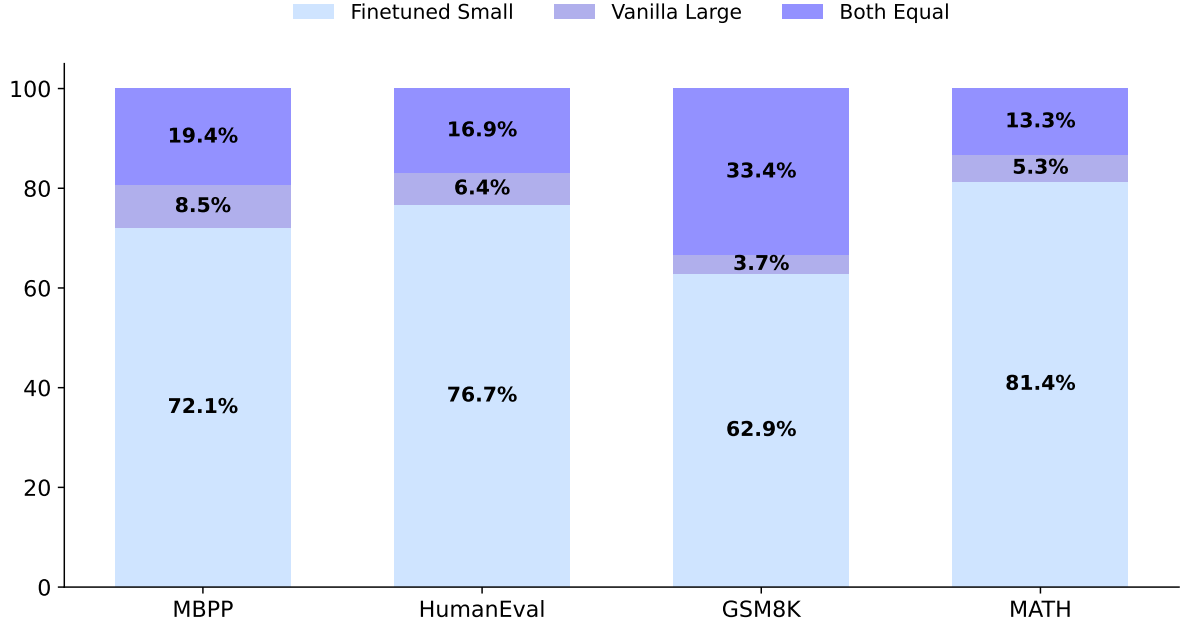


Figure 3: LLM-based judgement comparison for clarity, coherence, conciseness, and completeness between Granite-3.1-8B-Instruct fine-tuned on CRISP and Llama-3.1-70B-Instruct, which was also the judge. Each bar is divided into three sections representing the percentage of cases where the judge preferred the plan generated by one of the models or found both plans to be equally good.

from longer and more robust steps compared to the vanilla model’s high-level plans. However, when analyzing the average number of steps in the high-level plans across the four benchmarks, we observe that the large vanilla model generates, on average, 1.3 more steps than our small fine-tuned model on coding benchmarks and 0.8 more steps on math benchmarks. This disproves such speculation and suggests that fewer, more concise, and well-structured steps have a greater impact on the final solution.

4.5 Out-of-Domain Evaluation

We hypothesize that training a model to generate high-level plans in a specific domain—such as mathematics or coding—can provide it with transferable task decomposition capabilities that enhance performance in other domains. To investigate this, We evaluated our high-level plan generation models on out-of-distribution data by fine-tuning each model on one domain and testing it on the other, i.e. the mathematics-trained model to coding tasks and vice versa.

As shown in Table 3, training on out-of-domain data still substantially improves. Moreover, training improves on out-of-domain tasks almost as much as it does on in-domain,

and outperforms untrained large model. The mathematics-trained model, when applied to coding problems in the MBPP and HumanEval datasets, generated high-level plans that improved the final solution accuracy, achieving scores of 87.4 and 84.6, respectively. These results outperforms the vanilla models and are only marginally lower than those obtained by the specialized coding high-level plan generation model, which scored 87.9 and 85, respectively. Similarly, when the coding-trained model was tested on mathematical problems using the GSM8K and MATH benchmarks, it scored 95.2 and 71.9, compared to 96.9 and 73.1 achieved by the mathematics-trained model. These findings demonstrate that a model trained on one domain can indeed contribute effectively to problem-solving in another. Comparing those results with the ones in Table 2 shows that training on both the in-domain and the out-of-domain data further helps overall results. We take this to mean that diversity and or more high-level plan data are still valuable.

