LLMs as Translators, Not Thinkers: Structured Output Enables Stronger NP-Hard Problem Solving

₄ Haoyu Wang ⊠ 🗅

5 University of Pennsylvania, USA

🐻 Dan Roth 🖂 🕩

7 University of Pennsylvania, USA

8 — Abstract -

9 We study how large language models (LLMs) tackle NP-hard problems described in natural language,

¹⁰ comparing direct reasoning, search-based refinement, and translation to structured solver inputs.

¹¹ Using models like GPT-40 and OpenAI o1, we analyze trade-offs between accuracy and efficiency.

12 Our results highlight both the potential and limitations of LLMs in constraint-based tasks, offering

¹³ insights for hybrid LLM-solver systems¹.

- 14 2012 ACM Subject Classification Replace ccsdesc macro with valid one
- 15 Keywords and phrases LLM, NP-hard, Combinatorial Optimization
- ¹⁶ Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

17 **Introduction**

LLMs are increasingly being applied to 18 combinatorial problems, offering natural 19 language interfaces and flexible reasoning 20 capabilities. Yet, their behavior on classic-21 ally hard problems—such as Hamiltonian 22 Cycle and Traveling Salesman—remains 23 poorly understood. In this work, we in-24 vestigate how LLMs solve NP-hard prob-25 lems across three paradigms: (i) direct 26 reasoning via chain-of-thought, (ii) iterat-27 ive search guided by local feedback, and 28 (iii) translation-based solving where LLMs 29 convert textual descriptions of problems 30 into structured solver input. We bench-31 mark general-purpose models (e.g., GPT-32 40 [3]) and reasoning-optimized models 33



Figure 1 Accuracy vs. time cost (measured by # of generated tokens) for different LLM strategies. Translation-based solving yields the best trade-off, while iterative tree search improves over direct reasoning but is more expensive.

³⁵ guage inputs with varying phrasing and

34

(e.g., OpenAI o1 [4]), using natural lan-

 $_{36}$ $\,$ complexity. Our results highlight how solution quality depends on problem scale, amount

³⁷ of distractors in the problem description, and the strategy of using LLMs. We find that

 $_{\rm 38}$ $\,$ the most reliable way to solve NP-hard problems with LLMs is not by having them reason

- $_{39}$ directly or using them as heuristics to guide tree search, but by prompting them to translate
- $_{40}$ problem descriptions into structured formats compatible with classical solvers.

© Haoyu Wang and Dan Roth; licensed under Creative Commons License CC-BY 4.0 42nd Conference on Very Important Topics (CVIT 2016). Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:5 Leibniz International Proceedings in Informatics LIPICS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

 $^{^1\,}$ We plan to release a full version of this paper with complete results and extended discussion.

23:2 LLMs as Translators, Not Thinkers: Structured Output Enables Stronger NP-Hard Problem Solving

41 2 Methods

⁴² We compare three approaches for solving combinatorial problems with LLMs: direct reasoning,

43 translation with solver, and LLM-guided Monte Carlo tree search (MCTS) [1]. These methods

vary in how they leverage the LLM and how they handle problem constraints and search.

⁴⁵ Direct reasoning via chain-of-thought. In this approach, we prompt the LLM to directly
⁴⁶ generate a complete solution based on the natural language description of the problem, using
⁴⁷ a chain-of-thought (CoT) prompt² to encourage step-by-step reasoning.

LLM as Translator and Solver Pipeline. This approach employs the LLM to convert a
natural language problem description into a structured format (e.g., a graph definition, item
list with weights and values) that can be directly consumed by a conventional solver. Once
the structured input is produced, a deterministic algorithm—such as DFS for Hamiltonian
Cycle or dynamic programming for 0-1 Knapsack—is applied to obtain a solution.

LLM-Guided MCTS Framework. We develop a framework that integrates MCTS with 53 LLMs to optimize solutions for combinatorial tasks. The algorithm begins with a feasible 54 initial solution obtained from the direct reasoning method, represented as the root of the 55 search tree. At each expansion step, the LLM is prompted to generate a set of candidate 56 modifications to the current solution. These candidates are assessed by a task-specific 57 evaluator³ which computes the objective value and/or constraint satisfaction, according to 58 different NP-hard problems. Each valid candidate is added to the tree as a child node, and 59 MCTS selects promising branches using either UCT or PUCT scoring. During rollouts, 60 we simulate random paths through the tree to estimate the best achievable reward from 61 a given node. The maximal reward encountered is then backpropagated to update node 62 statistics. In this setup, the LLM provides a learned prior over the action space, while MCTS 63 systematically explores and refines high-quality solutions. 64

3 Experiments and Results

We evaluate the aforementioned methods on four NP-hard problems: Hamiltonian Cycle, Vertex Cover, 0-1 Knapsack, and TSP with time window, with instances varying in size and the presence of distractors—irrelevant or misleading information embedded in the textual description. For graph-based problems, we compare three translation output formats⁴: (1) edge list, (2) lexicographically ordered adjacency list, and (3) adjacency list following the order of appearance in the input. Evaluation metrics include accuracy or approximation score relative to optimal, and efficiency (measured by the number of generated tokens).

