

# ARE LLMs BETTER FORMALIZERS THAN SOLVERS ON COMPLEX PROBLEMS?

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

A trending line of recent work advocates for using large language models (LLMs) as formalizers instead of as end-to-end solvers for logical reasoning problems. Instead of generating the solution, the LLM generates a formal program that derives a solution via an external solver. While performance gain of the seemingly scalable LLM-as-formalizer over the seemingly unscalable LLM-as-solver has been widely reported, we show that this superiority does not hold on real-life constraint satisfaction problems. On 4 domains, we systematically evaluate 6 LLMs including 4 large reasoning models with inference-time scaling, paired with 5 pipelines including 2 types of formalism. We show that in few-shot settings, LLM-as-formalizer underperforms LLM-as-solver. While LLM-as-formalizer promises accuracy, robustness, faithfulness, and efficiency, we observe that the present LLMs do not yet deliver any of those, as their limited ability to generate formal programs leads to failure to scale with complexity, hard-coded solutions, and excessive reasoning tokens. We present our detailed analysis and actionable remedies to drive future research that improves LLM-as-formalizer.<sup>1</sup>

## 1 INTRODUCTION

Considerable efforts have been dedicated to improving the ability of large language models (LLMs) to solve logical reasoning tasks (Saxton et al., 2019; Clark et al., 2021), primarily in an end-to-end manner involving intermediate chain-of-thought tokens (Wei et al., 2022; Kojima et al., 2022). Despite their strong performance, these *LLM-as-solver* methods have been shown to lack faithfulness, verifiability, and formal guarantee (Lyu et al., 2023). To bridge this gap, an emerging line of neuro-symbolic methods use LLMs to translate the natural language problem description into a formal program, from which a solution can be derived using an external solver (Gao et al., 2023; Pan et al., 2023; Han et al., 2024). In addition to promising the desirable features above, these *LLM-as-formalizer* methods have been reported to achieve state-of-the-art in some tasks with high complexity where LLM-as-solver fails to scale (Valmeekam et al., 2024), while failing in others (Chen et al., 2025; Kagitha et al., 2025) due to the inherent challenge of code generation.

Recently, an enhanced chain-of-thought reasoning paradigm, inference-time scaling (Muennighoff et al., 2025; Guo et al., 2025) explicitly trains LLMs to generate more thinking tokens that scales with the problem complexity. Some have reported that these reasoning LLMs (LRMs) as planners start to outperform those as formalizers (Huang & Zhang, 2025), while others have reported that they do not scale beyond a certain complexity (Shojaee et al., 2025). To date, it remains unclear which methodology is preferred, hindering real-life application.

We attempt to provide clarity on the choice between *LLM-as-solver* and *LLM-as-formalizer* by performing a systematic evaluation and analysis. We focus on real-life constraint satisfaction problems (CSPs) including calendar scheduling, trip planning, and meeting planning, in addition to a logic grid puzzle. On these 4 domains, we evaluate 6 state-of-the-art LLMs including 4 LRMs (DeepSeek-R1, Qwen3-32B, o3-mini-high, GPT-5) and 2 non-reasoning counterparts (DeepSeek-V3, Qwen2.5-32B), both as solvers and as formalizers. To ensure maximal generalizability and fairness of comparison, we consider only few-shot settings where revision is allowed for *LLM-as-formalizer* only given solver errors. We consider two types of formalism as LLMs’ generation target, free-form Python and specific code to interface a Satisfiability Modulo Theories (SMT) solver.

<sup>1</sup>Our code and data are attached to the submission.

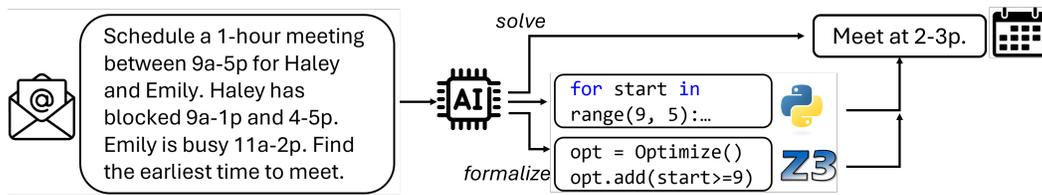


Figure 1: An example of the output for *LLM-as-solver* (top) and *LLM-as-formalizer* with two types of formalism (middle and bottom) in a real-life constraint satisfaction task, calendar scheduling.

While *LLM-as-formalizer* is posed to be accurate, robust, feasible, and efficient, we observe that the present LRMs and LLMs deliver none of the above, with the following key findings:

1. Both LRMs and regular LLMs are often worse formalizers than solvers;
2. Neither *LLM-as-solver* nor *LLM-as-formalizer* is robust to problem complexity;
3. *LLM-as-formalizer* suffers from mis-translating constraints even for the strongest LLMs;
4. *LLM-as-formalizer* does not generate fewer tokens than *LLM-as-solver* due to inefficient reasoning, or unnecessary solver-like reasoning that may lead to hard-coded solutions.

Based on these observations, we propose actionable recommendations for research in this area.

## 2 RELATED WORK

### 2.1 *LLM-as-solver*

*LLM-as-solver*, also known as informal reasoning or end-to-end reasoning, is a paradigm when one or more LLMs generate optional intermediate tokens and eventually the solution to a logical reasoning problem. As early problem solving ability comes from training, many major LLMs have been trained or evaluated as a solver (Rajani et al., 2019; Brown et al., 2020), while limited to low-complexity, common-sense problems. The introduction of chain-of-thought reasoning made *LLM-as-solver* feasible in complex tasks such as mathematics and logic puzzles (Wei et al., 2022; Kojima et al., 2022; Sel et al., 2024; Feng et al., 2023). Follow-up techniques to further improve *LLM-as-solver* based on chain-of-thought included self-verification (Weng et al., 2023), self-refine (Madaan et al., 2023), tree-of-thought (Yao et al., 2023), self-consistency (Wang et al., 2023b), etc.

Despite strong performance, *LLM-as-solver* has been shown to lack faithfulness, verifiability, and formal guarantees (Lyu et al., 2023; Turpin et al., 2023; Stechly et al., 2025). To bridge these gaps, neuro-symbolic methods were introduced to combine *LLM-as-solver* with formal tools (Jha et al., 2023; Saparov & He, 2023).

### 2.2 *LLM-as-formalizer*

*LLM-as-formalizer* is a subset of the above neuro-symbolic methods where LLMs generate neither the solution nor the chain-of-thought towards the solution, but rather translate the natural language problem description to an executable program based on LLM’s ability to generate code. This methodology has been reported to greatly outperform *LLM-as-solver* in classical planning (Liu et al., 2023; Xie et al., 2023; Hao et al., 2025), constraint satisfaction (Berman et al., 2024; Kesseli et al., 2025), mathematics (Wu et al., 2022; Yang et al., 2023; Jiang et al., 2023), and general logical reasoning tasks (Ye et al., 2023; Gao et al., 2023; Pan et al., 2023; Han et al., 2024). Despite the overwhelming report of its success, it has also been shown to be unstable, dependent on the choice of formal language and solver (Matthew Lam et al., 2024; Beiser et al., 2025).

### 2.3 *LLM-as-formalizer* vs. *LLM-as-solver*

For effective application, research like ours that systematically compares *LLM-as-formalizer* and *LLM-as-solver* is crucial but lacking, especially with the introduction of LRMs which reported dramatically improved performance on algorithmic tasks where *LLM-as-solver* used to fail (Muenighoff et al., 2025; Guo et al., 2025). The closest work to ours is Chen et al. (2025), which did

	Calendar Scheduling	Trip Planning	Meeting Planning	Zebra Logic
$X_i$	start hour $s$ , end hour $e$	day $d$ , city $c$	person $p$ , start $s$ , duration $d$	person $p$ , attribute $a$
$D_i$	$\{(s, e) \mid s, e \in [9, 17], e > s\}$	$\{(d, c) \mid d \in [1, \text{len}], c \in C\}$	$\{(p, s, d) \mid p \in P, s \in [9, 24], d \in [1, 12]\}$	$\{(p, a) \mid p \in P, a \in A\}$
$C_j$	Unavailable time range Meeting duration	Total number of days Allowed direction of travel Expected duration in a city Expected city on a day	Start time and place Travel time of two places Time and place of a person Min. duration with a person	Various provided clues about people and attributes

Table 1: Formal definition of variables  $X_i$ , domains  $D_i$ , and constraints  $C_j$  from the 4 domains in NaturalPlan and ZebraLogic.

not consider any open-source model or inference-time scaling LRMs whereas open-source LRMs are our focus. Moreover, their work is an initial investigation where conclusions are surface-level (formalizing is hard) and inconclusive (neither is optimal), while ours is a deep-dive resulting in an array of fine-grained, relatively definitive, and somewhat counterintuitive conclusions. The other work that considered both *LLM-as-formalizer* and *LLM-as-solver* for LRMs is Kagitha et al. (2025) which only focuses on classical planning tasks, similarly lacking decisive findings.

