



Combinatorial Optimization for All: Using LLMs to Aid Non-Experts in Improving Optimization Algorithms

Camilo Chacón Sartori 

Artificial Intelligence Research Institute (IIIA-CSIC)

Bellaterra, Spain

cchacon@iiaa.csic.es

Corresponding author: cchacon@iiaa.csic.es

Christian Blum 

Artificial Intelligence Research Institute (IIIA-CSIC)

Bellaterra, Spain

christian.blum@iiaa.csic.es

Abstract: We investigate whether Large Language Models (LLMs) can refine the given codebase of an optimization algorithm without requiring specialized user expertise. This is in contrast to works that study optimization algorithm code generation from scratch. To this end, 10 baseline algorithms covering metaheuristics, reinforcement learning, and exact methods are applied to the Traveling Salesman Problem. The results demonstrate that our simple methodology leads to improved algorithm variants in 9 out of the 10 cases analyzed. Notably, the LLMs autonomously incorporated advanced techniques—such as heuristic initializations in exact methods—leading to significant runtime reductions. Furthermore, this performance enhancement did not come at the cost of software quality; the generated code preserved a high maintainability index (averaging 53.40), and for certain models, coincided with simplified structures, observing a reduction in average cyclomatic complexity of up to 19.4%, all without requiring specialized optimization knowledge from the user.

Keywords: Algorithms, Combinatorial Optimization, Large Language Models, Travelling Salesman Problem.

1 Introduction

If we asked you how many optimization algorithms exist, would you be able to come up with an exact answer? Probably not, as there are simply too many algorithms (or algorithm variants) to count. A simple search for ‘optimization algorithm’ in databases like Scopus, IEEE Xplore, or GitHub returns thousands of results. And that is just the start. Classic algorithms, such as the Genetic Algorithm, have spawned so many variations that some barely resemble the original. On top of that, hybrid approaches combine heuristics with exact methods or machine learning techniques. New algorithms are introduced every day.

All these algorithms—whether open-source or proprietary—can be improved using modern technologies. Enhancing their implementations with advanced techniques would benefit not only optimization specialists but also non-expert users. In this paper, we define “non-expert” users as those who have programming skills but lack formal training in optimization theory.

The emergence of Large Language Models (LLMs)—AI systems trained on vast text corpora that can generate code and assist with complex tasks—showcased by innovations like OpenAI’s GPT-O1 [29], Anthropic’s Claude [1], Google’s Gemini [39], Meta’s Llama [40], and DeepSeek R1 [13], has significantly broadened possibilities for scientific innovation. These models have shown code generation capabilities through tools like GitHub Copilot¹ and Cursor AI², enhancing developer efficiency. This raises an intriguing possibility: could individuals without optimization expertise leverage LLMs to enhance existing optimization algorithms?

Current research on LLMs for algorithm design has largely followed two paths: evolving novel heuristics or specific components within complex frameworks [8, 20, 33, 38, 45], and discovering specific heuristic operators by leveraging the LLM’s semantic understanding of the code to identify parameters that are unused in one part but could be repurposed or utilized in another [35]. However, less attention has been given to the practical challenge of improving *complete existing algorithm implementations as a whole*, without requiring theoretical knowledge of combinatorial optimization.

Unlike previous approaches that focus on evolving isolated operators, our work proposes a holistic methodology where the LLM refines the *entire algorithm codebase*. This allows for structural optimizations and the integration of global heuristics which would not be possible when limited to specific algorithm components.

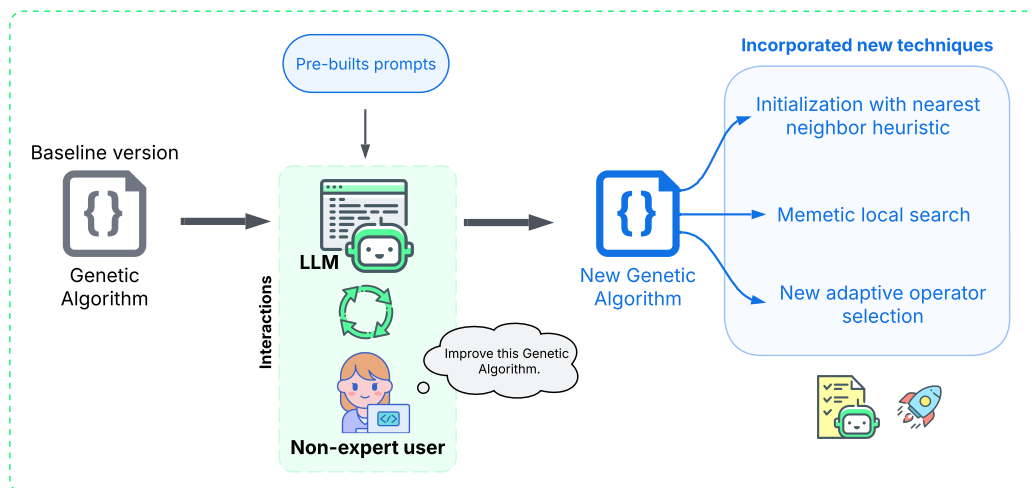


Figure 1: A non-expert user’s interaction with an LLM can enhance an existing genetic algorithm by incorporating modern techniques.

To bridge this gap, our work presents the first large-scale, systematic evaluation of the ability of LLMs to upgrade classical algorithms across a wide range of families—including metaheuristics, reinforcement learning, deterministic heuristics, and exact methods. Our study focuses on the canonical Travelling Salesman Problem (TSP), where we apply a simple and reproducible prompting strategy through a chatbot interface. While our experiments are centered on the TSP, its foundational nature as a problem of optimal sequencing makes it a highly relevant case study. The core challenge of finding the best permutation of nodes is shared by many other classic optimization problems, such as the Vehicle Routing Problem, the Hamiltonian Cycle Problem, and the Sequential Ordering Problem. Therefore, the strategies and findings from this work are likely applicable to any problem centered on the fundamental question: “In what order should a series of tasks be performed to minimize total cost?”

Our findings demonstrate that LLMs can propose and implement sophisticated improvements, such as incorporating modern heuristic components or engineering more efficient data structures, often resulting in a reduction in code complexity compared to the original implementations. Crucially, these enhancements are applied not to isolated code fragments (heuristics), as in previous work, but to the entirety of an

¹<https://github.com/features/copilot>

²<https://www.cursor.com/>

optimization algorithm’s codebase. Furthermore, we analyze cases where the improvements were not substantial, offering insights into the current limitations of this approach. Ultimately, this work showcases a practical methodology that lowers the barrier for practitioners without deep optimization expertise to access higher-performance algorithms, leveraging the vast knowledge embedded in modern LLMs. A schematic of our framework using a Genetic Algorithm is shown in Figure 1.

The paper unfolds as follows. In Section 2, we provide an overview of LLM advancements in combinatorial optimization, introduce the TSP problem, and briefly describe the 10 selected algorithms. Our methodology for enhancing existing optimization algorithms is detailed in Section 3. An exhaustive analysis of the 10 optimization algorithms improved through LLM application is presented in Section 4. Section 5 discusses implications and identifies directions for future research. We conclude by highlighting the key findings from our investigation.

2 Background

2.1 Large Language Models in Combinatorial Optimization

Large Language Models (LLMs) have recently shown promise in optimization tasks [15, 21, 44] by exploiting the vast knowledge gained during their pre-training phase. Beyond guiding the optimization process, LLMs excel at detecting patterns, identifying key features in problem instances, and refining search spaces. They have also shown to be able to generate new heuristics tailored to specific problems. Furthermore, LLMs offer valuable insights by explaining the results of optimization problems, making them versatile tools for both solving and interpreting complex tasks.

Recent research has convincingly demonstrated that LLMs can be powerful tools for creating and enhancing complex optimization algorithms [20, 27, 33, 35, 38, 45]. To situate our current work within this rapidly evolving field and clearly delineate its unique contributions, we provide a comprehensive comparative analysis of these state-of-the-art approaches in Table 1. Building upon our own initial proof-of-concept [35], this work significantly expands the investigation by conducting a large-scale, systematic study across 10 classical algorithms of four different types. Our findings confirm that the principle of LLM-driven improvement does not only work occasionally, but is a broadly applicable technique in the case of algorithms for the Travelling Salesman Problem.

The rapid advancements in using LLMs as black-box collaborators for optimization demand a clear positioning of new methodologies. Table 1 delineates our work’s unique contributions in this crowded landscape. Specifically, our approach is distinguished by:

- **Systematic Scope over Anecdotal Evidence.** Where prior work often provides proofs-of-concept on a limited set of algorithms, we deliver a systematic benchmark across ten algorithms spanning four distinct families. We prove that a straightforward, reproducible prompting strategy is robust enough to yield significant performance gains across this diverse set without compromising code stability.
- **Practitioner-Centric Design over Expert-Centric Frameworks.** Our methodology marks a fundamental departure from expert-centric systems that require designing complex evolutionary loops or prompt engineering strategies. It is uniquely designed for the practitioner—a user with programming skills but lacking deep theoretical expertise. Our work is the first to demonstrate a path for these users to upgrade complete, monolithic codebases, rather than just optimizing isolated functions. This holistic approach significantly lowers the barrier for applying state-of-the-art optimization techniques.

In this article, we focus on a specific type of optimization problem: combinatorial problems. These problems have unique properties that set them apart: many valid solutions (including the existence of equivalent or similar solutions), decomposability (some problems can be broken down into smaller, more manageable subproblems), constraint handling (a set of rules that define valid solutions), and search space structure (the presence of multiple local optima, requiring well-chosen search strategies to avoid getting trapped in suboptimal solutions), among others.

Table 1: Expanded Comparative Analysis of LLM-based Algorithm Design Approaches.