While not completely comparable, it seems that the mathematics-trained model, i.e. ‘Trained Small Math’ in table 3, demonstrates

PLANNER	MBPP Pass@1(Err ↓)	HUMANEVAL Pass@1(Err ↓)	GSM8K Acc.(Err ↓)	MATH Acc.(Err ↓)
CoT	73.7	80.8	94.3	67.2
VANILLA SMALL	73.4 (↑1.1%)	80.3 (↑1.5%)	93.6 (↑9.1%)	65.8 (↑6.5%)
VANILLA LARGE	74.1 (↓3.5%)	82.2 (↓7.3%)	94.8 (↓8.8%)	70.2 (↓9.1%)
TRAINED SMALL CODE	<u>76.1</u> (↓15.4%)	<u>85.0</u> (↓21.9%)	95.2 (↓15.8%)	71.9 (↓14.3%)
TRAINED SMALL MATH	75.4 (↓11.9%)	84.6 (↓19.8%)	<u>95.9</u> (↓28.1%)	<u>73.1</u> (↓18.0%)

Table 3: Out-of-domain evaluation with Llama-3.1-70B-Instruct as the solution generation model. ‘Trained Small Code/Math’ refers to Granite-3.1-8B-Instruct which was fine-tuned on one domain in CRISP. ‘Err↓’ refers to the error reduction in percentage compared to the CoT baseline.

stronger transfer performance on coding tasks compared to the coding-trained model’s performance on mathematical tasks. This is evident from its error reduction, which is relatively close to that of the coding-trained model. Specifically, the average difference in error reduction between the two was 2.8% in code generation benchmarks, compared to 8% in math benchmarks.

This asymmetry can be attributed to the intrinsic relationship between mathematical reasoning and coding. Algorithmic programming often relies on mathematical concepts such as logic, recursion, combinatorics, probability, and number theory. Consequently, a model trained on mathematical problems is likely to develop robust reasoning, pattern recognition, and generalization skills—attributes that are critically important for effective coding. In contrast, while a model trained on coding problems may acquire knowledge of syntax and common programming patterns, it might not cultivate the deeper mathematical reasoning skills that are essential for addressing abstract or algorithmically complex tasks.

In summary, our results suggest that high-level plan generation models possess a notable degree of domain generalizability; they improve scores substantially, and the data provided and its part all contribute to performance, surpassing generation from much stronger models. The abstract reasoning and general problem-solving strategies fostered by training on mathematical problems appear to be more readily transferable to coding tasks than the reverse. This observation underscores the potential benefits of leveraging cross-domain training to enhance the versatility and effectiveness of problem-

solving models. We believe this transferability also extends to other domains and topics, even those that are not inherently symbolic.

5 Conclusions

In this work, we introduced CRISP, a dataset for enhancing complex reasoning in large language models through structured high-level planning. CRISP was developed through a rigorous data generation process, leveraging existing problem-solving datasets to extract structured high-level plans, followed by an extensive filtering and validation pipeline. Our experiments demonstrated that fine-tuning on CRISP enables smaller models to generate higher-quality plans, outperforming much larger models across mathematical reasoning and code generation tasks. Additionally, our intrinsic evaluation revealed that plans generated by fine-tuned models were shorter, more concise, coherent, and complete compared to those from vanilla models. We also showed that high-level planning capabilities transfer effectively across domains, with fine-tuning in one domain improving performance in another. This highlights the generalizability of structured planning as a trainable capability that enhances reasoning efficiency across domains. By releasing CRISP, we aim to encourage further research into explicit planning mechanisms, structured reasoning, and their broader applications in NLP. Future work may explore expanding CRISP to additional domains and refining planning strategies to bridge the gap between human and machine reasoning further.

6 Limitations

While high-level planning was shown to benefit highly from training data, and while we do release a substantial amount of data for two domains, it is likely that other domains would benefit from such datasets and would require further work to apply our methods (or new ones) to them. Moreover, as we base our data on data that existed in other forms and for other purposes, this may not be available in other domains.

