# HIERARCHICAL REPRESENTATIONS FOR CROSS-TASK AUTOMATED HEURISTIC DESIGN USING LLMS

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

016

017

018

019

021

025

026027028

029

031

032

033

034

037

038

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

Designing heuristic algorithms for complex optimization problems is a timeconsuming and expert-driven process. Recently, Automated Heuristic Design (AHD) using Large Language Models (LLMs) has shown significant promise for automating algorithm development. However, existing works mainly rely on programs to represent heuristics, which are inherently task-specific and fail to generalize as effectively as established metaheuristics like tabu search or guided local search. To bridge this gap, we introduce Multi-Task Hierarchical Search (MTHS), an LLM-guided evolutionary method that co-designs general-purpose metaheuristics and task-specific programs. MTHS employs a hierarchical representation and adopts a two-level evolution framework to evolve task-agnostic metaheuristics and task-specific program implementations simultaneously across multiple heuristic design tasks. During this evolution, a knowledge transfer mechanism allows learning from elite programs designed for other tasks. We evaluated MTHS on distinct combinatorial optimization problems, where it outperforms both commonly-used heuristics and existing LLM-driven AHD approaches. Our results demonstrate that the hierarchical representations facilitate effective multi-task AHD, and the evolved metaheuristics exhibit strong generalization to related tasks.

## 1 Introduction

Designing high-performance heuristic algorithms for complex problem-solving tasks is a notoriously challenging endeavor, traditionally relying on a time-consuming, expert-driven process of trial and error. Recently, Large Language Model (LLM)-driven Automated Heuristic Design (AHD) (Liu et al., 2024b; Ye et al., 2024; Zheng et al., 2025; Ye et al., 2025) has emerged as a powerful paradigm to automate algorithm development and mitigate this tedious process. This approach has already demonstrated its potential by automating the design of high-performance heuristics in diverse optimization domains including combinatorial optimization (Liu et al., 2024b; Ye et al., 2024), blackbox optimization (van Stein & Bäck, 2024; Xie et al., 2025a), and Bayesian optimization (Yao et al., 2024).

A prevalent strategy in LLM-driven AHD is to embed LLMs as heuristic designers within iterative search frameworks (Zhang et al., 2024). Various search paradigms have been explored, from Evolutionary Computation (EC) (Liu et al., 2024b; Ye et al., 2024; Dat et al., 2025; Yao et al., 2025) to Monte Carlo Tree Search (MCTS) (Zheng et al., 2025). For instance, EoH (Liu et al., 2024b) evolves both natural language thoughts and executable code, ReEvo (Ye et al., 2024) integrates reflection strategies to refine the design process, and MCTS-AHD (Zheng et al., 2025) organizes heuristics in a tree to systematically explore the heuristic space.

However, a fundamental limitation persists in current AHD methods: they produce monolithic, task-specific heuristics. These approaches typically represent heuristics as either low-level programs (Zheng et al., 2025) or high-level thoughts (Liu et al., 2024b). Task-specific programs offer limited portability to new problems, while high-level thoughts are often too abstract to guarantee a direct correspondence with a high-performing implementation (Liu et al., 2024b). Consequently, existing systems must essentially restart the discovery process for each new problem, failing to institutionalize learning and generalize algorithmic knowledge across domains. This stands in stark contrast to human experts, who design and reuse metaheuristics, such as tabu search (Glover & Laguna, 1998) or simulated annealing (Van Laarhoven & Aarts, 1987), as general-purpose meta-

heuristics that are effective across a vast range of optimization tasks (Gendreau et al., 2010; Martí et al., 2025).

To bridge this gap, we argue that the key lies in creating hierarchical representations that separate general algorithmic logic from task-specific components, enabling cross-task automated heuristic design. We introduce the Multi-Task Hierarchical Search (MTHS), an LLM-guided hierarchical evolutionary framework designed to co-design general-purpose metaheuristics and their task-specific program implementations across multiple tasks simultaneously. MTHS leverages its hierarchical structure to explicitly transfer knowledge across tasks, allowing effective programs discovered in one task to inform and accelerate program design in others. Our primary contributions are threefold:

- We propose a hierarchical representation for LLM-driven AHD that consists of a task-agnostic metaheuristic and its task-specific program instantiations. This mirrors expert practice, where general metaheuristics are paired with tailored implementations.
- We introduce the MTHS framework, which jointly designs the general metaheuristic and its
  task-specific implementations across diverse optimization tasks. At the high level, MTHS
  evolves metaheuristics; at the low level, it creates and refines programs and their associated
  key functions for each task. A cross-task knowledge transfer is adopted to learn from elite
  programs from other tasks.
- We conduct extensive experiments on diverse combinatorial optimization problems. MTHS
  consistently discovers heuristics that outperform widely used heuristic baselines and stateof-the-art LLM-driven AHD methods. Crucially, the evolved metaheuristics exhibit strong
  generalization to related problems.