Feature	This Study	Improving Existing Optimization Algorithms [35]	AlphaEvolve [27]	LLaMEA/ReEvo [38, 45]	FunSearch/Evo-Heuristics [20, 33]
Main Goal	To provide a large-scale benchmark on enhancing <i>complete</i> implementations of <i>existing</i> algorithms from diverse families.	To demonstrate a proof of concept through the improvement of a <i>single</i> , complex heuristic in a state-of-the-art algorithm.	To autonomously evolve a population of algorithms starting from a seed program, leveraging LLM-based agents to discover improved variants.	To generate new metaheuristics or heuristic components using an evolutionary computation loop.	To discover specific mathematical functions or heuristics from scratch.
LLM Input	A complete, functional code file from a well-known framework [30].	The complete code of a single, high-performance algorithm.	An initial “seed” algorithm and a fitness function.	A set of prompts representing algorithm components (e.g., operators).	A problem description and a code skeleton for a specific function.
Improvement Process	Interactive LLM-driven prompting and refinement of full algorithms.	Interactive LLM prompting to refine a specific heuristic or function.	Evolutionary loop where the LLM acts as an advanced mutation operator to generate new programs.	Evolutionary loop where the LLM acts as a crossover/mutation operator on algorithm components.	Program search tree where the LLM proposes new function implementations.
User’s Role	Practitioner/Non-expert: Provides base code, validates final solution.	Expert: Seeks to enhance a specific, already-strong algorithm.	Expert: Designs the evolutionary process and fitness function.	Expert: Designs the evolutionary framework and component prompts.	Expert: Defines the problem, evaluator, and code structure.
Key Contribution	A systematic benchmark demonstrating that a single, prompt-based methodology can holistically improve algorithms from diverse families without specialized theoretical knowledge.	A case study demonstrating the feasibility of LLM-based enhancement on a complex, expert-level algorithm.	An evolutionary framework for automated program search and algorithm improvement.	A framework for the automatic generation of algorithms.	A method for the automatic discovery of functions/heuristics.

Table 2: Overview of Selected Algorithms for Solving the Travelling Salesman Problem (TSP).

Algorithm	Characteristics	Application to TSP
Metaheuristic		
Ant Colony Optimization (ACO) [10]	Probabilistic, pheromone-based learning	Simulates ants' foraging behavior where solutions (routes) are constructed based on pheromone trails left by previous solutions.
Genetic Algorithm (GA) [14, 31]	Population-based, crossover, mutation	Generates a population of routes and evolves them through selection, crossover, and mutation to find near-optimal solutions.
Adaptive Large Neighborhood Search (ALNS) [34]	Adaptive destruction and reconstruction	Iteratively destroys and reconstructs routes, dynamically adjusting strategies based on previous performance.
Tabu Search (TABU) [12]	Use of memory structures (tabu lists)	Iteratively modifies routes while keeping a list of features of previously visited solutions to prevent revisits.
Simulated Annealing (SA) [17]	Probabilistic, temperature-based search	Iteratively refines a route by accepting worse solutions with a decreasing probability to escape local optima.
Reinforcement Learning		
Q_Learning [41]	Value-based learning, exploration-exploitation trade-off	Learns an optimal routing policy by iteratively updating action-value functions based on rewards from different paths.
SARSA [41]	On-policy learning, continuous updates	Uses an on-policy approach to learning optimal routing strategies based on real-time interactions with the environment.
Deterministic Heuristic		
Christofides [7]	Guarantees 1.5-optimality, MST-based	Constructs a minimum spanning tree, finds perfect matching, and combines them to form a tour.
Convex Hull [11]	Geometric approach	Starts with a convex hull and incrementally inserts remaining points in a way that minimizes travel distance.
Exact		
Branch and Bound (BB) [26]	Systematic enumeration, pruning	Explores all possible solutions while pruning suboptimal paths to guarantee optimality.

As a result, researchers in this field must not only be knowledgeable about combinatorial optimization problems but also highly proficient in implementing optimization algorithms. Given the importance of computational efficiency, factors such as memory optimization, effective data structure management, minimizing unnecessary abstractions, and carefully selecting the right programming language play a crucial role in the design and development of these algorithms.

To introduce this novel research direction, the next two subsections present a classic combinatorial optimization problem along with ten traditional optimization algorithms, grouped into four distinct categories. Then, in the methodology section, Section 3, we demonstrate how LLMs can be leveraged to enhance these ten algorithms, improving both performance and efficiency while streamlining their implementations.

2.2 The Travelling Salesman Problem: A Brief Description

The Travelling Salesman Problem (TSP) is one of the most iconic problems in combinatorial optimization, serving as a foundational pillar in both Artificial Intelligence and Operations Research. Formally, let $\mathcal{C} = \{c_1, c_2, \dots, c_n\}$ be a set of n cities, and let $D : \mathcal{C} \times \mathcal{C} \rightarrow \mathbb{R}_{\geq 0}$ be a function that assigns a non-negative distance between each pair of cities. The objective of the TSP is to find a permutation π of the indices

$\{1, 2, \dots, n\}$ that minimizes the total travel distance of a closed tour, formally defined as:

$$\min_{\pi \in S_n} \left\{ \sum_{i=1}^{n-1} D(c_{\pi(i)}, c_{\pi(i+1)}) + D(c_{\pi(n)}, c_{\pi(1)}) \right\},$$

where S_n denotes the set of all permutations of $\{1, 2, \dots, n\}$.

For this study, we selected the TSP because of the vast amount of implementations available online—public code repositories, books, and scientific articles—indicating that LLMs likely possess extensive knowledge of techniques for solving it [24].

2.3 Traditional Optimization Algorithms for the TSP

Algorithms for solving the TSP come in a broad variety of forms and approaches. We have chosen ten distinct optimization algorithms, summarized in Table 2. These algorithms are categorized into four groups: metaheuristic methods (stochastic, iterative optimization algorithms), reinforcement learning (policy-driven), deterministic heuristics, and an exact algorithm. These four algorithm categories can be briefly characterized as follows:

1. **Metaheuristics** [4] are versatile optimization methods that use heuristic and stochastic principles to search solution spaces. While not guaranteeing global optima, they efficiently find high-quality solutions.
2. **Reinforcement Learning** [43] is a machine learning paradigm where an agent learns optimal decision-making by interacting with an environment and maximizing cumulative rewards.
3. **Deterministic heuristics** [25] are among the most basic algorithms for combinatorial optimization. They generate generally one solution from scratch by making a myopic, deterministic decision at each step.
4. **Exact algorithms** [25] ensure optimality by exhaustively exploring the solution space.

The choice of algorithms from these four categories guarantees that our LLM-based improvement framework is tested across a diverse spectrum of methodologies, as the selected algorithms—though addressing the same problem (TSP)—vary fundamentally in their underlying principles.

2.4 Selected Implementations

To minimize the possibility of implementation errors in the ten algorithms for solving the TSP, and to ensure that these implementations are being utilized by the community, we employed `pyCombinatorial` [30].³ This Python library, which includes a whole range of optimization algorithms for solving the TSP, has received over 100 stars on GitHub and was created by Valdecy Pereira. Each implementation is based on a standard algorithm variant, providing us with an excellent testing environment for attempting to improve them using LLMs. In the next section, we detail our methodology for improving optimization algorithms.

3 Methodology

3.1 Enhancing Traditional Optimization Algorithms with Large Language Models

We introduce a methodology based on LLM interactions to generate enhanced versions of the 10 before-mentioned algorithms, building upon and extending the approach proposed by [35] (see Figure 1). This process can be replicated using the chatbot available at the project URL.⁴

Given the initial set $\{A_i \mid i = 1, \dots, 10\}$ of 10 original algorithm codes and an LLM, the algorithm improvement process can technically be stated as follows. First, a prompt P_i is generated based on our

³<https://github.com/Valdecy/pyCombinatorial>

⁴Here will be the url.

general prompt template T (shown in the Prompt Template box), the algorithm name N_i , the signature of the main function S_i , and the algorithm code A_i :

$$P_i = \text{PRODUCE_PROMPT}(T, N_i, S_i, A_i), \quad i = 1, \dots, 10 \quad (1)$$

Then, the generated prompt P_i is executed by the LLM given a set of hyperparameters θ , resulting in a changed/updated implementation A'_i :

$$A'_i = \text{LLM}_{\text{execute}}(P_i, \theta), \quad i = 1, \dots, 10$$

To ensure correctness, each A'_i undergoes a validation process (see below). If an output fails validation, the refinement process is repeated iteratively until a valid version is obtained.

3.1.1 Code Validation

The validation of A'_i may fail for two reasons:

1. **Execution errors:** These lead to immediate failures during code runtime.
2. **Logical inconsistencies:** The algorithm executes without errors but produces invalid TSP solutions.

In the first case, an error message e is generated, and the LLM refines the code based on this feedback:

$$A'_i = \text{LLM}_{\text{execute}}(A'_i, \theta, e)$$

For the second case, where execution is error-free but the generated solutions are invalid, an explicit prompt requesting a correction is passed to the LLM, such as:

“The provided code generates invalid solutions; please verify and return a corrected version.”

The refinement loop proceeds until a valid code A'_i is obtained. This process is carried out through an interactive conversation with an LLM-based chatbot. To increase the chances of generating a valid code with each retry, we begin with a high-temperature setting, which is then progressively lowered in each iteration of the process.

In LLMs, *temperature* controls the randomness of the model’s output. A higher temperature (e.g., 1.0 or 2.0) makes the model’s responses more diverse and less predictable, while a lower temperature (e.g., 0.2) makes the output more deterministic and stable.