References

- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. [Large language models for mathematical reasoning: Progresses and challenges](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 225–237, St. Julian’s, Malta. Association for Computational Linguistics.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. [Graph of thoughts: Solving elaborate problems with large language models](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690.
- Yonatan Bitton, Hritik Bansal, Jack Hessel, Rulin Shao, Wanrong Zhu, Anas Awadalla, Josh Gardner, Rohan Taori, and Ludwig Schmidt. 2023. Visit-bench: A benchmark for vision-language instruction following inspired by real-world use. *arXiv preprint arXiv:2308.06595*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *Transactions on Machine Learning Research*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. [Mind2web: Towards a generalist agent for the web](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yangruibo Ding, Jinjun Peng, Marcus J. Min, Gail Kaiser, Junfeng Yang, and Baishakhi Ray. 2024. [Semcoder: Training code language models with comprehensive semantics reasoning](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2022. Successive prompting for decomposing complex questions. *arXiv preprint arXiv:2212.04092*.
- William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela Veloso, and Ruslan Salakhutdinov. 2019. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv preprint arXiv:1907.13440*.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023. [Reasoning with language model is planning with world model](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8154–8173, Singapore. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the MATH dataset](#). In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.
- Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. 2024a. [A survey on large language models for code generation](#). *Preprint*, arXiv:2406.00515.
- Xue Jiang, Yihong Dong, Lecheng Wang, Zheng Fang, Qiwei Shang, Ge Li, Zhi Jin, and Wenpin Jiao. 2024b. Self-planning code generation with large language models. *ACM Transactions on Software Engineering and Methodology*, 33(7):1–30.