### 3 TASK AND DATA

To systematically compare *LLM-as-solver* and *LLM-as-formalizer* including various formalisms, we prioritize depth over breadth by focusing on real-world CSPs. Such a problem includes the following elements described in natural language:

1. A set of variables  $\{X_i\}$ , to be assigned values
2. Domains of variables  $\{D_i\}$ , the possible values of the variables
3. A set of constraints  $\{C_j\}$ , boolean rules that govern one or more variables

The goal of planning is thus to find a value assignment of each variable  $X_i$  based on possibilities  $D_i$  that does not violate any constraints  $C_j$ .

We consider three domains including calendar scheduling, trip planning, and meeting planning provided by the NaturalPlan dataset (Zheng et al., 2024), where only *LLM-as-solver* methods have been evaluated (Lee et al., 2025; Parmar et al., 2025) but not *LLM-as-formalizer*. We additionally consider a logic grid puzzle domain from the ZebraLogic dataset (Lin et al., 2025) which is less natural but considered in past work (Berman et al., 2024; Kesseli et al., 2025) evaluating *LLM-as-formalizer* methods.

The formulation and examples of all three tasks are shown in Table 1. Unlike previous work that evaluated the predicted solution by matching the ground-truth solution, we manually and formally annotate all constraints  $\{C_j\}$ . We therefore count a solution as correct if it formally satisfies all the constraints. This design choice enables us to perform fine-grained analysis on problem complexity and model errors. Bound by the cost of annotation, we randomly sample 100 out of 1,000 examples for each domain. The tasks and prompts are exemplified in Appendix A.

### 4 METHODS

We prompt LLMs in a one-shot manner to generate three different types of output and formalism.

**Solution.** This is the *LLM-as-solver* pipeline where the LLM is given an input and directly outputs the answer, optionally after generating a reasoning chain. Noting that most of our LLMs are LRMs, we do not explicitly prompt any model to “think step by step.”

**Python code.** This is an *LLM-as-formalizer* pipeline where the LLM is given an input and generates a Python program that will be executed to output the answer. Conceptually, the model generates both the declarative component (i.e., the variables and constraints) and the search component (i.e., algorithm to search for the correct plan).

**SMT code.** Alternatively, the LLM may generate a specific program to invoke a CSP solver, such as an SMT solver. Following previous work, we prompt LLMs to generate code using the Z3 solver

162 Python wrapper<sup>2</sup>. Here, the model primarily generates the declarative component as the search  
 163 component is simply a call to the pre-defined solver function. Examples of the two formalisms are  
 164 juxtaposed in Appendix C.

165 For the two *LLM-as-formalizer* pipelines, we follow Pan et al. (2023); Ye et al. (2023) to optionally  
 166 include a revision-by-error module. If the solver, either implemented by the model in the case of  
 167 Python or provided by the Z3 library in case of SMT, returns an error message or cannot find any  
 168 plan, this signal is returned to the LLM to re-generate the program for up to 5 times. We consider  
 169 this addition sufficient to represent state-of-the-art methods without loss of the generality, as more  
 170 sophisticated pipelines such as (Hao et al., 2025) are specific to a particular formalism or language,  
 171 defeating our purpose of studying general, real-life applications.

172 We consider 6 state-of-the-art LLMs or LLM products spanning 3 dimensions: size, whether it is  
 173 open-source, and whether it is trained to generate a scaling reasoning chain (an LRM). The open-  
 174 source, non-reasoning LLMs include **Qwen-2.5-Coder-32B** and **DeepSeek-V3-671B**. The open-  
 175 source LRMs include **Qwen-3-32B** and **DeepSeek-R1-671B**. The closed-source LRMs include **o3-**  
 176 **mini-high-2025-01-31** and **gpt-5-2025-08-07**. We run the OpenAI and DeepSeek models via their  
 177 APIs respectively. We run the Qwen models using KANI (Zhu et al., 2023) and HuggingFace<sup>3</sup> lo-  
 178 cally on 8 H100 GPUs with default hyperparameters. For easy extraction of the generated program,  
 179 we constrain decoding to JSON when possible via APIs and Outlines<sup>4</sup>.

180 All models are prompted in a one-shot manner for the most generalizable, cost-effective setting,  
 181 as we show in preliminary studies (Appendix B) that 5-shot prompting does not lead to system-  
 182 atic performance gains for LLM-as-planner. The one in-context exemplar includes the description  
 183 and solution in the expected format of a handpicked, representative problem in each domain. The  
 184 descriptive prompt includes an instruction to generate each target and some reminders about the  
 185 idiosyncrasy of each domain that a human solver would otherwise not consider. For example, in  
 186 trip planning, if the plan involves flying from one location to another, the dataset requires tying both  
 187 locations to the day of the flight. Without describing and exemplifying this rule explicitly, all models  
 188 in our preliminary studies consistently fail, leading to unfair underestimation of their performance.

## 190 5 RESULTS

191 We present our findings guided by a series of research questions.

### 194 5.1 RQ1: ARE LLMs BETTER FORMALIZERS THAN SOLVERS?

195 As discussed in Section 2.2, much existing work has shown the efficacy of *LLM-as-formalizer* on  
 196 various logical reasoning tasks. However, we observe that **even non-reasoning LLMs do not per-**  
 197 **form well as a formalizer** on our CSP domains. As shown in Figure 2, on relatively simpler  
 198 domains such as Calendar Scheduling and Zebra Logic, Qwen2.5, the weakest model we consider,  
 199 does exhibit a better performance as either a Python or SMT formalizer than a solver. The perfor-  
 200 mance gain more significant given the revision-by-error module which emulates the state-of-the-art  
 201 *LLM-as-formalizer* pipelines. However, the alignment of our results with reported results in other  
 202 tasks does not extend to more challenging tasks like Meeting Planning. A stronger non-reasoning  
 203 LLM, DeepSeek-V3, achieves the best performance as a solver than a formalizer of any formal-  
 204 ism by a large margin on every domain. Although non-reasoning LLMs are generally worse than  
 205 their reasoning counterparts, practitioners should not assume that they can generate high-quality  
 206 formalism that leads to the correct solution.

207 The trend of unsatisfactory performance of *LLM-as-formalizer* is even clearer for LRMs, as we  
 208 clearly show that **LRMs are consistently stronger as a planner than as a formalizer**. All 4  
 209 LRMs we consider, both open- and closed-source including Qwen3, DeepSeek-R1, o3-mini,  
 210 and gpt-5, achieve the best performance in every domain when generating the solution rather  
 211 than any kind of formalism even when revision is allowed. Understandably, the ability of LRMs  
 212 to reason about complex problems like those in our CSP domains has been greatly improved via  
 213

214 <sup>2</sup>[pypi.org/project/z3-solver](https://pypi.org/project/z3-solver)

215 <sup>3</sup>[huggingface.co](https://huggingface.co)

<sup>4</sup>[dottxt-ai.github.io/outlines](https://dottxt-ai.github.io/outlines)