Table 3 shows the results of this procedure for the five selected LLMs.⁵ A green checkmark (✓) indicates that the first obtained A'_i was valid, while a red cross (✗) signifies that corrections were necessary due to either of the two code failures. A double red cross (✘) indicates that both failures occurred. Moreover, in the case of corrections, the number of necessary corrections is provided in a second column. Note that Table 3, below the LLM names, also indicates the initial temperature setting.

3.2 Prompt Design: A Focus on Simplicity and Accessibility

A central hypothesis of this study is that significant algorithmic improvements can be achieved without requiring users to possess deep expertise in prompt engineering. To test this, we developed a single, standardized prompt template. The design of this prompt is *intentionally basic*.

Our objective is not to engage in an exhaustive search for the “optimal” algorithm variant through elaborate prompting, which constitutes a separate research direction. Instead, our goal is to establish a reproducible baseline for what LLMs can achieve when prompted by a non-expert user. This approach directly aligns with our focus on accessibility and allows us to isolate the LLM’s intrinsic ability to enhance code.

Despite its simplicity, the prompt’s formulation provides clear, high-level directives. It focuses the LLM on two critical performance axes—(1) improving solution quality and (2) accelerating convergence—and explicitly encourages the integration of state-of-the-art techniques to achieve these goals. The template is shown below:

⁵The reasoning behind selecting these five LLMs (rather than others) will be explained in Section 4.

Table 3: Success rates of code generation across five LLMs. The table indicates whether the first attempt yielded valid code (✓) or required corrections (✗), along with the number of attempts needed. ‘✗✗’ indicates failure in both execution and logic.

Algorithm	Claude-3.5-Sonnet (temp = 1.0)		Gemini-exp-1206 (temp = 2.0)		Llama-3.3-70B (temp = 1.0)		GPT-O1 (temp = 1.0)		DeepSeek-R1 (temp = 1.0)	
	Success		Success		Success		Success		Success	
	1st Try	# Attempts	1st Try	# Attempts	1st Try	# Attempts	1st Try	# Attempts	1st Try	# Attempts
ACO	✗	1	✓	-	✓	-	✓	-	✓	-
GA	✓	-	✗	1	✓	-	✗✗	3	✓	-
ALNS	✗	1	✓	-	✓	-	✓	-	✗	3
TABU	✓	-	✓	-	✓	-	✓	-	✓	-
SA	✓	-	✓	-	✓	-	✓	-	✓	-
Q_Learning	✓	-	✓	-	✓	-	✓	-	✓	-
SARSA	✓	-	✓	-	✓	-	✓	-	✓	-
Christofides	✓	-	✓	-	✓	-	✓	-	✓	-
Convex Hull	✗	1	✓	-	✓	-	✓	-	✓	-
Branch and Bound	✓	-	✗✗	4	✓	-	✓	-	✓	-

Prompt Template

You are an optimization algorithm expert.

I need to improve this {{algorithm name}} implementation for the travelling salesman problem (TSP) by incorporating state-of-the-art techniques. Focus on:

1. Finding better quality solutions
2. Faster convergence time

Requirements:

- Keep the main function signature: {{the signature of an the main function}}
- Include detailed docstrings explaining:
 - * What improvement is implemented
 - * How it enhances performance
 - * Which state-of-the-art technique it is based on
- All explanations must be within docstrings, no additional text
- Check that there are no errors in the code

IMPORTANT:

- Return ONLY Python code
- Any explanation or discussion must be inside docstrings
- At the end, include a comment block listing unmodified functions from the original code

Current implementation:

{{algorithm code}}

The prompt template begins with “*You are an optimization algorithm expert,*” a technique known as “role prompting”, which has been empirically shown to guide LLMs toward a specific behavior or specialization [3, 18]. By setting a clear context from the outset, it enhances both the relevance and quality of the model’s response.

This approach aligns with “in-context learning” [9, 19, 36], combining an external context (the user-provided code in the prompt) with an internal one (the LLM’s knowledge of various techniques to improve the TSP algorithm).

An important insight: supplying a complete algorithm code within the prompt works like a ‘map,’ guiding the model on how to update the code. The LLM must preserve the overall structure, making modifications only in the relevant sections without breaking the logic. Without this external context (the provided algorithm), the model’s solution would likely be more constrained and less effective, as it would have to generate everything from scratch. In contrast, with an initial codebase, the model can focus on refining and improving specific areas rather than rebuilding the entire algorithm code from scratch. In other words, the provided code *influences* the update proposed by the LLM [9, 35].

The effectiveness of this simple prompting strategy stems from the LLMs’ extensive exposure to standard optimization code during pre-training. Since the selected baseline algorithms are classic and widely represented in open-source repositories, the models possess a strong “prior” or latent knowledge regarding their implementation and potential improvements. Consequently, the prompt does not need to teach the model optimization theory; rather, it acts as a catalyst to retrieve and adapt this pre-existing knowledge to the specific context of the provided code skeleton.

As technically indicated already in Eq. 1, the prompt template receives three dynamic variables that are placed in the corresponding positions in the prompt template (enclosed within `{ { ... } }`):

- **Algorithm’s name:** steers the LLM toward a specific context.
- **Main function’s signature:** ensures the initial function’s input arguments, output values, and name remain unchanged, preventing unintended modifications that could affect compatibility with the original code.
- **Algorithm code:** the optimization algorithm’s original implementation in Python.

Additionally, we explicitly instruct the LLM to report the modifications it makes (*“Include detailed docstrings explaining: ...”*). This step is essential, as it enables us—as shown in Section 4—to understand why a particular LLM’ code outperforms the original or another model’s code.

4 Experimental evaluation

In this section, we present experiments with the 10 previously mentioned optimization algorithm codes taken from `pyCombinatorial`. We describe our setup for utilizing LLMs, the parameter tuning of the stochastic algorithms, the TSP datasets and evaluation metrics employed, and the comparative analysis of the results. We highlight key details from the generation process and conclude with an analysis concerning code complexity.

4.1 Setup

4.1.1 LLMs Environment

We selected five leading code-generation LLMs: Anthropic Claude-3.5-Sonnet [1], Google Gemini-exp-1206 [39], Meta Llama-3.3-70b [40], OpenAI GPT-O1 [29], and DeepSeek-R1 [13]. For simplicity, we will refer to the models as Claude, Gemini, Llama, O1, and R1 throughout the remainder of this paper. Note that these LLMs rank among the top models in the LiveBench benchmark [42], which is immune to both test set contamination and the biases of LLM-based and human crowdsourced evaluations (**as of February 2025**). Using the OpenRouter API, we executed identical prompts across all models, enabling straightforward model switching for transparent experimentation. This produced 50 new algorithm codes which, combined with the 10 original algorithm codes from the `pyCombinatorial` framework, gave us 60 Python files ready for evaluation.

4.1.2 Hardware Environment

All experiments, including parameter tuning, are conducted on a cluster equipped with Intel® Xeon® CPU 5670 processors (12 cores at 2.933 GHz) and 32 GB of RAM.

Table 4: Parameter values obtained by tuning with `irace`. Ranges show minimum/maximum values considered for tuning.

Algorithm	Parameter	Range	Original	Claude-3.5-Sonnet	Gemini-exp-1206	Llama-3.3-70B	GPT-O1	DeepSeek-R1
Metaheuristics								
ACO	m (ants)	(2, 20)	7	4	2	3	17	20
	α (alpha)	(1.0, 2.0)	1.34	1.72	1.67	1.46	1.72	1.22
	β (beta)	(1.0, 2.0)	1.59	1.24	1.98	1.97	1.93	1.55
	ρ (decay)	(0.01, 0.3)	0.24	0.12	0.24	0.29	0.06	0.05
GA	N (population size)	(5, 100)	97	14	97	84	55	58
	μ (mutation rate)	(0.01, 0.2)	0.02	0.04	0.16	0.16	0.01	0.13
	e (elite)	(1, 5)	4	2	5	2	3	5
ALNS	λ (removal fraction)	(0.05, 0.3)	0.27	0.05	0.26	0.29	0.22	0.29
	ρ (rho)	(0.01, 0.3)	0.27	0.25	0.04	0.2	0.27	0.02
TABU	T (tabu tenure)	(3, 30)	8	12	30	15	10	9
SA	T_0 (initial temperature)	(1, 50)	12	49	9	30	35	50
	T_f (final temperature)	(0.0001, 0.1)	0.0547	0.0464	0.074	0.056	0.0433	0.048
	α (cooling rate)	(0.8, 0.99)	0.9895	0.8732	0.8956	0.8131	0.9154	0.8777
Reinforcement Learning								
RL_QL	lr (learning rate)	(0.01, 0.5)	0.44	0.15	0.26	0.49	0.46	0.34
	df (decay factor)	(0.8, 0.99)	0.97	0.82	0.98	0.87	0.98	0.82
	ϵ (epsilon)	(0.01, 0.3)	0.09	0.28	0.03	0.24	0.21	0.13
	E (episodes)	(1000, 10000)	4266	1082	4906	2474	1294	1989
SARSA	lr (learning rate)	(0.01, 0.5)	0.04	0.36	0.49	0.19	0.41	0.29
	df (decay factor)	(0.8, 0.99)	0.86	0.91	0.80	0.88	0.83	0.87
	ϵ (epsilon)	(0.01, 0.3)	0.23	0.18	0.16	0.08	0.12	0.16
	E (episodes)	(100, 5000)	105	156	1850	137	124	1711

4.1.3 Parameter Tuning

While the deterministic heuristics and the branch and bound method are parameter-less, the seven probabilistic approaches (five metaheuristics and two reinforcement learning algorithms) require careful parameter tuning to perform well. Consequently, we tuned all 42 stochastic algorithm codes: the 7 original versions and the 35 new variants generated by the LLMs. To ensure a fair and robust comparison, we employed `irace` [23], a well-established automatic algorithm configuration tool. The tuning process was executed as parallel jobs on the SLURM cluster described in Section 4.1.2. For each algorithm execution during the tuning phase, the CPU time limit was set to the number of cities in the instance (in seconds). Table 4 details the parameter ranges considered for tuning each algorithm variant, as well as the final best configurations selected by `irace`.