715	Shyam Sundar Kannan, Vishnunandan L. N.	Archiki Prasad, Alexander Koller, Mareike Hart-	773
716	Venkatesh, and Byung-Cheol Min. 2024. Smart-	mann, Peter Clark, Ashish Sabharwal, Mohit	774
717	llm: Smart multi-agent robot task planning using	Bansal, and Tushar Khot. 2024. ADaPT: As-	775
718	large language models . In <i>2024 IEEE/RSJ In-</i>	needed decomposition and planning with lan-	776
719	<i>ternational Conference on Intelligent Robots and</i>	guage models. In <i>Findings of the Association for</i>	777
720	<i>Systems (IROS)</i> , pages 12140–12147.	<i>Computational Linguistics: NAACL 2024</i> , pages	778
721	Tushar Khot, Harsh Trivedi, Matthew Finlayson,	4226–4252, Mexico City, Mexico. Association for	779
722	Yao Fu, Kyle Richardson, Peter Clark, and	Computational Linguistics.	780
723	Ashish Sabharwal. 2023. Decomposed prompt-	Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li,	781
724	ing: A modular approach for solving complex	Tingwu Wang, Sanja Fidler, and Antonio Tor-	782
725	tasks . In <i>The Eleventh International Conference</i>	ralba. 2018. Virtualhome: Simulating household	783
726	<i>on Learning Representations</i> .	activities via programs. In <i>Proceedings of the</i>	784
727	Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn	<i>IEEE conference on computer vision and pattern</i>	785
728	Park, Zae Myung Kim, and Dongyeop Kang.	<i>recognition</i> , pages 8494–8502.	786
729	2023. Benchmarking cognitive biases in large	Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan	787
730	language models as evaluators. <i>arXiv preprint</i>	Hernandez, and Percy Liang. 2017. World of	788
731	<i>arXiv:2309.17012</i> .	bits: An open-domain platform for web-based	789
732	Boyi Li, Philipp Wu, Pieter Abbeel, and Jitendra	agents . In <i>Proceedings of the 34th International</i>	790
733	Malik. 2023. Interactive task planning with lan-	<i>Conference on Machine Learning</i> , volume 70 of	791
734	guage models . In <i>2nd Workshop on Language</i>	<i>Proceedings of Machine Learning Research</i> , pages	792
735	<i>and Robot Learning: Language as Grounding</i> .	3135–3144. PMLR.	793
736	Evan Zheran Liu, Kelvin Guu, Panupong Pasupat,	Noah Shinn, Federico Cassano, Ashwin Gopinath,	794
737	and Percy Liang. 2018. Reinforcement learning	Karthik R Narasimhan, and Shunyu Yao. 2023.	795
738	on web interfaces using workflow-guided explo-	Reflexion: language agents with verbal reinforce-	796
739	ration . In <i>International Conference on Learning</i>	ment learning . In <i>Thirty-seventh Conference on</i>	797
740	<i>Representations</i> .	<i>Neural Information Processing Systems</i> .	798
741	Aman Madaan, Niket Tandon, Prakhar Gupta,	Mohit Shridhar, Jesse Thomason, Daniel Gordon,	799
742	Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri	Yonatan Bisk, Winson Han, Roozbeh Mottaghi,	800
743	Alon, Nouha Dziri, Shrimai Prabhumoye, Yim-	Luke Zettlemoyer, and Dieter Fox. 2020. Alfred:	801
744	ing Yang, Shashank Gupta, Bodhisattwa Prasad	A benchmark for interpreting grounded instruc-	802
745	Majumder, Katherine Hermann, Sean Welleck,	tions for everyday tasks. In <i>Proceedings of the</i>	803
746	Amir Yazdanbakhsh, and Peter Clark. 2023. Self-	<i>IEEE/CVF conference on computer vision and</i>	804
747	refine: Iterative refinement with self-feedback . In	<i>pattern recognition</i> , pages 10740–10749.	805
748	<i>Thirty-seventh Conference on Neural Information</i>	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre	806
749	<i>Processing Systems</i> .	Cote, Yonatan Bisk, Adam Trischler, and	807
750	Mayank Mishra, Matt Stallone, Gaoyuan Zhang,	Matthew Hausknecht. 2021. {ALFW}orld:	808
751	Yikang Shen, Aditya Prasad, Adriana Meza	Aligning text and embodied environments for	809
752	Soria, Michele Merler, Parameswaran Selvam,	interactive learning . In <i>International Conference</i>	810
753	Saptha Surendran, Shivdeep Singh, Manish Sethi,	<i>on Learning Representations</i> .	811
754	Xuan-Hong Dang, Pengyuan Li, Kun-Lung Wu,	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu,	812
755	Syed Zawad, Andrew Coleman, Matthew White,	Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim.	813
756	Mark Lewis, Raju Pavuluri, Yan Koyfman, Boris	2023a. Plan-and-solve prompting: Improving	814
757	Lublinksky, Maximilien de Bayser, Ibrahim Ab-	zero-shot chain-of-thought reasoning by large	815
758	delaziz, Kinjal Basu, Mayank Agarwal, Yi Zhou,	language models . In <i>Proceedings of the 61st</i>	816
759	Chris Johnson, Aanchal Goyal, Hima Patel,	<i>Annual Meeting of the Association for Compu-</i>	817
760	S. Yousaf Shah, Petros Zerfos, Heiko Ludwig,	<i>tational Linguistics (Volume 1: Long Papers)</i> ,	818
761	Asim Munawar, Maxwell Crouse, Pavan Kapani-	pages 2609–2634, Toronto, Canada. Association	819
762	pathi, Shweta Salaria, Bob Calio, Sophia Wen,	for Computational Linguistics.	820
763	Seetharami Seelam, Brian Belgodere, Carlos A.	Xuezhi Wang, Jason Wei, Dale Schuurmans,	821
764	Fonseca, Amith Singhee, Nirmal Desai, David D.	Quoc V Le, Ed H. Chi, Sharan Narang,	822
765	Cox, Ruchir Puri, and Rameswar Panda. 2024.	Aakanksha Chowdhery, and Denny Zhou. 2023b.	823
766	Granite code models: A family of open foun-	Self-consistency improves chain of thought rea-	824
767	dation models for code intelligence . <i>CoRR</i> ,	soning in language models . In <i>The Eleventh</i>	825
768	abs/2405.04324.	<i>International Conference on Learning Represen-</i>	826
769	Aske Plaat, Annie Wong, Suzan Verberne, Joost	<i>tations</i> .	827
770	Broekens, Niki van Stein, and Thomas Back.	Jason Wei, Xuezhi Wang, Dale Schuurmans,	828
771	2024. Reasoning with large language models, a	Maarten Bosma, brian ichter, Fei Xia, Ed H.	829
772	survey . <i>Preprint</i> , arXiv:2407.11511.		

Chi, Quoc V Le, and Denny Zhou. 2022a. [Chain of thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022b. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.

Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, Min-Yen Kan, Junxian He, and Qizhe Xie. 2023. [Self-evaluation guided beam search for reasoning](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2025. [Magpie: Alignment data synthesis from scratch by prompting aligned LLMs with nothing](#). In *The Thirteenth International Conference on Learning Representations*.

John Yang, Akshara Prabhakar, Karthik R Narasimhan, and Shunyu Yao. 2023. [Intercode: Standardizing and benchmarking interactive coding with execution feedback](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. [Webshop: Towards scalable real-world web interaction with grounded language agents](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 20744–20757. Curran Associates, Inc.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Jesse Zhang, Jiahui Zhang, Karl Pertsch, Ziyi Liu, Xiang Ren, Minsuk Chang, Shao-Hua Sun, and Joseph J Lim. 2023. [Bootstrap your own skills: Learning to solve new tasks with large language model guidance](#). In *7th Annual Conference on Robot Learning*.

Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2024. [Language agent tree search unifies reasoning acting and planning in language models](#).

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. [Least-to-most prompting enables complex reasoning in large language models](#). In *The Eleventh International Conference on Learning Representations*.

A Appendix

A.1 LoRA Finetuning

We trained Granite-3.1-8B-Instruct¹ on CRISP with LoRA for 5 epochs, with $R = 32$, $\alpha = 16$, dropout ratio of 0.05%, a learning rate of 1e-5, a cosine learning rate scheduler and a 0.0 weight decay.

A.2 Results with Llama-3.1-8B-instruct

Here are the results for another small model that we experimented with: Llama-3.1-8B-Instruct. We did that to make sure that the results obtained with Granite-3.1-8B-Instruct are indeed representative.

A.3 Prompt for Plan Generation in CRISP

To generate plans for each domain in CRISP, we crafted a few-shot prompt for each domain. Here is the prompt we used for the generation of plans in the Math domain of Magpie-reasoning-V1-150K. The overall objective was to extract the logical strategy needed to solve a problem without relying on specific equations, function names, or detailed computations.

System Message:

You are a helpful and concise assistant. You have access to:

1. A **Problem Description** that explains the problem at hand.
2. A **Detailed Solution** that fully works out how to solve the problem step-by-step.

Your goal is to produce a short, high-level plan describing how to solve the problem logically.

This plan must not include any specific equations, function names, or detailed numerical computations.

It should be purely indicative and helpful, outlining the logical strategy in 3-5 simple steps.

User Message:

Task

1. Read and understand the **Problem Description** below.
2. Review the **Detailed Solution** below (do not copy it).
3. From these, generate a concise, 3-5 step high-level plan that explains the logical approach needed to solve the problem.

¹<https://huggingface.co/ibm-granite/granite-3.1-8b-instruct>

PLANNER MODEL	MBPP Pass@1(Err ↓)	HUMAN EVAL Pass@1(Err ↓)	GSM8K Acc.(Err ↓)	MATH Acc.(Err ↓)
CoT	60.6	72.4	84.3	51.5
VANILLA LLAMA3.1-8B	60.9 (0.8%↓)	72.6 (0.7%↓)	84.7 (2.5%↓)	51.9 (0.8%↓)
VANILLA LLAMA-70B	62.3 (4.3%↓)	73.5 (4.0%↓)	85.6 (8.3%↓)	52.7 (2.5%↓)
TRAINED LLAMA	64.9 (10.9%↓)	76.1 (13.4%↓)	86.3 (12.7%↓)	54.5 (6.2%↓)

Table 4: Results with Llama-3.1-8B-Instruct as a generator model with different planners, including a fine-tuned model of the aforementioned model. Notably, the trends are similar to the trends seen with Granite-3.1-8B-Instruct.