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

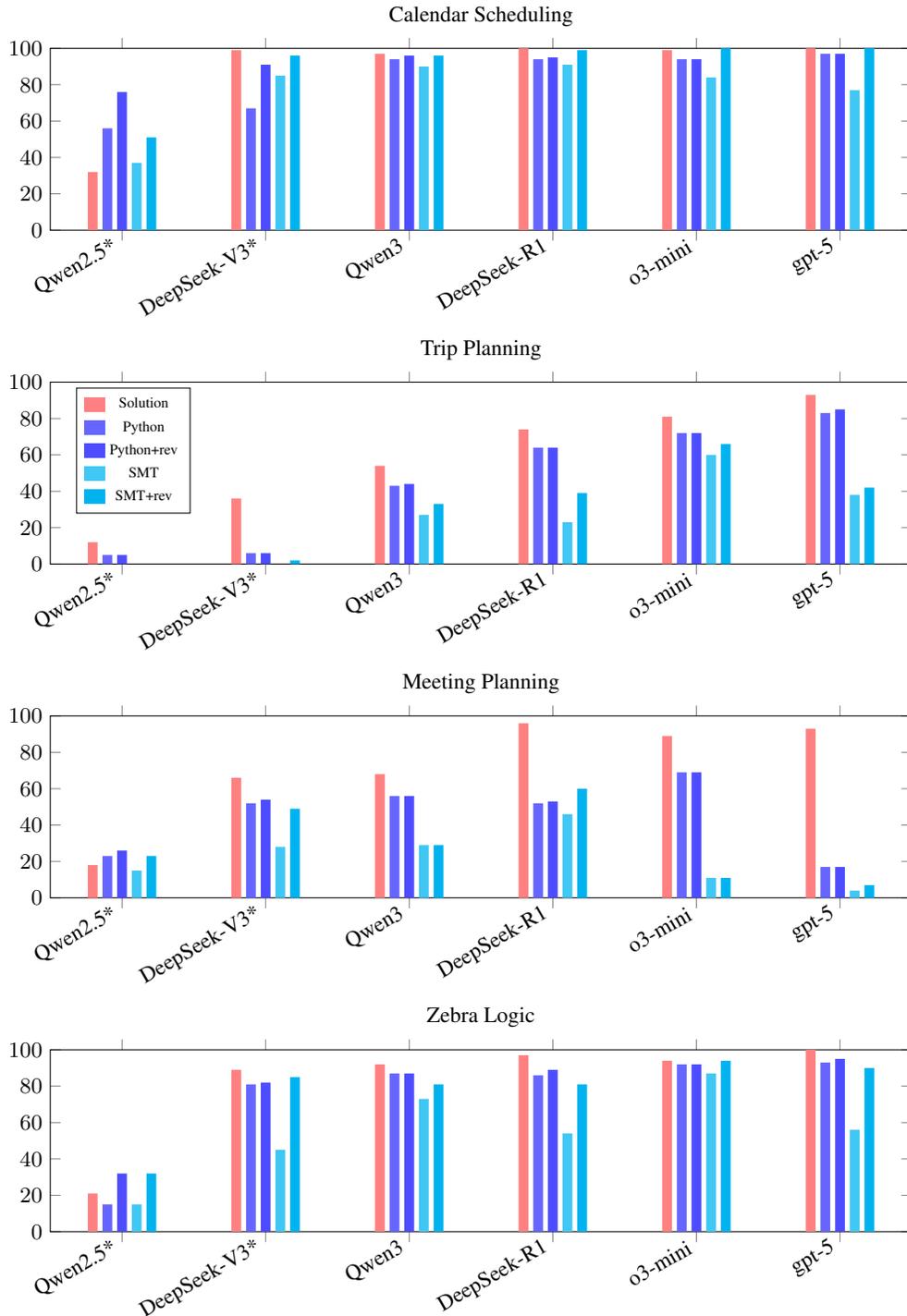


Figure 2: The percentage of correct plans (defined by formally passing all annotated constraints) by *LLM-as-solver* and *LLM-as-formalizer* generating both Python and SMT code of various LLMs on all 4 CSP domains. In the settings with revision, only solver errors including inability to find a plan induce revisions. LLMs that are not LRMs are marked by an asterisk (\*).

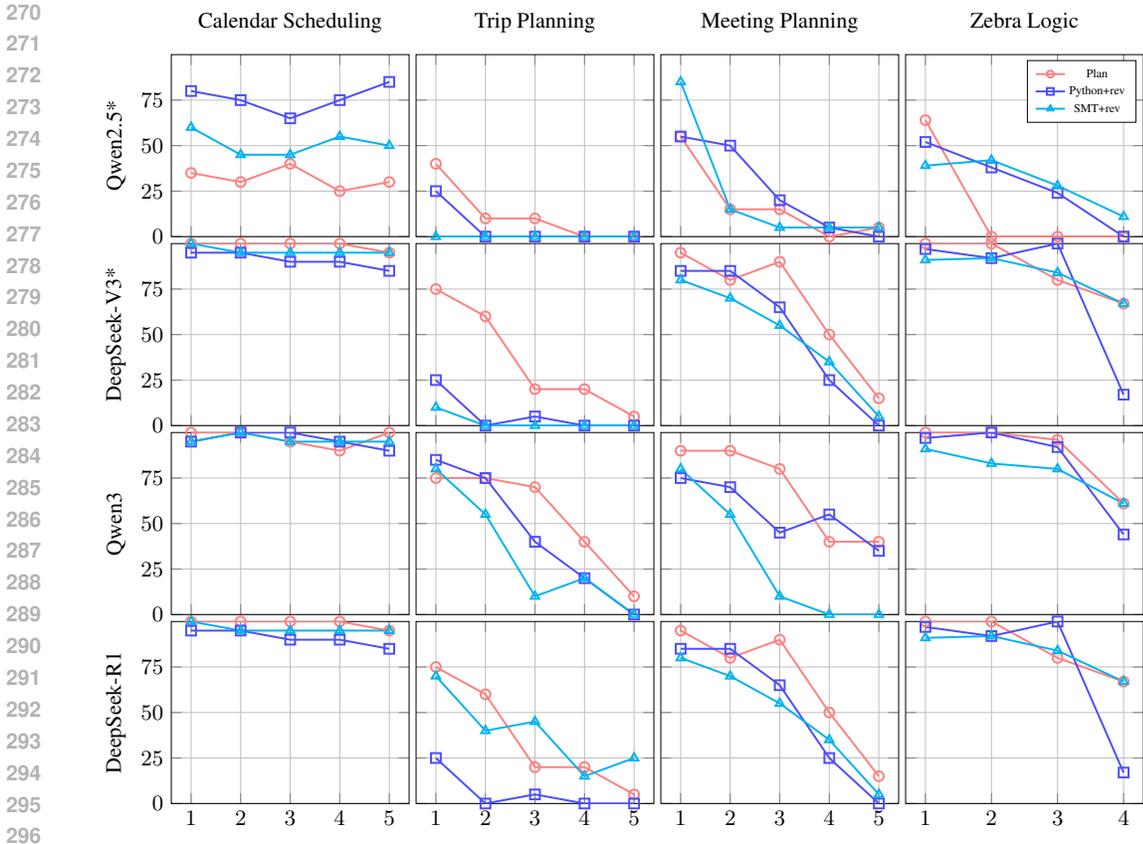


Figure 3: The change in the percentage of correct plans by *LLM-as-solver* and *LLM-as-formalizer* over buckets of examples stratified by complexity measured by the number of constraints. Data in the NaturalPlan domains is roughly equally partitioned into 5 percentiles, while that in ZebraLogic is partitioned based on the 4 strata provided by the dataset.

training. While their ability to generate code has also been reported to improve, the latter ability has not caught up with the former, calling for a more robust yet flexible model to generate formal representations beyond a particular language.

### 5.2 RQ2: IS *LLM-as-formalizer* MORE ROBUST TO COMPLEXITY THAN *LLM-as-solver*?

Even if *LLM-as-formalizer* is worse performing than *LLM-as-solver*, the former still promises many benefits over the latter, such as robustness to complexity. Most of our cited publications have pointed out non-reasoning LLMs’ failure to scale as a solver when the search space of a problem becomes intractable. We validate this claim as shown in Figure 3, as both Qwen2.5 and DeepSeek-V3 as planners (red lines) quickly fail as the problems become complex in all domains other than the easiest Calendar Scheduling. We also validate Shojaei et al. (2025) which claimed that LLMs also do not scale well as a solver, as both Qwen3 and DeepSeek-R1 display a similar trend, though their performance degradation tends to happen later and smoother than their non-reasoning counterparts.