Parameter tuning in stochastic algorithms with parameters is essential, as suboptimal configurations can lead to poor performance regardless of the algorithm’s inherent quality. For instance, in ACO, we tune key parameters— m (number of ants), α , β , and ρ (pheromone decay)—for all six code variants (five LLM-generated ones plus the original) to ensure they operate under optimal conditions for the TSP problem. This guarantees a fair evaluation, as each code variant is assessed using its best possible configuration.

It is important to note that the crossover operator for the Genetic Algorithm was not subject to automatic tuning due to the specific design of the original implementation. The pyCombinatorial library employs a hardcoded hybrid crossover strategy that stochastically alternates between Best Cost Route (BCR) Crossover and Edge Recombination (ER) Crossover with equal probability (0.5) at each step. Since this mechanism does not come with a tuneable parameter in the baseline code, we maintained this hard-coded strategy to ensure a faithful comparison with the original library. Similarly, regarding the Reinforcement Learning algorithms, the range of episodes (E) differs between Q-Learning and SARSA. This decision reflects their distinct convergence characteristics observed during preliminary testing—where SARSA typically required fewer episodes to stabilize—and aligns with the default parameterization logic of the original pyCombinatorial implementation.

4.2 Benchmark Datasets and Evaluation Metrics

To evaluate all algorithm codes, we use problem instances from the well-known TSPLib library [32]. We select 10 instances from the available ones, ranging from a small instance with 99 cities to a large one

with 1084 cities.⁶ This selection ensures a comparison across a diverse set of problem instance sizes.

As an evaluation metric, we used the objective function value of the best-found solution in all cases except for the Branch and Bound (BB) codes. This is because BB is an exact algorithm that, if given enough computation time, will always find an optimal solution. Therefore, we use runtime as the evaluation metric in the case of BB. Moreover, as the runtime of BB for the 10 selected problem instances is very high, we instead generate 10 random TSP instances with 10 to 15 cities for the evaluation of the BB codes. Interestingly, as we will see in the comparative analysis subsection, some LLM-generated versions of BB incorporate heuristic mechanisms during algorithm initialization, leading to significant improvements in runtime performance.

4.3 Experimental Design

The experiments were designed as follows:

- **Stochastic Algorithms.** Each of the metaheuristics and reinforcement learning codes is applied 30 times independently to each of the 10 problem instances. The output of each run is the best solution found. Performing 30 independent algorithm executions for each problem instance is a common practice in the optimization community to obtain a reliable estimate of the algorithm’s performance. Moreover, the CPU time limit for each algorithm execution is set to the number of cities (in seconds) of the tackled problem instance. For example, the run-time limit for RAT99 is 99 seconds. This method, which aligns execution time with the instance size, is a common practice for comparing algorithms that solve the TSP.
- **Deterministic Heuristics.** Christofides and Convex Hull, since they are deterministic heuristics, always yield the same result for a given problem instance. Therefore, all corresponding codes are executed exactly once per instance.
- **Branch and Bound.** As previously mentioned, since Branch and Bound is an exact algorithm, the focus is on runtime rather than solution quality. All Branch and Bound codes are applied 30 times to each of the small problem instances specifically generated for the evaluation of the Branch and Bound codes.

4.4 Comparative Analysis with Original Algorithm Codes

The comparative analysis between the LLM-updated algorithm codes and the original algorithm codes (referred to as ‘original’ from now on) is studied in the following in a separate way depending on the type of optimization algorithm.

4.4.1 Metaheuristics

The results of all metaheuristic codes are shown by means of boxplots in Figure 2. Note that the y-axes are shown in a logarithmic scale.

1. **ACO:** Apart from the first problem instance (RAT99) where the LLM-generated codes of Gemini, O1, and R1 perform similarly to the original code, in all other problem instances the LLM-generated codes of the mentioned three LLMs outperform the original code with statistical significance. It also appears that the LLM-generated codes of O1 and R1 are somewhat more robust than the original code, which can be seen in smaller boxes. In contrast, the code generated by Claude, apart from the first problem instance, performs always worse than the original code. Finally, the Llama-generated code generally performs similarly to the original code, with the exception of the first problem instance.
2. **GA:** The original code exhibits very low robustness for the first two problem instances, which can be seen by the large boxes. Generally, most codes (except for R1) are quickly trapped in local

⁶TSPLib names of the selected TSP instances: RAT99, BIER127, D198, A280, F417, ALI535, GR666, U724, PR1002, and VM1084.

optima which they cannot escape. The R1-generated code clearly outperforms all others, including the original.

3. **ALNS**: Generally, the best-performing codes are those by Claude and Gemini, with a slight advantage for Gemini in the last two instances. Another noteworthy aspect is the low robustness of the O1-generated code in this case.
4. **TABU**: Like in the case of GA, also the TABU code generated by R1 significantly outperforms the remaining codes. Only the O1-generated code can compete for the smallest two problem instances. This suggests that R1 excels at generating efficient optimization algorithms.
5. **SA**: The Claude-generated codes clearly show the weakest performance here. In contrast, the R1-generated code again outperforms the remaining ones.

4.4.2 Reinforcement Learning (RL)

In Figure 3, the RL codes generated by the LLMs are compared with the original RL codes. The displayed results differ from the metaheuristics case presented before in two key aspects. First, some LLM-generated codes fail to produce a result for all problem instances within the time limit. These cases are marked as N/A. Second, especially in the case of Q_Learning the original code performs more competitively. We analyze these aspects below:

6. **Q_learning**: The original code is actually the best-performing one in this case. The Llama-generated code is the only one that achieves a nearly comparable performance—an unexpected result given Llama’s poor performance in the case of the metaheuristics. In addition, the Claude-generated code outperforms the original one on the RAT99 and A280 instances.
7. **SARSA**: The general picture here aligns more with that observed in the case of the metaheuristics. The original code struggles (in comparison to some LLM-generated codes) as instance sizes grow larger. The R1-generated code fails to produce results for the largest three instances. The only code maintaining a stable and strong performance across all 10 instances is the O1-generated one.

4.4.3 Deterministic/Heuristic

Remember that, as the chosen heuristics are deterministic, there is no need to analyze the distribution of their results over multiple runs. Therefore, in Figure 4, we simply compare the GAP of the results produced by the LLM-generated codes (in percent) relative to the results of the original codes. A positive value indicates that the respective LLM-generated code outperforms the original, while a negative value suggests the opposite.

8. **Christofides** (see Figure 4 (a)): While the Llama-generated code produces very similar results to the original code over the whole instance range, the Gemini-generated code is (apart from instance D198) always inferior. Moreover, its relative performance decreases as instance size grows. Concerning O1 and R1, it can be stated that the performance of their codes is slightly inferior to the one of the original code for rather small problem instances. However, with growing instance size, they clearly outperform the original code.
9. **Convex Hull** (see Figure 4 (b)): In contrast to Christofides, the Convex Hull codes generated by O1 and R1 perform rather poorly. In fact, the best code for Convex Hull is the one generated by Gemini. This code has slight disadvantages for smaller instances but increasingly outperforms the original code with growing instance size. The code generated by Claude shows the opposite pattern. While it outperforms the original code for smaller problem instances, its performance strongly decreases with growing instance size.