4. The plan should be abstract and conceptual-avoid quoting or revealing detailed equations, formulas, function names, or code.	**Detailed Solution** Let's break it down step by step. Step 1: Triangles with 3 vertices from the octagon (choose any 3 of 8). Step 2: Triangles with 2 vertices from the octagon plus the center. Then sum the totals from Step 1 and Step 2.
5. Focus on the reasoning steps rather than low-level implementation.	
Problem Description {problem_description\}	
Detailed Solution {detailed_solution\}	**High-Level Plan** 1. Recognize the two types of triangles possible: those with only octagon vertices and those that use the center as one vertex. 2. Conceptually determine how to count each type of triangle without going into specific combinations. 3. Combine the counts logically to get the total number of different triangles.
Formatting Requirements 1. Your final answer should be 3-5 bullet points (or numbered steps). 2. Each bullet/step should be brief, logical, and to the point. 3. Do not include specific equations or code references. 4. Do not include extraneous commentary or repeat large sections from the solution. 5. Focus on a clear, conceptual strategy that someone could follow to solve the problem at a high level.	**Example 2** **Problem Description** Write a function that merges two sorted lists into a single sorted list, without using any built-in sorting functions. The time complexity should be $O(n)$, where n is the total number of elements in both lists.
Example Output Structure 1. Identify the main elements, quantities, or variables in the problem. 2. Determine the key relationships or principles that connect these elements. 3. Outline a general strategy for combining or manipulating these elements to get closer to a solution. 4. Check or validate the approach by ensuring it aligns with the key requirements. 5. Summarize the final reasoning step or expected result in broad terms.	**Detailed Solution** def merge_sorted_lists(list1, list2): i, j = 0, 0 result = [] while i < len(list1) and j < len(list2): if list1[i] < list2[j]: result.append(list1[i]) i += 1 else: result.append(list2[j]) j += 1 while i < len(list1): result.append(list1[i]) i += 1 while j < len(list2): result.append(list2[j]) j += 1 return result
Few-Shot Examples	
Example 1	
Problem Description Consider a regular octagon. How many different triangles can be formed if the octagon is placed inside a circle and we can also use the center of the circle as a vertex for the triangles?	**High-Level Plan**

1047	1. Recognize the need to keep track of where	- Are all necessary steps included to fully	1112
1048	we are in each list as we form the new	address the problem and derive the	1113
1049	list.	solution?	1114
1050	2. Conceptually compare the front elements	- Does the plan leave out any critical	1115
1051	from both lists to decide which goes next.	information or assume unstated knowledge?	1116
1052	3. Continue until one list is exhausted, then	- Are there any logical leaps where a step is	1117
1053	add any remaining elements from the other.	missing between two points?	1118
1054	4. Return the combined list as the final	---	1119
1055	merged sequence.		1120
1056		### **Output Format**	1121
1057		Your evaluation must be returned **as a	1122
1058	Now, please provide your high-level plan in	single JSON object** containing exactly	1123
1059	3-5 steps.	**two fields**:	1124
			1125
1060	A.4 LLM-based Judgement Prompt	1. **`explanation`**: A detailed assessment,	1126
		addressing how well the plan meets each	1127
1061	We used the following prompt for judging the	of the four criteria above. Reference	1128
1062	four attributes of generated plans with LLM-	specific strengths and weaknesses.	1129
1063	based judgement.	2. **`judgement`**: A string set to `true`	1130
		if the plan **fully satisfies all four	1131
		criteria**, or `false` otherwise.	1132
		---	1133
1064	You are an intelligent, knowledgeable, and		1134
1065	impartial judge. Your task is to evaluate	### **Strict Output Requirements:**	1135
1066	whether a **High-Level Plan** effectively	- **Do not** include any extra keys or fields.	1136
1067	outlines the logical steps required to	- **Do not** output any additional text	1137
1068	address a given **Problem Description**	outside the JSON structure.	1138
1069	and reach a **Detailed Solution**.	- The final output must strictly match the	1139
1070		following format:	1140
1071	You are provided with three components:		1141
1072	1. **Problem Description:**	```json	1142
1073	{problem_description}	{	1143
1074	2. **High-Level Plan:** {high_level_plan}	"explanation": "Your detailed reasoning	1144
1075	3. **Detailed Solution:** {solution}	here.",	1145
1076		"judgement": "true or false"	1146
1077	---	}	1147
1078			1148
1079	### **Evaluation Criteria**		
1080	Assess whether the High-Level Plan	A.5 Prompt for Intrinsic Evaluation	1149
1081	sufficiently and logically bridges the	We attach here the prompt we used to com-	1150
1082	Problem Description and the Detailed	pare two plans based on clarity, conciseness,	1151
1083	Solution based on the following four	coherence, and completeness.	1152
1084	criteria:		
1085		You are an impartial and expert judge. Your	1153
1086	#### **1. Clarity**	task is to evaluate two plans that each	1154
1087	- Are the steps written in a way that is easy	aim to solve the same problem.	1155
1088	to understand?	They both rely on the same problem	1156
1089	- Does the plan avoid ambiguity and vague	description and reach the same final	1157
1090	language?	solution, but they may differ in how they	1158
1091	- Are complex ideas broken down into	outline the logical steps to get from the	1159
1092	comprehensible components?	problem statement to the solution.	1160
1093			1161
1094	#### **2. Conciseness**	### Your Goal	1162
1095	- Does the plan avoid unnecessary repetition	Read the problem description, the detailed	1163
1096	or overly verbose explanations?	solution, and both Plan A and Plan B	1164
1097	- Are the steps described efficiently without	carefully.	1165
1098	losing essential details?	Then, compare and evaluate Plan A and Plan B	1166
1099	- Is there any redundant or overly wordy	according to four specific criteria:	1167
1100	content that could be simplified?		1168
1101		1. **Clarity**	1169
1102	#### **3. Coherence (Logical Flow &	- Are the steps written in a way that is	1170
1103	Structure)**	easy to understand?	1171
1104	- Do the steps follow a clear and logical	- Does the plan avoid ambiguity and vague	1172
1105	progression from problem to solution?	language?	1173
1106	- Are there any gaps, abrupt transitions, or	- Are complex ideas broken down into	1174
1107	missing links in the reasoning?	comprehensible components?	1175
1108	- Is the structure intuitive, making it easy		1176
1109	to follow the problem-solving approach?		1177
1110			
1111	#### **4. Completeness**	2. **Conciseness**	