We note that **neither does *LLM-as-formalizer* scale well with complexity**. Regardless of the formalism, the percentage of correctly generated Python (blue line) or SMT (cyan line) code decreases almost as much as *LLM-as-solver* for all models on all domains. This finding is an antithesis to the *LLM-as-formalizer* methodology, which is based on the premise that only using LLMs for translation and using formal solvers for the search for solution increases robustness. Our results show that LLMs are no more robust in translating a problem description into a formal program than in solving the problem directly, regardless of whether the model is an LRM or not.

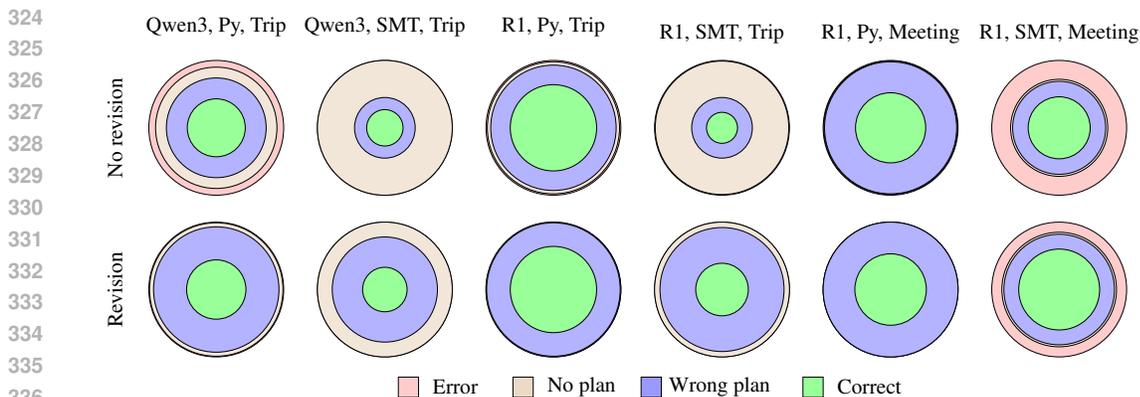


Figure 4: Error analysis of two LLMs Qwen3 and DeepSeek-R1 as both Python and SMT formalizers. Revision is only allowed in case of errors (red ring) and not being able to find a plan (brown ring) and thus have no bearing on existing wrong plans (blue ring).

### 5.3 RQ3: WHY DOES *LLM-as-formalizer* UNDERPERFORM AND FAIL TO SCALE?

To understand why LLMs cannot reliably generate formal programs, we consider 4 outcomes:

- **Error:** the program is unable to be executed due to an error such as `TypeError`, `RuntimeError`, or an error specific to the Z3 library.
- **No plan:** the solver, either defined by the model in the case of Python or provided by the Z3 library in case of SMT, outputs a custom message that no plan can be found.
- **Wrong plan:** the solver outputs a plan, but the plan does not pass all the constraint checks for `NaturalPlan` or does not exactly match the ground truth for `ZebraLogic`.
- **Correct plan:** the plan is correct.

Figure 4 shows a breakdown of errors of *LLM-as-formalizer* on Trip Planning and Meeting Planning, the two most challenging domains. We only consider the two open-source LLMs but not their non-reasoning counterparts due to close-to-zero performance. **Syntax errors are rare but still existent**, while the majority of them are eliminated given revision. Common syntax errors are diverse and include import errors, value error, and run-time errors, but are primarily dominated by erroneous use of APIs in the Z3 library for SMT generation. Similarly, the inability to find a plan, which is a kind of semantic errors, can be largely addressed with revision. Despite its effect on both kinds of errors, **revision does not increase overall correctness of the solution**. The persistent error is wrong plan, the other kind of semantic errors, which is not influenced by revision since the solver cannot provide signal on the correctness of a plan. In real-life applications, validating the correctness of plans assumes much more resources and risks than validating the existence of a plan.

To further study the root cause of semantic errors including both cases of “no plan” and “wrong plan”, we manually annotate 5 examples per domain, per Python or SMT, per Qwen3 or gpt-5 as a formalizer, 80 examples in total. We break down the semantic errors into 3 finer-grained categories:

- **Missing constraint:** the program misses the definition of a constraint otherwise present in the problem description.
- **Wrong constraint:** the program wrongly defines a constraint as presented in the problem description.
- **Wrong solver:** the program wrongly defines or calls the solver.

We count that across domains and formalisms, **the majority of the semantic errors are due to wrongly defining constraints** (62% of gpt-5 and 95% of Qwen3). Examples include treating 2 : 16 as an end time while it should be the start time for Calendar Scheduling, defining meeting time to be 9 : 16 while it should be 9 : 18 for Meeting plan, and so on. Such trivial errors of information extraction still plague the state-of-the-art LLMs and cannot be caught by the error message provided by the solver, posing great safety concern in real-life applications. This calls for the need for a specific translation module beyond solely relying on LLMs to perform the translation from the

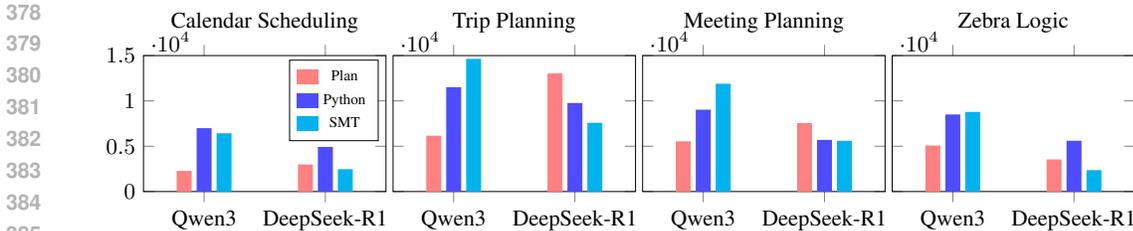


Figure 5: Number of reasoning tokens generated LLMs as a solver and as a formalizer.

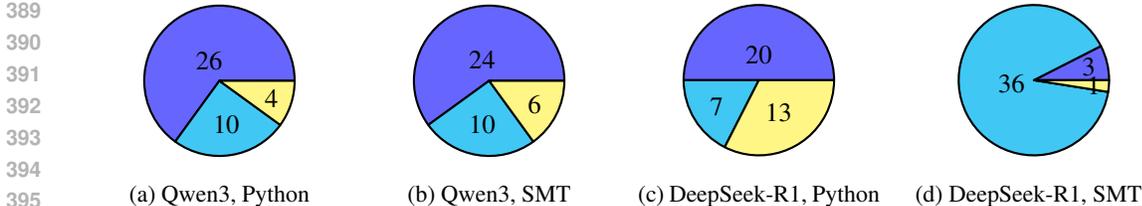


Figure 6: The distribution of different types of reasoning chains (sound, unnecessary, and spurious) in all domains for both LLMs and both formalisms.

problem description to the formal program. The other two categories, missing constraint and wrong solver are much less common though still existent.

#### 5.4 RQ4: ARE LLMs MORE EFFICIENT AS A FORMALIZER THAN AS A PLANNER?

It is known that LLMs’ impressive performance as a planner comes at a cost of excessive and often inefficient reasoning chains (Sui et al., 2025). On the other hand, *LLM-as-formalizer* may be expected to generate fewer reasoning tokens because the task of formalizing is  $O(n + d)$  with  $n$  variables and a domain size of  $d$ , as opposed to the worst case  $O(d^n)$  for an NP-complete problem like CSP. However, we show that surprisingly **LLMs often generate more reasoning tokens as a formalizer than as a solver**. From Figure 5, Qwen3 generates double the amount of reasoning tokens as an SMT formalizer compared to as a solver, likely because the Z3 formalism and syntax are unfamiliar given the pre-training. The number of tokens in the case of Python is less, but still significantly more than the case of *LLM-as-solver*. On the other hand, DeepSeek-R1 as a formalizer is more efficient in Trip Planning and Meeting Planning despite worse performance.