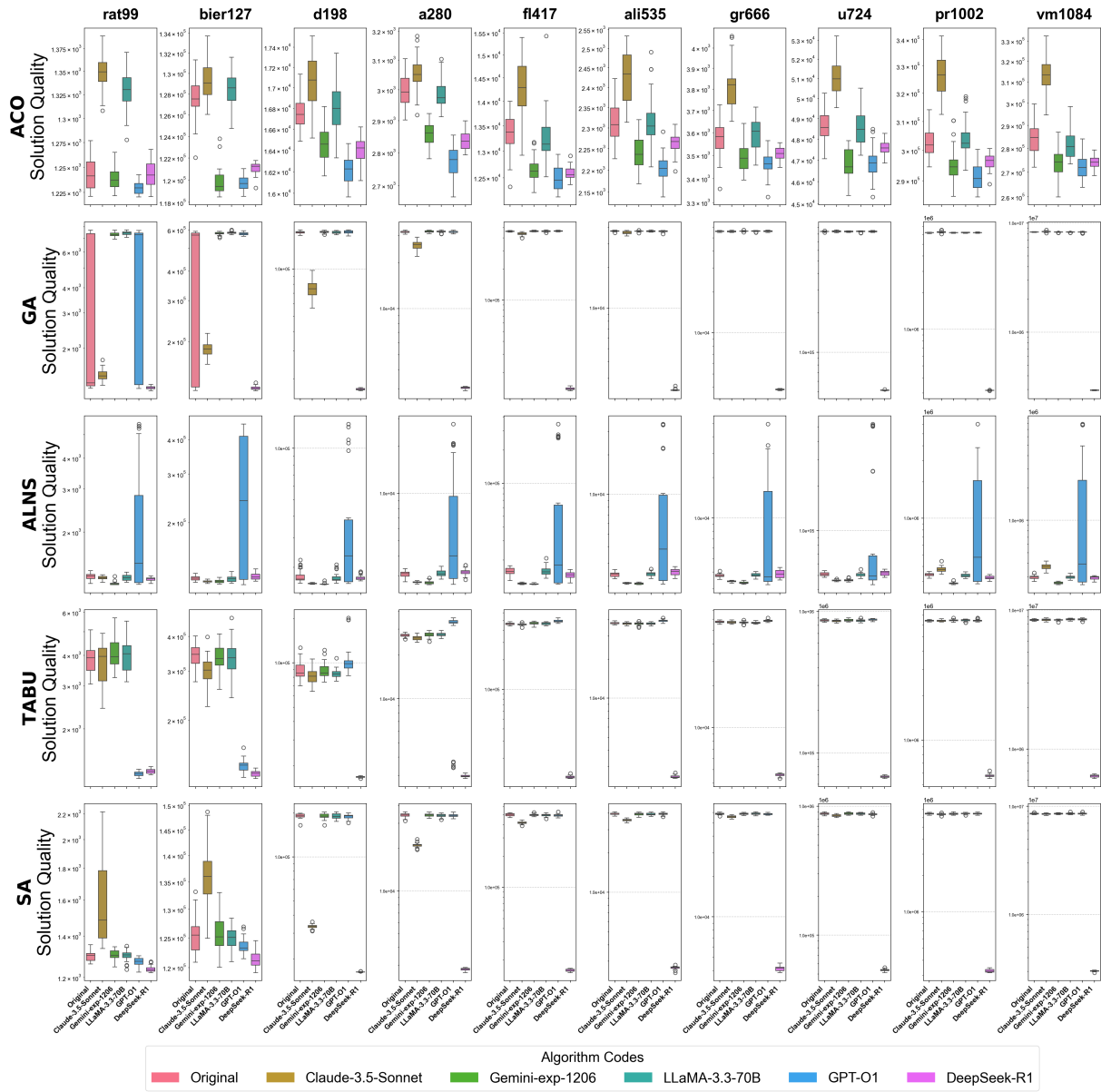


Figure 2: Performance comparison between LLM-generated metaheuristic codes and original implementations. The boxplots illustrate the distribution of the best solution values obtained across 30 independent runs for each TSP instance. Since the TSP is a minimization problem, lower values indicate better performance. Note that the y-axes employ a logarithmic scale to accommodate the significant variance in solution quality among the algorithms.

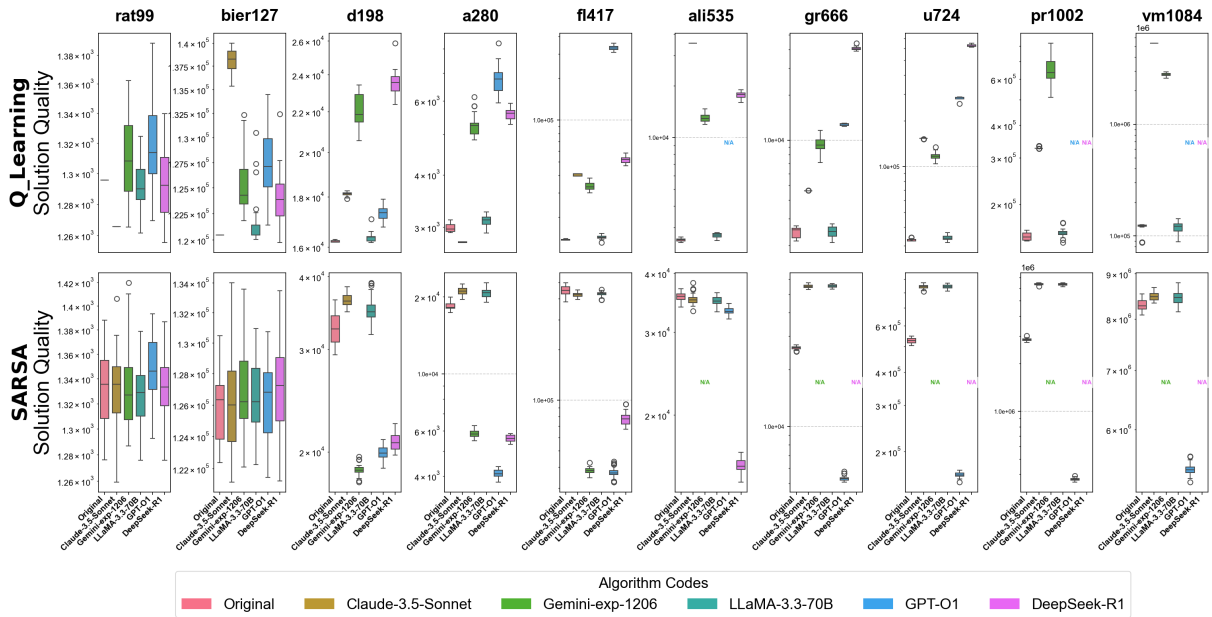


Figure 3: Performance comparison between LLM-generated Reinforcement Learning (RL) codes and original implementations. The boxplots illustrate the distribution of the best solution values obtained across 30 independent runs. Since the TSP is a minimization problem, lower values indicate better performance. Note that the y-axes employ a logarithmic scale.

4.4.4 Exact Approach

As already mentioned in Section 4.3 (Experimental Design), in the context of the exact BB method, the comparison is based on computation time.

- Branch and Bound:** Figure 5 shows that the codes generated by O1 and R1 outperform both the original code. Notably, R1—an open-weight LLM—achieves the best performance, surpassing all proprietary models. In contrast to O1 and R1, the other LLM-generated codes perform worse than the original code.

Summary: A notable conclusion is that LLMs can produce improved versions of baseline algorithms, resulting in performance improvements without necessitating specialized expertise in each algorithm. In the following subsection, we showcase examples of code improvements achieved, for example, by integrating more sophisticated algorithmic components.

4.5 Key Insights in Code Generation

Next, we explore why certain LLM-generated (or LLM-updated) codes outperform the original ones. Our focus was to understand if this was due to optimized data structures, for example, or due to adding different algorithmic components. In particular, we analyze four cases to address these questions on the basis of the LLM-generated codes.

4.5.1 Case 1: GA (R1-generated code)

The R1-generated version of GA features the following improvements, as stated by the model itself by means of a docstring in the code, as requested in the prompt.

Improvements:
 1. Hybrid initialization with nearest neighbor heuristic

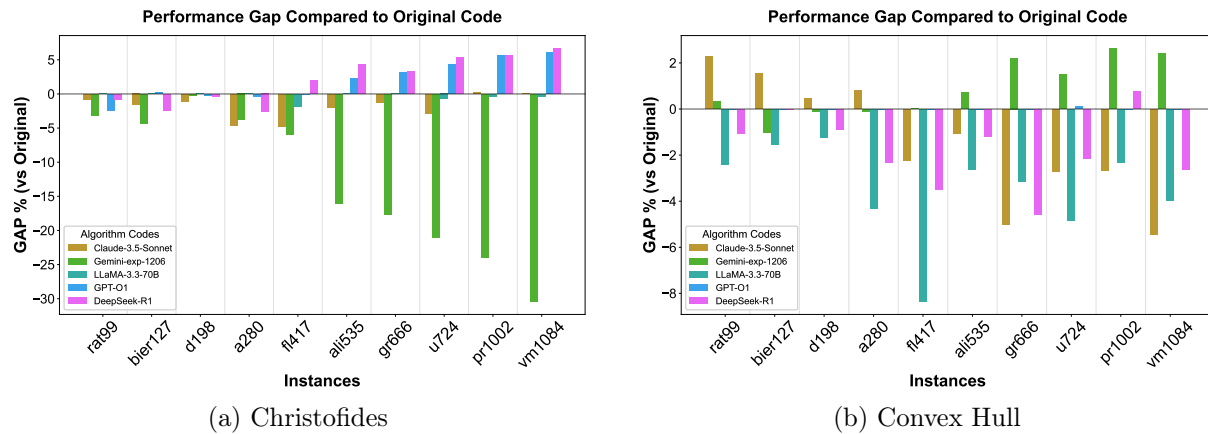


Figure 4: Comparison of the deterministic heuristic codes generated by the five LLMs with the original codes. The bar plots show the performance gaps (in percent) relative to the original codes. Note that a positive value indicates that the LLM-generated code produces a better solution.

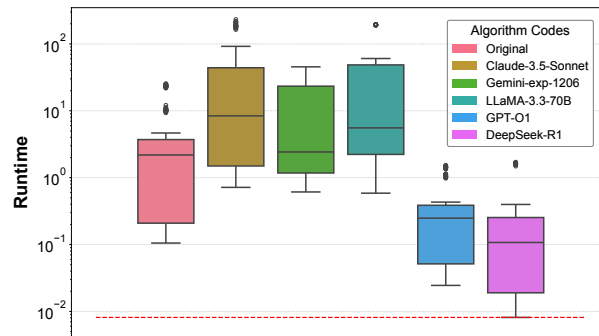


Figure 5: Comparison of the BB codes generated by the five chosen LLMs with the original BB code (in terms of computation time). Each code was applied 30 times, and the y-axis is shown in a logarithmic scale.

2. Rank-based fitness + tournament selection
3. Adaptive operator selection (OX, ER, BCR)
4. Memetic local search with stochastic 2-opt
5. Diversity preservation mechanisms

In particular, R1-generated GA is the only LLM-generated code that introduces a modification to the population initialization by incorporating the nearest neighbor heuristic. In contrast, both the original code and all other LLM-generated variants use the following initialization function:

```
# Function: Initial Population
def initial_population(population_size, distance_matrix):
    population = []
    for i in range(0, population_size):
        seed = seed_function(distance_matrix)
        population.append(seed)
    return population
```

Instead, R1-generated GA features the following initialization that leads to an improved performance.⁷

```
1 def initial_population(population_size, distance_matrix):
2     """Initialize population with mix of random and heuristic solutions. Combines diversity (random) with quality (NN) for
3     better exploration. Implements hybrid population initialization from modern metaheuristics."""
4     population = []
5     if population_size >= 5: # Include 20% NN seeds
6         for _ in range(max(1, population_size//5)):
7             population.append(nearest_neighbor_seed(distance_matrix))
```

⁷Note that, in all Python code snippets shown in this paper, '...' indicates omitted parts that are not relevant.

```

7
8     ...
9     return population
10 def nearest_neighbor_seed(distance_matrix):
11     """Generate initial solution using Nearest Neighbor heuristic. Provides high-quality initial seeds to accelerate
12     convergence. Based on constructive heuristic methods commonly used in TSP."""
    ...