1178	- Does the plan avoid unnecessary	indicating why you believe	1248
1179	repetition or overly verbose	Plan A or Plan B is better, or why they are	1249
1180	explanations?	the same. Please point to relevant	1250
1181	- Are the steps described efficiently,	details from each plan	1251
1182	without omitting crucial details?	when forming your reasoning.	1252
1183	- Is there any redundant or overly wordy		1253
1184	content that could be simplified?	- **judgement** : Must be exactly one of:	1254
1185		- "A" (if Plan A is judged superior	1255
1186	3. **Coherence (Logical Flow & Structure)**	overall),	1256
1187	- Do the steps follow a clear and logical	- "B" (if Plan B is judged superior	1257
1188	progression from problem to solution?	overall),	1258
1189	- Are there any gaps, abrupt transitions,	- "Same" (if they are equally good).	1259
1190	or missing links in the reasoning?		1260
1191	- Is the structure intuitive and easy to	Ensure you base your judgment only on the	1261
1192	follow?	given criteria and the content of the	1262
1193		plans. Output **only** the JSON with	1263
1194	4. **Completeness**	no additional text, headers, or explanations.	1264
1195	- Are all the essential steps included to		
1196	fully address the problem and derive		
1197	the solution?		
1198	- Does the plan omit any critical		
1199	information or assume unstated		
1200	knowledge?		
1201	- Are there any logical leaps or missing		
1202	transitions between key points?		
1203			
1204	### Inputs		
1205	**Problem Description:**		
1206	{problem_description}		
1207			
1208	**Detailed Solution:**		
1209	{solution}		
1210			
1211	**Plan A:**		
1212	{planA}		
1213			
1214	**Plan B:**		
1215	{planB}		
1216			
1217	Each plan proposes a logical sequence of		
1218	steps to move from the problem		
1219	description to the final solution.		
1220			
1221	### What to Do		
1222	1. Examine each plan in relation to the		
1223	problem and the solution.		
1224	2. Assess Plan A and Plan B based on the four		
1225	criteria: Clarity, Conciseness,		
1226	Coherence, and Completeness.		
1227	3. Decide whether Plan A is superior, Plan B		
1228	is superior, or they are equally good		
1229	overall.		
1230			
1231	### How to Report		
1232	Provide your final output as a single JSON		
1233	object in the exact format below:		
1234			
1235	{		
1236	"explanation": "Explain your comparison		
1237	referencing each of the four criteria		
1238	as needed.		
1239	Describe strengths,		
1240	weaknesses, and the		
1241	reasoning leading to your		
1242	final verdict.",		
1243	"judgement": "A or B or Same"		
1244	}		
1245			
1246	- **explanation** : Briefly but		
1247	comprehensively summarize the comparison,		