Beyond token counts, we also examine the reasoning chains of *LLM-as-formalizer* by breaking them down into the following granular categories:

- **Sound reasoning** focuses on coding and anticipates errors, edge cases, fixes.
- **Unnecessary reasoning** unnecessarily tries to solve the problem like *LLM-as-solver*, but still generates plausible code that declares the constraints and calls a solver.
- **Spurious reasoning** hard-codes the answer within the generated code, without declaring the constraints and calling a solver. This is a special case of unnecessary reasoning and takes precedence when counting.

We manually annotate the reasoning chains of DeepSeek-R1 and Qwen3 for 10 examples per domain per formalism, resulting in 160 labels in total. Figure 6 reveals a concerning phenomenon where only less than 75% of the reasoning chains of Qwen3 and less than 50% of that of DeepSeek-R1 are sound. **Unnecessary reasoning is common, where models “think like a solver”** to perform search, enumeration, and backtracking to attempt to solve the problem rather than formalizing it. We count that the average number of tokens generated for unnecessary reasoning is 5,933 over 3,624 for sound reasoning in the case of DeepSeek-R1, and 13,010 over 8,537 in the case of Qwen3. This more than 1.5 times of token consumption sheds light on the low performance and low efficiency of *LLM-as-formalizer* on our CSP domains.

```

432 1 # Reasoning chain: ... Let me think about a possible sequence. Maybe the
433 2 user should visit Rome on the first day? Wait, but that will violate the
434 3 constraint of being in Dubrovnik for a conference. Let me think again...
435 4 import json
436 5 def main():
437 6     itinerary = [
438 7         {"day_range": "Day 1-2", "place": "Dubrovnik"},
439 8         {"day_range": "Day 3", "place": "Rome (travel from Dubrovnik)"},
440 9         {"day_range": "Day 4", "place": "Riga (travel from Rome)"},
441 10        # More days omitted
442 11    ]
443 12    print(json.dumps({"itinerary": itinerary}))
444 13 if __name__ == "__main__":
445 14     main()

```

Figure 7: An example of hard-coded program resulting from a spurious reasoning chain (abridged) of DeepSeek-R1 as a Python formalizer on Trip Planning.

Despite those shortcomings, one of the built-in advantages of *LLM-as-formalizer* is its improved interpretability, verifiability, and faithfulness, as previously discussed. However, spurious reasoning, the most insidious category of reasoning chains, nullifies this advantage. Figure 7 shows an example of spurious reasoning, where the model not only generates a solver-like, unnecessary reasoning chain, but also **hard-codes the entire proposed solution in the output program** without performing any declaration of constraints or implementation of search. Regardless of the correctness, a practitioner would have no way of interpreting and verifying such a hard-coded program. Spurious reasoning constitutes as much as 90% of SMT programs generated by DeepSeek-R1 in all domains, revealing a considerable concern in safety.

## 6 DISCUSSION

In response to the current line of effort in the community to develop *LLM-as-formalizer*, our findings from a systematic evaluation on CSP domains indicate its present shortcomings. While *LLM-as-formalizer* is intended to be accurate, robust, feasible, and efficient, we observe that even state-of-the-art LLMs currently fall short on all four dimensions. At a coarse-grained level, the root cause of failure is LLMs' inability to generate low-resource, domain-specific, and formal languages (Joel et al., 2024). Ongoing work is using techniques like intermediate representation (Zhang et al., 2025), grammar (Wang et al., 2023a) to improve performance, which is likely to benefit *LLM-as-formalizer* as a whole. At a fine-grained level, major issues like mis-specifying constraints could be mitigated through modular generation and local verification. However, this effort must be cautioned against overfitting, as different formal languages may require different techniques to generate them. This explains why we deliberately choose to only evaluate a general pipeline with few-shot generation and revision-by-error. The advance of LRMs brings a promising outlook to not only *LLM-as-solver*, but also *LLM-as-formalizer* due to their reported improved performance in code generation. While we also observe such improvements in our evaluation compared to non-reasoning LLMs, we also identify serious issues such as inefficient and spurious reasoning before generating programs. The root cause of these issues is that LRMs are trained more to solve a problem, not to formalize it. Present work has just begun to explore ways to steer the models' behavior from problem solving to formalization (Chen et al., 2025). In addition to fine-tuning and prompting, attention steering might also be a reasonable approach for future work (Tian & Zhang, 2025).

## 7 CONCLUSION

We provide a comprehensive evaluation of LLMs' ability to solve real-life constraint satisfaction problems both as solvers and as formalizers. We show that current techniques of *LLM-as-formalizer* fall short of their promise to be high-performing and trustworthy. Our diverse findings, fine-grained analyses, and actionable remedies are intended to drive the community's progress on further development of this promising neuro-symbolic methodology of complex-problem solving.

## REFERENCES

- 486  
487  
488 Alexander Beiser, David Penz, and Nysret Musliu. Intermediate languages matter: Formal choice  
489 drives neurosymbolic llm reasoning, 2025. URL <https://arxiv.org/abs/2502.17216>.
- 490 Shmuel Berman, Kathleen McKeown, and Baishakhi Ray. Solving zebra puzzles using constraint-  
491 guided multi-agent systems, 2024. URL <https://arxiv.org/abs/2407.03956>.
- 492  
493 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
494 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,  
495 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.  
496 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,  
497 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,  
498 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proceedings of the*  
499 *34th International Conference on Neural Information Processing Systems, NIPS '20*, Red Hook,  
500 NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 501 Yongchao Chen, Harsh Jhamtani, Srinagesh Sharma, Chuchu Fan, and Chi Wang. Steering large  
502 language models between code execution and textual reasoning. In *The Thirteenth International*  
503 *Conference on Learning Representations*, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=5X5Z7Ffrjb)  
504 [id=5X5Z7Ffrjb](https://openreview.net/forum?id=5X5Z7Ffrjb).
- 505  
506 Peter Clark, Oyvind Tafjord, and Kyle Richardson. Transformers as soft reasoners over language.  
507 In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJ-*  
508 *CAI'20*, 2021. ISBN 9780999241165.
- 509 Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing  
510 the mystery behind chain of thought: A theoretical perspective. In *Thirty-seventh Conference on*  
511 *Neural Information Processing Systems*, 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=qHrADgAdYu)  
512 [id=qHrADgAdYu](https://openreview.net/forum?id=qHrADgAdYu).
- 513 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and  
514 Graham Neubig. Pal: program-aided language models. In *Proceedings of the 40th International*  
515 *Conference on Machine Learning, ICML'23*. JMLR.org, 2023.
- 516  
517 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu  
518 Zhang, Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-  
519 hong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng,  
520 Chengda Lu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Deli Chen, Dongjie  
521 Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei  
522 Li, H. Zhang, Hanwei Xu, Honghui Ding, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Ji-  
523 ashu Li, Jingchang Chen, Jingyang Yuan, Jinhao Tu, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi  
524 Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaichao You, Kaige Gao, Kang Guan, Kexin  
525 Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu,  
526 Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingxu Zhou, Meng Li, Miaojun  
527 Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen,  
528 Qiusi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi  
529 Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuip-  
530 ing Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Tao Yun, Tian  
531 Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin  
532 Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xi-  
533 aotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng  
534 Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaox-  
535 iang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang,  
536 Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu,  
537 Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan,  
538 Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zhou,  
539 Yujia He, Yufan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu,  
Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yut-  
ing Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang,  
Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li,