```

In particular, the GA is initialized with 20% nearest neighbor solutions for population sizes of at least five individuals. This well-known TSP heuristic significantly speeds up convergence. In this way, R1 shows its ability to ‘dig’ into its knowledge base to choose an alternative population initialization method and implement it effectively.

4.5.2 Case 2: SA (R1-generated code)

Also in the case of SA, R1 identifies and utilizes two well-known mechanisms recognized for their efficiency in solving the TSP: (1) the Lundy-Mees adaptive cooling schedule for improved temperature control, introduced years after the original SA [22], and (2) the nearest neighbor heuristic for TSP. In the latter case, R1 integrates the nearest neighbor heuristic for the TSP in a way similar to what it did in the case of the GA, demonstrating a consistent pattern in leveraging effective initialization strategies.

4.5.3 Case 3: SARSA (O1-generated code)

The O1-generated SARSA code achieved the best results among the competitors. This is due to being the only code to make use of *Boltzmann Exploration* (see code below). Unfortunately, LLMs do not have the capacity to identify the exact source (book, scientific article, etc) from which the information about Boltzmann Exploration was extracted. However, after reviewing the code, it is likely that it was sourced from a 2017 paper (see [2]), which suggests a Boltzmann operator for SARSA applied to the TSP.

In fact, the code below shows that, unlike the original code, O1 not only applies a random operator to select the next unvisited city but also assigns a probability—derived from the `q_table` data structure—to this choice (line 8), making the selection more dynamic. Moreover, it avoids unnecessary abstractions (e.g., extra data structures) that could slow down the Python code.

```

1 ...
2 while len(visited) < num_cities:
3     unvisited = [city for city in range(num_cities) if city not in visited]
4     # Boltzmann exploration
5     q_values = q_table[current_city, unvisited]
6     exp_q = np.exp(q_values / temperature)
7     probabilities = exp_q / np.sum(exp_q)
8     next_city = np.random.choice(unvisited, p=probabilities)
9
10    reward = -distance_matrix_normalized[current_city, next_city]
11    visited.add(next_city)
12    route.append(next_city)
13    ...
14 ...

```

4.5.4 Case 4: BB (R1-generated code)

When studying why the BB code of R1 was faster than the original code, first we noticed that, like in cases 4.5.1 and 4.5.2, R1 made use of the nearest neighbor heuristic for initialization. Moreover, R1 modified the `explore_path` function of BB by dynamically sorting the next candidates by edge weight to prioritize the cheapest/nearest extensions first. Both updates are not trivial. R1 notes the following in the code comments: “*Enhancements reduce unnecessary branching and accelerate convergence through early solution bias.*”⁸

In addition, the code snippet below is not present in the original code. O1 introduces `current_node` and `candidates` efficiently, using slicing (line 4) and sorting with a lambda function (line 5) to enhance path exploration in BB. This new array-based data structure is both efficient and implemented in a Pythonic style to improve performance.

⁸This can be seen in line 48 of file `bb_deepseek_r1.py` of our online repository URL.

Table 5: Average Cyclomatic Complexity of the codes.

	Algorithm Codes	Average Complexity	Risk Category
	Original	6.95	B (Low - Well structured)
LLMs	Claude-3.5-Sonnet	5.60	A (Low - Simple)
	Gemini-exp-1206	7.34	B (Low - Well structured)
	Llama-3.3-70b	7.38	B (Low - Well structured)
	GPT-O1	6.84	B (Low - Well structured)
	DeepSeek-R1	7.51	B (Low - Well structured)

```

1 def explore_path(route, distance, distance_matrix, bound, weight, level, path, visited, min1_list, min2_list):
2     ...
3     current_node = path[level - 1]
4     candidates = [i for i in range(distance_matrix.shape[0]) if distance_matrix[current_node, i] > 0 and not visited[i]]
5     candidates = sorted(candidates, key=lambda x: distance_matrix[current_node, x])
6     ...

```

4.6 Code complexity

In the previous subsection, we analyzed the LLM-generated codes in terms of their performance. But do these codes also offer better readability and reduced complexity in comparison to the original codes? To address this question, we evaluate their *cyclomatic complexity*—a metric that quantifies the number of independent paths through a program’s source code. Through empirical research, Chen [6] demonstrated that a high cyclomatic complexity correlates with an increased bug prevalence. For our measurements, we employ the `radon` library for Python.⁹

As shown in Table 5, the Claude-generated codes have the lowest average cyclomatic complexity (5.60 points), which improves code readability but, as shown before, comes at the cost of performance. The other models’ codes and the original code have complexity scores between 6.84 (O1) and 7.51 (R1), which is still considered low and well-structured according to standard software engineering metrics. The values in the Risk Category column are taken from the documentation of the `radon` library.¹⁰

Finally, Figure 6 reveals that there are cases—such as the R1-generates codes in the case of GA and Christofides, or the O1-generated code for SARSA—in which the LLM-generated codes not only outperform the original codes, but also decrease the cyclomatic complexity.

Summary: Based on all evaluations presented in this paper, we can state that among the five LLMs tested, R1 generally produces the best results, followed by O1. Gemini performed well in certain cases, such as ACO and ALNS; however, it underperformed in others, such as, for example, Christofides. Among all tested models, Claude showed the lowest performance. In summary, *LLM-enhanced code versions clearly outperformed the original implementations in nine out of ten cases/algorithms*. Only for Q_Learning none of the models was able to improve the original code. In this case, Llama matched the performance of the original code.

4.7 Code Quality Analysis

In this subsection, we introduce two additional metrics to assess the quality of the generated code. Beyond the cyclomatic complexity analysis presented in the previous section, these metrics capture complementary and deeper aspects of code complexity and software maintainability.

⁹<https://pypi.org/project/radon/>.

¹⁰<https://radon.readthedocs.io/en/latest/commandline.html#the-cc-command>



Figure 6: Heatmap analysis of Cyclomatic Complexity across the ten optimization algorithms. The chart compares the complexity scores of the original implementations against the variants generated by the five LLMs. Lower values (indicated in green) represent reduced complexity and potentially higher readability, while higher values (indicated in orange/red) denote more complex code structures.

Maintainability Index (MI) is a composite metric ranging from 0 to 100 that combines Halstead volume, cyclomatic complexity, and lines of code into a single score indicating how maintainable the code is [28]. Higher values indicate more maintainable code.¹¹

Cognitive Complexity (CC) assesses the relative difficulty of understanding a codebase, specifically penalizing structural nesting and non-linear logic [5]. By weighting nesting depth more heavily than standard Cyclomatic Complexity, it provides a more realistic proxy for the mental effort required by developers; thus, lower values correlate with higher code readability.¹²

While Cyclomatic Complexity quantifies the number of linearly independent paths through the code, it fails to distinguish between flat and deeply nested structures of equivalent branch count. These two metrics address this limitation from different angles: MI provides a holistic maintainability score that aggregates multiple code properties into a single actionable index, whereas CC specifically isolates human readability—revealing how structural choices impact comprehension independently of path count.

As DeepSeek-R1 outperformed the other LLMs overall, we adopt it as a case study to compare the 10 baseline/original algorithms using these two metrics.

4.7.1 Case Study: DeepSeek-R1

Figure 7(A) shows the MI comparison. DeepSeek-R1 produced code that maintained or exceeded the baseline MI scores in the majority of cases, with the most notable improvement observed in Tabu Search. Although slight regressions were observed in three algorithms (ACO, ALNS, and Branch and Bound) and marginally in RL-QL, the overall average MI increased from 52.99 (Baseline) to 53.40 (DeepSeek-R1). This result suggests that the LLM successfully enhanced the functional performance of the algorithms without degrading the structural health of the codebase, preserving a level of maintainability consistent with the original implementations.

In contrast, Figure 7(B) presents the CC analysis, which measures the mental effort required for human comprehension. Here, the average score rose from 45.7 (Baseline) to 69.1 (DeepSeek-R1), representing an increase of approximately 51%. This sharp increase provides a critical insight into the optimization mechanism: DeepSeek-R1 did not merely “clean” or refactor the code; rather, it injected sophisticated logic—such as memetic local searches, heuristic initializations, and hybrid strategies—which inevitably complicates the implementation. Consequently, the generated code is algorithmically “smarter” and more efficient, but this comes at the cost of higher cognitive load for the human reader.

However, there are exceptions to this complexity-performance trade-off. One of these can be seen in the context of ALNS, where DeepSeek-R1 degraded both MI and CC metrics. Crucially, as observed in Figure 2, this increased structural complexity did not translate into performance gains over the baseline implementation. This highlights a specific failure case where the LLM introduced unnecessary complexity without achieving the intended algorithmic optimization, resulting in code that is harder to maintain yet functionally equivalent or inferior to the original one.

While the lower cyclomatic complexity observed in many LLM-generated variants suggests improved maintainability and readability—a finding further reinforced by the calculation of the MI—the stochastic nature of LLM generation can occasionally produce “brittle” code. Such code might exhibit instability across different runs or edge cases despite being structurally simple. This observation reinforces the critical importance of the iterative validation and refinement process described in our methodology (Section 3). Furthermore, to maximize security and stability, we recommend performing these validation executions within isolated environments, such as sandboxes or containers. This precaution not only protects the host system from potential security risks or infinite loops but also ensures a controlled and reproducible testing environment.