As shown in Fig 1, translation-based method with the third format yields the best performance. Iterative methods improve over direct LLM guesses but remain limited by token budgets and rollout depth. Reasoning-optimized models occasionally succeed on small, well-structured instances but degrade with problem size or added noise. We find that using LLMs as translators into structured formats, followed by solvers, is the most reliable and efficient strategy for tackling NP-hard problems described in natural language⁵. Yet the effectiveness of the approach is highly sensitive to the format choice.

 $^{^2}$ Details see Appendix A.

³ We translate the problem description into a structured format as in the second method. The evaluator operates solely on this structured representation to assess objective value and constraint satisfaction.
⁴ Datails and Appendix R

⁴ Details see Appendix B.

 $^{^{5}}$ We find that finetuning does not show significant improvement, which is consistent with findings in [2].

80 — References

 Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In International conference on computers and games, pages 72–83. Springer, 2006.

Lizhou Fan, Wenyue Hua, Lingyao Li, Haoyang Ling, and Yongfeng Zhang. NPHardEval:
 Dynamic benchmark on reasoning ability of large language models via complexity classes. In
 Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
 4092-4114, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL:
 https://aclanthology.org/2024.acl-long.225/, doi:10.18653/v1/2024.acl-long.225.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,
 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-40 system card. arXiv
 preprint arXiv:2410.21276, 2024.

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low,
 Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card.
 arXiv preprint arXiv:2412.16720, 2024.

A Example Prompts for Direct Reasoning

Example Prompt: Hamiltonian Cycle in a Real-World Graph

System: You are a helpful assistant that solves graph problems. Think step-by-step and use <think> tag to enclose your thinking thoughts. Give me an answer with your own reasoning. Refrain from writing code to solve it. Always verify the constraints and provide a valid solution. User: In a metropolitan area, there exist 30 delivery hubs identified as A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P, Q, R, S, T, U, V, W, X, Y, Z, AA, AB, AC, AD. Hub AA processes around 90 packages daily. Station F has direct connections to D, E, G. Hub H is linked directly to A, B, C. ... E connects directly to D, O, P, AB, and AC. Create a route that commences from one station, visits each station precisely once, and ultimately returns to the original station.

23:4 LLMs as Translators, Not Thinkers: Structured Output Enables Stronger NP-Hard Problem Solving

⁹⁷ **B** Prompts for Translation-based Methods

Example Prompt: Translation to Edge List Format

System: You are an expert at converting real-world problem descriptions into structured formats for Python solvers. Your task is to extract graph edges from the description below. Output a single line of JSON with exactly the following keys: "sha", "labels", and "edges". sha" is a made-up 12-character string (e.g., "abc123ef4567") representing the graph ID.

- **"labels"** is the full list of nodes in lexicographic order.
- **"edges"** is a list of two-element lists, each representing an undirected edge as a pair of node labels in lexicographic order.

Output nothing but the JSON.

User: In a metropolitan area, there exist 30 delivery hubs identified as A, B, C, ..., AD. Hub F is connected to D, E, G.

Hub H is linked to A, B, C.

...

E connects directly to D, O, P, AB, and AC.

Create a route that starts and ends at the same station, visiting each exactly once.

98

Example Prompt: Translation to Lexicographic Adjacency List

System: You are an expert at converting real-world problem descriptions into structured formats for Python solvers. Your task is to extract lexicographically ordered adjacency list from the description below. Output a single JSON object with: "sha", "labels", and "adjacency".

- **"sha"** is a made-up 12-character string.
- **"labels"** lists all nodes in lexicographic order.

"adjacency" maps each node (as key) to a list of neighbors, also in lexicographic order. Ensure that all edges are bidirectional, and no duplicates exist.

Output nothing else but the JSON.

User: In a metropolitan area, there exist 30 delivery hubs identified as A, B, C, ..., AD. Hub F is connected to D, E, G.

Hub H is linked to A, B, C.

• • •

qq

E connects directly to D, O, P, AB, and AC.

Create a route that starts and ends at the same station, visiting each exactly once.

Example Prompt: Translation to Lexicographic Adjacency List

System: You are an expert at converting real-world problem descriptions into structured formats for Python solvers. Parse the delivery station connections and return an adjacency list that respects the order in which nodes appear in the description. Output one line of JSON with the following fields: "sha", "labels", "adjacency".

- **"labels**" should list station names in the order they first appear in the description.
- **"adjacency"** is a dictionary that maps each station to its neighbors, preserving the order of mention.

Each connection is undirected. Output strictly one line of JSON, no explanations.

User: In a metropolitan area, there exist 30 delivery hubs identified as A, B, C, ..., AD. Hub F is connected to D, E, G.

Hub H is linked to A, B, C.

. . .

E connects directly to D, O, P, AB, and AC.

Create a route that starts and ends at the same station, visiting each exactly once.

100