- 540 Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen  
541 Zhang. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*,  
542 645(8081):633–638, Sep 2025. ISSN 1476-4687. doi: 10.1038/s41586-025-09422-z. URL  
543 <https://doi.org/10.1038/s41586-025-09422-z>.  
544
- 545 Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Wenfei Zhou, James  
546 Coady, David Peng, Yujie Qiao, Luke Benson, Lucy Sun, Alexander Wardle-Solano, Hannah  
547 Szabó, Ekaterina Zubova, Matthew Burtell, Jonathan Fan, Yixin Liu, Brian Wong, Malcolm  
548 Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Alexander Fabbri, Woj-  
549 ciech Maciej Kryscinski, Semih Yavuz, Ye Liu, Xi Victoria Lin, Shafiq Joty, Yingbo Zhou, Caim-  
550 ing Xiong, Rex Ying, Arman Cohan, and Dragomir Radev. FOLIO: Natural language reasoning  
551 with first-order logic. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceed-  
552 ings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 22017–  
553 22031, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL  
554 <https://aclanthology.org/2024.emnlp-main.1229>.
- 555 Yilun Hao, Yang Zhang, and Chuchu Fan. Planning anything with rigor: General-purpose zero-shot  
556 planning with LLM-based formalized programming. In *The Thirteenth International Confer-  
557 ence on Learning Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=0K10aL6XuK)  
558 [0K10aL6XuK](https://openreview.net/forum?id=0K10aL6XuK).
- 559 Cassie Huang and Li Zhang. On the limit of language models as planning formalizers. In Wanxi-  
560 ang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings  
561 of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
562 Papers)*, pp. 4880–4904, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. URL [https://aclanthology.org/2025.acl-long.](https://aclanthology.org/2025.acl-long.242/)  
564 [242/](https://aclanthology.org/2025.acl-long.242/).
- 565 Susmit Jha, Sumit Kumar Jha, Patrick Lincoln, Nathaniel D. Bastian, Alvaro Velasquez, and  
566 Sandeep Neema. Dehallucinating large language models using formal methods guided iterative  
567 prompting. In *2023 IEEE International Conference on Assured Autonomy (ICAA)*, pp. 149–152,  
568 2023. doi: 10.1109/ICAA58325.2023.00029.  
569
- 570 Albert Qiaochu Jiang, Sean Welleck, Jin Peng Zhou, Timothee Lacroix, Jiacheng Liu, Wenda Li,  
571 Mateja Jamnik, Guillaume Lample, and Yuhuai Wu. Draft, sketch, and prove: Guiding formal  
572 theorem provers with informal proofs. In *The Eleventh International Conference on Learning  
573 Representations*, 2023. URL <https://openreview.net/forum?id=SMa9EAovKMC>.
- 574 Sathvik Joel, Jie JW Wu, and Fatemeh H. Fard. A survey on llm-based code generation for low-  
575 resource and domain-specific programming languages, 2024. URL [https://arxiv.org/  
576 abs/2410.03981](https://arxiv.org/abs/2410.03981).
- 577
- 578 Prabhu Prakash Kagitha, Andrew Zhu, and Li Zhang. Addressing the challenges of planning lan-  
579 guage generation, 2025. URL <https://arxiv.org/abs/2505.14763>.
- 580 Pascal Kesseli, Peter O’Hearn, and Ricardo Silveira Cabral. Logic.py: Bridging the gap between  
581 llms and constraint solvers, 2025. URL <https://arxiv.org/abs/2502.15776>.  
582
- 583 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
584 language models are zero-shot reasoners. In *Proceedings of the 36th International Conference on  
585 Neural Information Processing Systems, NIPS ’22*, Red Hook, NY, USA, 2022. Curran Associates  
586 Inc. ISBN 9781713871088.
- 587 Kuang-Huei Lee, Ian Fischer, Yueh-Hua Wu, Dave Marwood, Shumeet Baluja, Dale Schuurmans,  
588 and Xinyun Chen. Evolving deeper llm thinking, 2025. URL [https://arxiv.org/abs/  
589 2501.09891](https://arxiv.org/abs/2501.09891).
- 590
- 591 Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter  
592 Clark, and Yejin Choi. Zebralogic: On the scaling limits of LLMs for logical reasoning. In *Forty-  
593 second International Conference on Machine Learning*, 2025. URL [https://openreview.](https://openreview.net/forum?id=sTAJ9QyA6l)  
[net/forum?id=sTAJ9QyA6l](https://openreview.net/forum?id=sTAJ9QyA6l).

- 594 Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone.  
595 Llm+ p: Empowering large language models with optimal planning proficiency. *arXiv preprint*  
596 *arXiv:2304.11477*, 2023.
- 597  
598 Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki,  
599 and Chris Callison-Burch. Faithful chain-of-thought reasoning. In Jong C. Park, Yuki Arase, Bao-  
600 tian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi (eds.), *Proceedings of*  
601 *the 13th International Joint Conference on Natural Language Processing and the 3rd Conference*  
602 *of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long*  
603 *Papers)*, pp. 305–329, Nusa Dua, Bali, November 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.ijcnlp-main.20. URL <https://aclanthology.org/2023.ijcnlp-main.20>.
- 604  
605 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri  
606 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad  
607 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine:  
608 Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Pro-*  
609 *cessing Systems*, 2023. URL <https://openreview.net/forum?id=S37hOerQLB>.
- 610  
611 Long Hei Matthew Lam, Ramya Keerthy Thatikonda, and Ehsan Shareghi. A closer look at tool-  
612 based logical reasoning with LLMs: The choice of tool matters. In Tim Baldwin, Sergio José  
613 Rodríguez Méndez, and Nicholas Kuo (eds.), *Proceedings of the 22nd Annual Workshop of the*  
614 *Australasian Language Technology Association*, pp. 41–63, Canberra, Australia, December 2024.  
615 Association for Computational Linguistics. URL <https://aclanthology.org/2024.alt-1.4/>.
- 616  
617 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke  
618 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time  
619 scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- 620  
621 Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. Logic-LM: Empowering large  
622 language models with symbolic solvers for faithful logical reasoning. In Houda Bouamor, Juan  
623 Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP*  
624 *2023*, pp. 3806–3824, Singapore, December 2023. Association for Computational Linguistics.  
625 doi: 10.18653/v1/2023.findings-emnlp.248. URL <https://aclanthology.org/2023.findings-emnlp.248>.
- 626  
627 Mihir Parmar, Xin Liu, Palash Goyal, Yanfei Chen, Long Le, Swaroop Mishra, Hossein Mobahi,  
628 Jindong Gu, Zifeng Wang, Hootan Nakhost, Chitta Baral, Chen-Yu Lee, Tomas Pfister, and Hamid  
629 Palangi. Plangen: A multi-agent framework for generating planning and reasoning trajectories for  
630 complex problem solving, 2025. URL <https://arxiv.org/abs/2502.16111>.
- 631  
632 Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. Explain yourself!  
633 leveraging language models for commonsense reasoning. In Anna Korhonen, David Traum, and  
634 Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computa-*  
635 *tional Linguistics*, pp. 4932–4942, Florence, Italy, July 2019. Association for Computational Lin-  
636 guistics. doi: 10.18653/v1/P19-1487. URL <https://aclanthology.org/P19-1487>.
- 637  
638 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis  
639 of chain-of-thought. In *The Eleventh International Conference on Learning Representations*,  
640 2023. URL <https://openreview.net/forum?id=qFVVBzXxR2V>.
- 641  
642 David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical rea-  
643 soning abilities of neural models. In *International Conference on Learning Representations*, 2019.  
644 URL <https://openreview.net/forum?id=H1gR5iR5FX>.
- 645  
646 Bilgehan Sel, Ahmad Al-Tawaha, Vanshaj Khattar, Ruoxi Jia, and Ming Jin. Algorithm of thoughts:  
647 enhancing exploration of ideas in large language models. In *Proceedings of the 41st International*  
*Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning

- 648 models via the lens of problem complexity, 2025. URL <https://arxiv.org/abs/2506.06941>.
- 649
- 650
- 651 Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. On the self-verification limita-  
652 tions of large language models on reasoning and planning tasks. In *The Thirteenth International*  
653 *Conference on Learning Representations*, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=400v4s3IzY)  
654 [id=400v4s3IzY](https://openreview.net/forum?id=400v4s3IzY).
- 655 Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu,  
656 Andrew Wen, Shaochen Zhong, Na Zou, Hanjie Chen, and Xia Hu. Stop overthinking: A survey  
657 on efficient reasoning for large language models. *Transactions on Machine Learning Research*,  
658 2025. ISSN 2835-8856. URL <https://openreview.net/forum?id=HvoG8SxggZ>.
- 659 Yuan Tian and Tianyi Zhang. Selective prompt anchoring for code generation, 2025. URL <https://arxiv.org/abs/2408.09121>.
- 660
- 661 Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don’t always  
662 say what they think: unfaithful explanations in chain-of-thought prompting. In *Proceedings of the*  
663 *37th International Conference on Neural Information Processing Systems*, NIPS ’23, Red Hook,  
664 NY, USA, 2023. Curran Associates Inc.
- 665
- 666 Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kamb-  
667 hampati. Planbench: An extensible benchmark for evaluating large language models on planning  
668 and reasoning about change. *Advances in Neural Information Processing Systems*, 36, 2024.
- 669
- 670 Bailin Wang, Zi Wang, Xuezhi Wang, Yuan Cao, Rif A. Saurous, and Yoon Kim. Grammar prompt-  
671 ing for domain-specific language generation with large language models. In *Proceedings of the*  
672 *37th International Conference on Neural Information Processing Systems*, NIPS ’23, Red Hook,  
673 NY, USA, 2023a. Curran Associates Inc.
- 674
- 675 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha  
676 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language  
677 models. In *The Eleventh International Conference on Learning Representations*, 2023b. URL  
<https://openreview.net/forum?id=1PL1NIMMrw>.
- 678
- 679 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
680 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
681 models. In *Proceedings of the 36th International Conference on Neural Information Processing*  
*Systems*, NIPS ’22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- 682
- 683 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun  
684 Zhao. Large language models are better reasoners with self-verification. In Houda Bouamor, Juan  
685 Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP*  
686 *2023*, pp. 2550–2575, Singapore, December 2023. Association for Computational Linguistics.  
687 doi: 10.18653/v1/2023.findings-emnlp.167. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-emnlp.167)  
[findings-emnlp.167](https://aclanthology.org/2023.findings-emnlp.167).
- 688
- 689 Yuhuai Wu, Albert Q. Jiang, Wenda Li, Markus N. Rabe, Charles Staats, Mateja Jamnik, and Chris-  
690 tian Szegedy. Autoformalization with large language models. In *Proceedings of the 36th Inter-*  
691 *national Conference on Neural Information Processing Systems*, NIPS ’22, Red Hook, NY, USA,  
2022. Curran Associates Inc. ISBN 9781713871088.
- 692
- 693 Yaqi Xie, Chen Yu, Tongyao Zhu, Jinbin Bai, Ze Gong, and Harold Soh. Translating natural lan-  
694 guage to planning goals with large-language models. *arXiv preprint arXiv:2302.05128*, 2023.
- 695
- 696 Kaiyu Yang, Aidan M. Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil,  
697 Ryan Prenger, and Anima Anandkumar. Leandojo: theorem proving with retrieval-augmented  
698 language models. In *Proceedings of the 37th International Conference on Neural Information*  
*Processing Systems*, NIPS ’23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- 699
- 700 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik  
701 Narasimhan. Tree of thoughts: deliberate problem solving with large language models. In *Pro-*  
*ceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS  
’23, Red Hook, NY, USA, 2023. Curran Associates Inc.

702 Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. Satlm: satisfiability-aided language models  
703 using declarative prompting. In *Proceedings of the 37th International Conference on Neural*  
704 *Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2023. Curran Associates Inc.

706 Jipeng Zhang, Jianshu Zhang, Yuanzhe Li, Renjie Pi, Rui Pan, Runtao Liu, Zheng Ziqiang,  
707 and Tong Zhang. Bridge-coder: Transferring model capabilities from high-resource to low-  
708 resource programming language. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and  
709 Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics:*  
710 *ACL 2025*, pp. 10865–10882, Vienna, Austria, July 2025. Association for Computational Lin-  
711 guistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.567. URL [https://](https://aclanthology.org/2025.findings-acl.567/)  
712 [aclanthology.org/2025.findings-acl.567/](https://aclanthology.org/2025.findings-acl.567/).

713 Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang, Xinyun Chen, Minmin Chen, Azade Nova,  
714 Le Hou, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and Denny Zhou. Natural plan: Benchmarking  
715 llms on natural language planning. 2024. URL <https://arxiv.org/abs/2406.04520>.

716 Andrew Zhu, Liam Dugan, Alyssa Hwang, and Chris Callison-Burch. Kani: A lightweight and  
717 highly hackable framework for building language model applications, 2023. URL [https://](https://arxiv.org/abs/2309.05542)  
718 [arxiv.org/abs/2309.05542](https://arxiv.org/abs/2309.05542).

## 720 A DATA AND PROMPT EXAMPLES

722 We use the prompt wordings from NaturalPlan as is, but provide some examples here for the  
723 paper to be self-contained. Here is an example from the calendar scheduling task:

724 You are an expert at scheduling meetings. You are given a few  
725 constraints on the existing schedule of each participant, the meeting  
726 duration, and possibly some preferences on the meeting time. Note there  
727 exists a solution that works with existing schedule of every  
728 participant. Here are a few example tasks and solutions: TASK: You need  
729 to schedule a meeting for James and John for one hour between the work  
730 hours of 9:00 to 17:00 on Monday. Here are the existing schedules for  
731 everyone during the day: James has blocked their calendar on Monday  
732 during 11:30 to 12:00, 14:30 to 15:00; John is busy on Monday during  
733 9:30 to 11:00, 11:30 to 12:00, 12:30 to 13:30, 14:30 to 16:30; Find a  
734 time that works for everyone’s schedule and constraints. Please provide  
735 your solution in a JSON format as as  
736 {"start":{"day":"Monday","time"13:30"},  
737 "end":{"day":"Monday","time"14:30}}.

738 To facilitate parsing, we add the last sentence to encourage an output in JSON. Here is an example  
739 from the trip planning task:

740 You plan to visit 3 European cities for 7 days in total. You only take  
741 direct flights to commute between cities. You want to spend 4 days in  
742 Madrid. You would like to visit Dublin for 3 days. You want to spend 2  
743 days in Tallinn. You have to attend a workshop in Tallinn between day 6  
744 and day 7. Here are the cities that have direct flights: Madrid and  
745 Dublin, Dublin and Tallinn. Find a trip plan of visiting the cities for  
746 7 days by taking direct flights to commute between them. Please provide  
747 your solution in a JSON format as as {"itinerary": [{"day\_range": "Day  
748 1-2", "place": "Reykjavik"}, {"day\_range": "Day 2-4", "place":  
749 "Stockholm"}.....]}.

750 Here is an example from the meeting planning task:

751 You are visiting San Francisco for the day and want to meet as many  
752 friends as possible. Solve the problem by considering various different  
753 schedules and picking the best one to optimize your goals. Travel  
754 distances (in minutes): Sunset District to Chinatown: 30. Sunset  
755 District to Russian Hill: 24. Sunset District to North Beach: 29.  
Chinatown to Sunset District: 29. Chinatown to Russian Hill: 7.  
Chinatown to North Beach: 3. Russian Hill to Sunset District: 23.

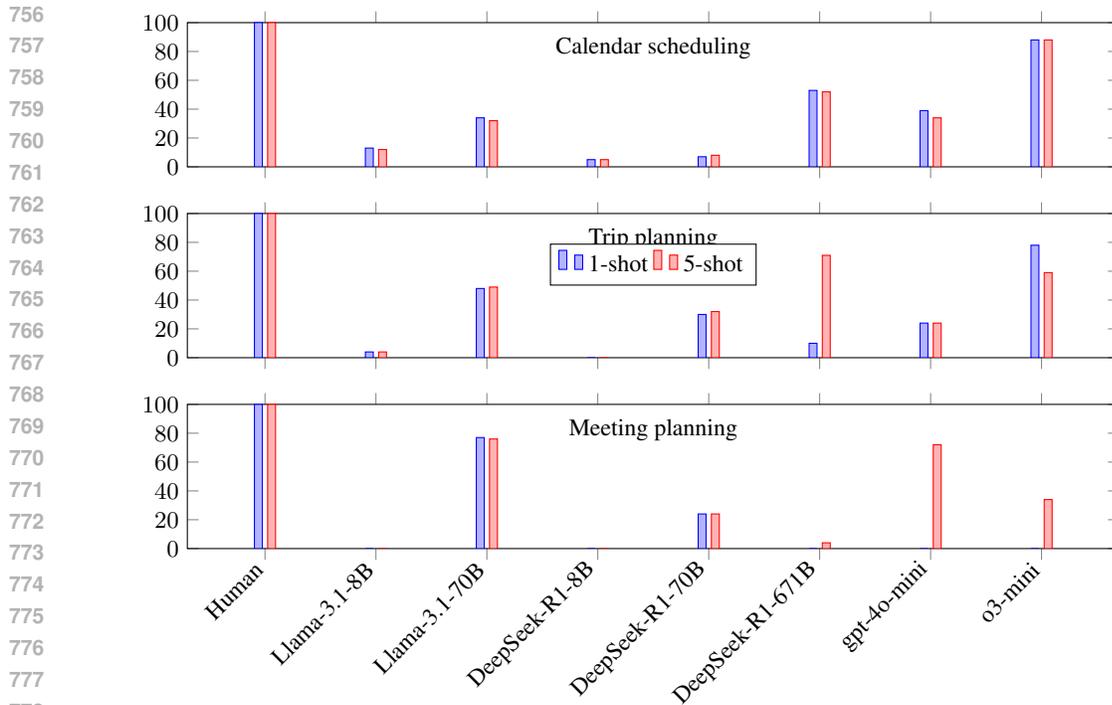


Figure 8: The performance of preliminary studies of 1-shot and 5-shot prompting on all 3 domains of NaturalPlan for *LLM-as-solver*.