Implications for Automated Optimization Loops These results suggest that maintainability and cognitive metrics should be explicitly incorporated into the algorithm improvement loop. In this context, it is not sufficient for an algorithm to be syntactically and logically correct (see Table 3); it should ideally

¹¹As in the cyclomatic complexity analysis, we use the `radon` package for computing this metric.

¹²Computed using the `cognitive-complexity` Python package: https://github.com/Melevir/cognitive_complexity

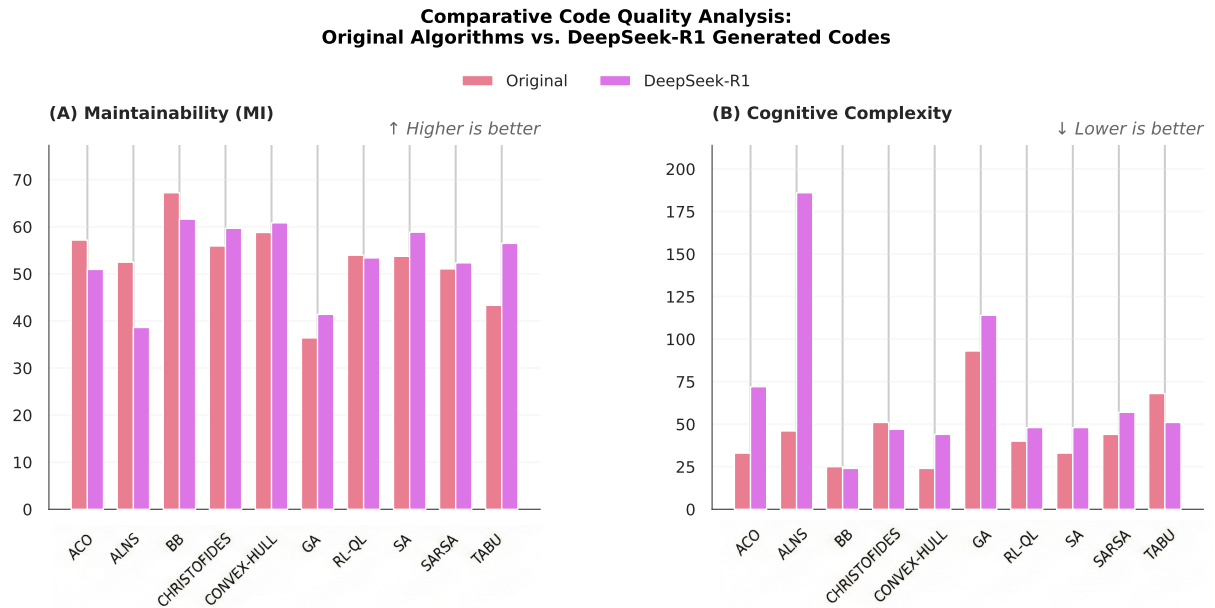


Figure 7: Comparative Code Quality Analysis between the original implementations and DeepSeek-R1 generated variants. **(A) Maintainability Index (MI):** Higher values indicate better software maintainability. The results show that DeepSeek-R1 preserved the structural health of the code, achieving an average score of 53.40, comparable to the baseline. **(B) Cognitive Complexity (CC):** Lower values indicate easier human comprehension. The observed increase in CC for the generated variants reflects the injection of sophisticated algorithmic logic—such as hybrid heuristics and local searches—which significantly improves performance but increases the mental effort required to understand the implementation.

preserve software maintainability. In other words, superior performance in terms of output quality (e.g., solution quality for the TSP) must be weighed against the potential difficulty for a human programmer to understand, modify, or maintain the resulting software.

Future iterations of this methodology could incorporate MI and CC metrics directly into the prompt, comparing them with the original algorithm to allow the model to make more informed decisions during generation. However, this refinement process must account for economic constraints, as including detailed metric analysis in the context window increases token consumption.

5 Discussion

5.1 Limitations and Methodological Considerations

While our results are promising, a nuanced understanding of the challenges is crucial for the practical application and future development of this approach.

- **The “Black Box” Nature and Its Risks.** A primary challenge, highlighted in [16], is that LLMs cannot precisely trace the sources of their suggestions. This opacity makes it difficult to pinpoint which specific modifications led to performance gains. More importantly, it introduces the risk of the LLM generating “hallucinated” algorithmic components or incorrectly combining concepts from different sources, potentially leading to subtle logical flaws that are not immediately apparent through basic testing. Exploring the question of “*where specifically does the generated code come from?*” is therefore not just a fascinating research avenue but a critical step towards building more reliable systems.
- **The Practical Costs of Iterative Refinement.** Our “simple prompt” strategy successfully lowers

the barrier to entry, but the overall process is not without cost. As shown in Table 3, several algorithms required multiple manual correction rounds to become functional. This iterative cycle of generating, testing, and providing feedback demands significant user time and effort. This reveals a key trade-off: the ease of initial prompting versus the manual cost of validation. For this methodology to be truly effective in practice, the efficiency of this human-in-the-loop refinement process must be considered.

- **Performance Variability and Failure Analysis.** Our results reveal significant performance variability among LLMs. An interesting insight comes from the failure of all models to enhance the Q-Learning baseline (Figure 3). We attribute this to the fundamental difference between refining heuristic logic and optimizing Reinforcement Learning (RL) agents. While LLMs excel at proposing structural improvements for heuristics (e.g., integrating a 2-opt local search), improving RL algorithms often requires fine-tuning the delicate interplay between state-space representations, reward functions, and hyperparameters. These elements are highly sensitive to the specific problem instance and typically require dynamic interaction with the environment to optimize, rather than static code refactoring. Consequently, our standard prompting strategy, which relies on static code analysis, proved insufficient for this specific paradigm. Similarly, the weaker performance on complex algorithms like ALNS suggests a scarcity of high-quality training data for such intricate heuristics compared to more common methods like GA or SA.
- **Generalization and Training Data Dependence.** A critical question regarding the scope of this study is the transferability of findings beyond the TSP. While our methodology—specifically the prompting strategy—is technically problem-agnostic, its effectiveness is intrinsically linked to the representation of the specific problem within the LLM’s training corpora. Since the TSP is a foundational problem with vast open-source coverage, models possess a strong “prior” on algorithmic components for TSP algorithms. We hypothesize that this success readily transfers to other well-known combinatorial optimization problems such as Vehicle Routing Problems (VRPs) or Knapsack Problems (KPs). However, performance may degrade for highly niche, novel, or proprietary problems where the model lacks sufficient pre-training examples to draw upon, representing a boundary condition for this approach. Nevertheless, this limitation could potentially be mitigated by prompting the model to explicitly abstract general heuristics from known problems (like the TSP) that are applicable to the new domain. Instead of transferring all pre-training knowledge, the model would focus on extracting only those operators relevant to the unseen problem structure.
- **Positioning within the State-of-the-Art.** As detailed in Table 1, our work is methodologically distinct from other recent approaches. While prominent frameworks like LLaMEA/ReEvo, AlphaEvolve, and FunSearch employ complex, automated systems—such as evolutionary loops or program search trees—to generate or discover *new* algorithmic components, our approach uses a simple, interactive prompting strategy. This reflects a fundamental difference in objective and user role: expert-centric frameworks focus on automated algorithm discovery or generation, requiring users to design a sophisticated search process. In contrast, our practitioner-centric method centers on the collaborative enhancement of *complete, existing codebases*. Furthermore, while the work [35] established the feasibility of this concept in a single case study, the present research validates its effectiveness systematically across a broad and diverse range of algorithms.

5.2 Directions for Future Research

The limitations identified above naturally lead to several exciting directions for future research.

- **Broadening the Problem Scope.** A crucial next step, addressing a limitation of our current study, is to validate these findings beyond the TSP. Applying this methodology to combinatorial optimization problems with different structures, such as the Vehicle Routing Problem (VRP) or the Knapsack Problem, is essential to determine the true generality of this approach.
- **Automating the Refinement Loop.** To mitigate the manual effort of validation and correction, future work should focus on automating the feedback cycle. This could involve developing systems

where an LLM agent can not only generate code but also autonomously execute it against a benchmark suite, analyze the results (including errors and performance metrics), and iteratively refine its own suggestions.

- **Specialized Models and Benchmarks.** As LLMs evolve, it will be necessary to continuously test their capabilities in a structured way. The creation of specialized benchmarks with baseline implementations for multiple optimization problems would be invaluable. Furthermore, developing fine-tuned LLMs specialized in optimization could overcome the data scarcity issues observed for complex or less-common algorithms, potentially leading to even more significant and reliable improvements. This aligns with trends towards “reasoning models” that prioritize response quality over speed [37].

6 Conclusion

Our systematic benchmark demonstrates that Large Language Models (LLMs) can be powerful collaborators in enhancing classical optimization algorithms. We have shown that a simple, reproducible prompting strategy is sufficient to improve 10 baseline algorithms for the Travelling Salesman Problem, often yielding higher-quality solutions, faster computation times, and reduced code complexity. Crucially, these improvements were achieved holistically on complete codebases through the incorporation of modern heuristics and data structures. This entire methodology is fully reproducible via the chatbot on our project website (*Here will be the url*), providing a practical pathway for practitioners to upgrade complex algorithms without deep theoretical expertise.

Building on this foundation, future work will address the limitations identified in our discussion. To validate the generality of our findings, these methods will be extended to other combinatorial optimization problems with diverse structures. Furthermore, to overcome the practical costs of manual refinement, we plan to explore the integration of LLM-based agents to create a fully automated and continuous improvement cycle for existing algorithms.

References

- [1] Anthropic Team. Introducing Claude 3.5 Sonnet — anthropic.com. <https://www.anthropic.com/news/claude-3-5-sonnet>, 2024. [Accessed 02-11-2024].
- [2] Kavosh Asadi and Michael L Littman. An alternative softmax operator for reinforcement learning. In *International Conference on Machine Learning*, pages 243–252. PMLR, 2017.
- [3] Liam Barkley and Brink van der Merwe. Investigating the Role of Prompting and External Tools in Hallucination Rates of Large Language Models, 2024. URL <https://arxiv.org/abs/2410.19385>.
- [4] Christian Blum and Andrea Roli. Metaheuristics in combinatorial optimization: Overview and conceptual comparison. *ACM Computing Surveys*, 35(3):268–308, 2003.
- [5] G. Ann Campbell. Cognitive complexity: an overview and evaluation. In *Proceedings of the 2018 International Conference on Technical Debt*, TechDebt '18, page 57–58, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450357135. doi: 10.1145/3194164.3194186. URL <https://doi.org/10.1145/3194164.3194186>.
- [6] Changqi Chen. An Empirical Investigation of Correlation between Code Complexity and Bugs, 2019. URL <https://arxiv.org/abs/1912.01142>.
- [7] Nicos Christofides. Worst-Case Analysis of a New Heuristic for the Travelling Salesman Problem. *Operations Research Forum*, 3(1):20, Mar 2022. ISSN 2662-2556. doi: 10.1007/s43069-021-00101-z. URL <https://doi.org/10.1007/s43069-021-00101-z>.
- [8] Pham Vu Tuan Dat, Long Doan, and Huynh Thi Thanh Binh. HSEvo: elevating automatic heuristic design with diversity-driven harmony search and genetic algorithm using LLMs. In *Proceedings*

- of the Thirty-Ninth AAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'25/IAAI'25/EAAI'25. AAAI Press, 2025. ISBN 978-1-57735-897-8. doi: 10.1609/aaai.v39i25.34898. URL <https://doi.org/10.1609/aaai.v39i25.34898>.
- [9] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, et al. A survey on in-context learning. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pages 1107–1128, 2024.
- [10] M. Dorigo and L.M. Gambardella. Ant colony system: a cooperative learning approach to the traveling salesman problem. *IEEE Transactions on Evolutionary Computation*, 1(1):53–66, 1997. doi: 10.1109/4235.585892.
- [11] Samuel Eilon, C. D. T. Watson-Gandy, Nicos Christofides, and Richard de Neufville. Distribution Management-Mathematical Modelling and Practical Analysis. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-4(6):589–589, 1974. doi: 10.1109/TSMC.1974.4309370.
- [12] Fred W. Glover. Tabu Search - Part I. *INFORMS J. Comput.*, 1:190–206, 1989. URL <https://api.semanticscholar.org/CorpusID:5617719>.
- [13] D. Guo, D. Yang, H. Zhang, and et al. DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning. *Nature*, 645(8081):633–638, September 2025.
- [14] John H. Holland. *Adaptation in Natural and Artificial Systems*. University of Michigan Press, Ann Arbor, MI, 1975. URL <https://mitpress.mit.edu/books/adaptation-natural-and-artificial-systems>.
- [15] Sen Huang, Kaixiang Yang, Sheng Qi, and Rui Wang. When large language model meets optimization. *Swarm Evol. Comput.*, 90:101663, 2024. URL <https://doi.org/10.1016/j.swevo.2024.101663>.
- [16] Muhammad Khalifa, David Wadden, Emma Strubell, Honglak Lee, Lu Wang, Iz Beltagy, and Hao Peng. Source-Aware Training Enables Knowledge Attribution in Language Models, 2024. URL <https://arxiv.org/abs/2404.01019>.
- [17] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by Simulated Annealing. *Science*, 220(4598):671–680, 2025/02/28/ 1983. URL <http://www.jstor.org/stable/1690046>. Full publication date: May 13, 1983.
- [18] Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and Xiaohang Dong. Better zero-shot reasoning with role-play prompting. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4099–4113, 2024.
- [19] Jia Li, Chongyang Tao, Jia Li, Ge Li, Zhi Jin, Huangzhao Zhang, Zheng Fang, and Fang Liu. Large language model-aware in-context learning for code generation. *ACM Trans. Softw. Eng. Methodol.*, 34(7), August 2025. ISSN 1049-331X. doi: 10.1145/3715908. URL <https://doi.org/10.1145/3715908>.
- [20] Fei Liu, Xialiang Tong, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu Zhang. Evolution of heuristics: towards efficient automatic algorithm design using large language model. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org, 2024.
- [21] Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew Soon Ong. Large language models as evolutionary optimizers. *2024 IEEE Congress on Evolutionary Computation (CEC)*, pages 1–8, 2023. URL <https://api.semanticscholar.org/CorpusID:264829031>.
- [22] M. Lundy and A. Mees. Convergence of an annealing algorithm. *Mathematical Programming*, 34(1): 111–124, Jan 1986. ISSN 1436-4646. doi: 10.1007/BF01582166. URL <https://doi.org/10.1007/BF01582166>.

- [23] Manuel López-Ibáñez, Jérémie Dubois-Lacoste, Leslie Pérez Cáceres, Thomas Stützle, and Mauro Birattari. The irace package: Iterated Racing for Automatic Algorithm Configuration. *Operations Research Perspectives*, 3:43–58, 2016. doi: 10.1016/j.orp.2016.09.002.
- [24] YINGWEI MA, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. At which training stage does code data help LLMs reasoning? In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=KIPJKST4gw>.
- [25] R. Martí and G. Reinelt. *The Linear Ordering Problem: Exact and Heuristic Methods in Combinatorial Optimization*. Applied Mathematical Sciences. Springer Berlin Heidelberg, 2011.
- [26] David R. Morrison, Sheldon H. Jacobson, Jason J. Sauppe, and Edward C. Sewell. Branch-and-bound algorithms: A survey of recent advances in searching, branching, and pruning. *Discrete Optimization*, 19:79–102, 2016. ISSN 1572-5286. doi: <https://doi.org/10.1016/j.disopt.2016.01.005>. URL <https://www.sciencedirect.com/science/article/pii/S1572528616000062>.
- [27] Alexander Novikov, Ngàn Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Pushmeet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and algorithmic discovery, 2025. URL <https://arxiv.org/abs/2506.13131>.
- [28] P. Oman and J. Hagemester. Metrics for assessing a software system’s maintainability. In *Proceedings Conference on Software Maintenance 1992*, pages 337–344, 1992. doi: 10.1109/ICSM.1992.242525.
- [29] OpenAI et al. GPT-4 Technical Report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- [30] V. Pereira. pyCombinatorial - A library to solve TSP (Travelling Salesman Problem) using Exact Algorithms, Heuristics and Metaheuristics. <https://github.com/Valdecy/pyCombinatorial>, 2022. Accessed: 2025-03-10.
- [31] Jean-Yves Potvin. Genetic algorithms for the traveling salesman problem. *Annals of Operations Research*, 63(3):337–370, Jun 1996. ISSN 1572-9338. doi: 10.1007/BF02125403. URL <https://doi.org/10.1007/BF02125403>.
- [32] Gerhard Reinelt. TSPLIB - A Traveling Salesman Problem Library. *INFORMS J. Comput.*, 3(4):376–384, 1991. URL <http://dblp.uni-trier.de/db/journals/informs/informs3.html#Reinelt91>.
- [33] Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M. Pawan Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan S. Ellenberg, Pengming Wang, Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, Jan 2024. ISSN 1476-4687. doi: 10.1038/s41586-023-06924-6. URL <https://doi.org/10.1038/s41586-023-06924-6>.
- [34] Stefan Ropke and David Pisinger. An Adaptive Large Neighborhood Search Heuristic for the Pickup and Delivery Problem with Time Windows. *Transportation Science*, 40(4):455–472, 2025/02/28/2006. URL <http://www.jstor.org/stable/25769321>. Full publication date: November 2006.
- [35] Camilo Chacón Sartori and Christian Blum. Optimizing the Optimizer: An Example Showing the Power of LLM Code Generation. In *2025 20th Conference on Computer Science and Intelligence Systems (FedCSIS)*, pages 47–57, 2025. doi: 10.15439/2025F1481.
- [36] Sander Schulhoff, Michael Ilie, and et al. The Prompt Report: A Systematic Survey of Prompt Engineering Techniques, 2025. URL <https://arxiv.org/abs/2406.06608>.
- [37] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.

-
- [38] Niki van Stein and Thomas Bäck. LLaMEA: A large language model evolutionary algorithm for automatically generating metaheuristics. *IEEE Transactions on Evolutionary Computation*, 29(2): 331–345, 2025. doi: 10.1109/TEVC.2024.3497793.
- [39] Gemini Team and et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL <https://arxiv.org/abs/2403.05530>.
- [40] Meta Team. The Llama 3 Herd of Models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- [41] Jiaying Wang, Chenglong Xiao, Shanshan Wang, and Yaqi Ruan. Reinforcement learning for the traveling salesman problem: Performance comparison of three algorithms. *The Journal of Engineering*, 2023(9):e12303, 2023. URL <https://api.semanticscholar.org/CorpusID:261504600>.
- [42] Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Siddartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. LiveBench: A Challenging, Contamination-Free LLM Benchmark, 2024. URL <https://arxiv.org/abs/2406.19314>.
- [43] Marco Wiering and Martijn van Otterlo. *Reinforcement Learning: State-of-the-Art*. Springer Publishing Company, Incorporated, 2014.
- [44] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Bb4VGOWELI>.
- [45] Haoran Ye, Jiarui Wang, Zhiguang Cao, Federico Berto, Chuanbo Hua, Haeyeon Kim, Jinkyoo Park, and Guojie Song. Reevo: Large language models as hyper-heuristics with reflective evolution. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 43571–43608. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/4ced59d480e07d290b6f29fc8798f195-Paper-Conference.pdf.