Russian Hill to Chinatown: 9. Russian Hill to North Beach: 5. North Beach to Sunset District: 27. North Beach to Chinatown: 6. North Beach to Russian Hill: 4. CONSTRAINTS: You arrive at Sunset District at 9:00AM. Anthony will be at Chinatown from 1:15PM to 2:30PM. You'd like to meet Anthony for a minimum of 60 minutes. Rebecca will be at Russian Hill from 7:30PM to 9:15PM. You'd like to meet Rebecca for a minimum of 105 minutes. Melissa will be at North Beach from 8:15AM to 1:30PM. You'd like to meet Melissa for a minimum of 105 minutes. Please provide your solution in a JSON format as as {"itinerary": [{"action": "meet", "location": "Golden Gate Park", "person": "David", "start\_time": "13:00", "end\_time": "14:00"}, ...]}.

For program output, we prompt the LLM to write code which produces the output in the format specified above, though this cannot be guaranteed. We perform simple post-processing with the aid of a small LLM, gpt-4.1-mini.

## B FEW-SHOT VERSUS ZERO-SHOT

In our preliminary experiment, we compared both 5-shot prompting and 1-shot prompting (default of this paper). For the former, we use the 5-shot prompts provided by NaturalPlan. We evaluate only generating plan on all 3 tasks, because in-context exemplars for code generation requires expertise in annotation and is thus unrealistic. From Figure 8, 5-shot prompting does not increase performance in most cases, while occasionally hurting performance. We inspected its performance gain for some models in meeting planning, and conclude that many errors in 1-shot settings are due to not conforming to a specific output format. These errors are mitigated via 5-shot prompting, while adding little value to insights into reasoning. In later experiments, we relaxed the stringency on output format and rely more on post-processing for a fairer evaluation.

810	Listing (1) Python solution	Listing (2) Z3 solution
811		<code>from z3 import Int, Optimize</code>
812		
813	<code>meeting_duration = 1</code>	<code>meeting_duration = 1</code>
814	<code>working_start = 9</code>	<code>working_start = 9</code>
815	<code>working_end = 17</code>	<code>working_end = 17</code>
816	<code># declare disallowed time</code>	
817	<code>unavailable = [(9, 10)]</code>	<code># declare the start time variable</code>
818		<code>start = Int('start')</code>
819	<code># declare the start time variable</code>	
820	<code># declare allowed time</code>	<code>opt = Optimize()</code>
821	<code>for start in range(working_start,</code>	<code># declare disallowed time</code>
822	<code>working_end - meeting_duration +</code>	<code>opt.add(start != 9)</code>
823	<code>1):</code>	
824	<code>    conflict = False</code>	<code># declare allowed time</code>
825	<code>    # check disallowed time</code>	<code>opt.add(start &gt;= working_start)</code>
826	<code>    for u_start, u_end in</code>	<code>opt.add(start + meeting_duration</code>
827	<code>unavailable:</code>	<code>&lt;= working_end)</code>
828	<code>        if not (start +</code>	
829	<code>meeting_duration &lt;= u_start or</code>	<code># find the earliest time</code>
830	<code>start &gt;= u_end):</code>	<code>opt.minimize(start)</code>
831	<code>            conflict = True</code>	
832	<code>            break</code>	<code>if opt.check() == 0:</code>
833	<code>        if conflict:</code>	<code>    m = opt.model()</code>
834	<code>            continue</code>	<code>    s = m[start].as_long()</code>
835	<code>        # find the earliest time</code>	<code>    print(f"Earliest meeting:</code>
836	<code>        print(f"Earliest meeting:</code>	<code>{s}:00 to {s +</code>
837	<code>{start}:00 to {start +</code>	<code>meeting_duration):00")</code>
838	<code>meeting_duration}:00")</code>	<code>else:</code>
839	<code>        break</code>	<code>    print("No solution found.")</code>

Figure 9: Side-by-side comparison of pure Python and Z3-Python code.

## C PYTHON AND Z3 CODE

To showcase the difference between the general Python syntax and the Z3-Python syntax, consider the following example input:

Find the earliest 1-hour meeting time between 9AM-5PM. 9AM-10AM is not available.

The correct solution (generated by o3-mini) is shown in Figure 9. Both include two distinctive components of declarative code and algorithmic code.

## D MORE ERROR ANALYSIS

Another example of a hard-coded solution due to spurious reasoning is shown in Figure 10.

## E LICENSE

The NaturalPlan dataset we use is licensed under Apache License 2.0.

```

864
865 import json
866 def compute_itinerary():
867     # Input parameters
868     total_days = 16
869     days_in_lyon = 7
870     days_in_bucharest = 7
871     days_in_porto = 4
872     # Available direct flights:
873     # Bucharest <-> Lyon, Lyon <-> Porto
874     # Additional constraint:
875     # - There is a wedding in Bucharest between day 1 and day 7, so one
876     #   must be in Bucharest
877     #   during that period.
878     # - Note: if one flies from a city A to city B on a day, that day
879     #   counts as a day in both A and B.
880     # We have three cities to visit:
881     #   Bucharest, Lyon, and Porto.
882     # With 2 flights (transitions), the overlapping days count for both
883     #   cities.
884     # Let the flight from Bucharest to Lyon be on day f1 and the flight
885     #   from Lyon to Porto be on day f2.
886     # The days spent in each city are:
887     #   Bucharest: days 1 to f1 (inclusive)  -> count = f1
888     #   Lyon: day f1 (arrival day) + full days from (f1+1) to (f2-1) +
889     #   day f2 (departure day)
890     #       -> count = 1 + (f2 - f1 - 1) + 1 = f2 - f1 + 1
891     #   Porto: day f2 (arrival day) + days (f2+1) to total_days
892     #       -> count = 1 + (total_days - f2)
893     # We want:
894     #   f1 = days_in_bucharest = 7 (so that days 1-7 cover Bucharest,
895     #     ensuring the wedding is attended in the first 7 days)
896     #   For Lyon: f2 - f1 + 1 = days_in_lyon = 7
897     #       => f2 - 7 + 1 = 7  -> f2 = 7 + 6 = 13
898     #   Then Porto: 1 + (total_days - f2) = 1 + (16 - 13) = 4 =
899     #   days_in_porto
900     # This gives:
901     flight_day_bucharest_to_lyon = 7 # f1
902     flight_day_lyon_to_porto = 13 # f2
903     # Building the itinerary:
904     # The trip plan segments:
905     # 1. Bucharest from day 1 to day 7 (day 7 is used for the flight and
906     #   counts as Bucharest)
907     # 2. Lyon from day 7 to day 13 (day 7 from arrival flight, day 13
908     #   for departure flight)
909     # 3. Porto from day 13 to day 16 (day 13 counts as the arrival day)
910     itinerary = [
911         {"day_range": "1-7", "place": "Bucharest"},
912         {"day_range": "7-13", "place": "Lyon"},
913         {"day_range": "13-16", "place": "Porto"}
914     ]
915     # Output result as a JSON-formatted dictionary
916     result = {"itinerary": itinerary}
917     return result
918 def main():
919     itinerary_plan = compute_itinerary()
920     print(json.dumps(itinerary_plan))
921 if __name__ == '__main__':
922     main()
923
924
925
926
927

```

Figure 10: Another Python program generated by o3-mini on trip planning which hard-codes the solution.