## 2 Multi-task Hierarchical Search

#### 2.1 HIERARCHICAL REPRESENTATION

This work addresses the problem of automated heuristic design across multiple, related tasks. The central goal is to discover high-level, general-purpose metaheuristics that can be specialized to achieve superior performance across multiple tasks. Formally, we are given a set of m tasks,  $\mathcal{T} = \{T_1, \ldots, T_m\}$ . Each task  $T_t$  is defined by a concise natural language description  $D_t$ , a program template  $Temp_t$  providing the necessary inputs and outputs for execution, and a black-box evaluation function  $E_t(\cdot)$  that returns a scalar performance score for a given program. Without loss of generality, we consider minimization problems in this paper.

Our representation for a candidate, which we term an *individual*  $I_i$ , is composed of two hierarchical levels: To be precise, each individual represents a complete metaheuristic for our multi-task AHD, encompassing both its high-level *metaheuristic* description and its task-specific program implementations. This structure is composed of two levels:

- 1. Task-Agnostic Metaheuristic  $(MH_i)$ : At the highest level is a general-purpose *metaheuristic*,  $MH_i$ . Represented as a high-level algorithmic description, it captures the core problem-solving logic independent of any specific task.
- 2. Task-Specific Programs  $(X_{i,t})$ : For each task  $T_t$ , the metaheuristic  $MH_i$  is instantiated into a concrete, executable program,  $X_{i,t}$ . This program adapts the general logic of  $MH_i$  to the specific requirements of task  $T_t$ . Within each program  $X_{i,t}$ , we identify a performance-critical key function, denoted  $F_{i,t}$ .

The performance of an individual  $I_i$  is evaluated based on the collective performance of its instantiated programs. Let  $S_{i,t} = E_t(X_{i,t})$  be the score obtained by program  $X_{i,t}$  on task  $T_t$ . The score list  $\{S_{i,1}, \ldots, S_{i,m}\}$  is assigned to each individual  $I_i$ , which will be used in population management.

To ensure clarity, we will use the term *individual* to refer to this entire hierarchical entity and *metaheuristic* to refer specifically to the high-level description within it. While the distinction between "heuristic" and "metaheuristic" lacks a universal consensus in the literature (Gendreau et al., 2010; Martí et al., 2025), we adopt the view that metaheuristics represent a more general problem-solving paradigm. Nevertheless, given their conceptual overlap, we may use these terms interchangeably where the context allows.

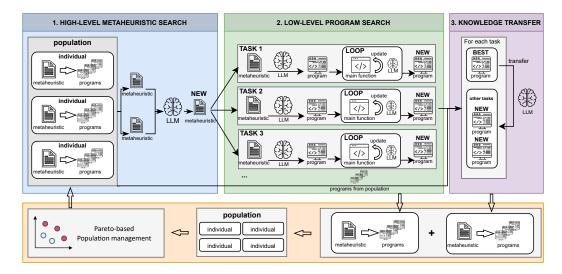


Figure 1: Overview of the MTHS pipeline. The pipeline consists of three main components: (1) High-level metaheuristic search, (2) Low-level program search, and (3) Knowledge transfer. In the high-level search, a population of metaheuristics is evolved, where each individual contains one metaheuristic paired with m task-specific programs. For each newly generated metaheuristic, a low-level program search is performed for each task: a program is created for each task, and its key function is identified and refined using LLMs. This produces a new candidate comprising the metaheuristic and its m programs. Next, a knowledge transfer phase is applied for each task: the best-performing program across both the existing population and the new candidate is identified and used to update the other m-1 programs within the same metaheuristic. Candidates produced from both low-level search and knowledge transfer are added to the population. Finally, a Pareto-based population management step selects individuals to form the next generation.

## 2.2 Framework

We introduce Multi-Task Heuristic Search (MTHS), a framework that automates heuristic design across a set of related tasks using a two-level evolutionary algorithm (see Figure 1 and Algorithm 1). The process begins by prompting LLMs with descriptions of all tasks to seed an initial high-level population ( $\mathcal{P}_H$ ) of diverse, task-agnostic metaheuristics. Each metaheuristic represents a general problem-solving strategy intended to be effective across multiple tasks. For each of these metaheuristics, MTHS initiates a distinct low-level search for every individual task. This low-level process evolves separate populations ( $\mathcal{P}_{L,t}$ ) of task-specific programs. A knowledge transfer mechanism then shares insights from the best-performing programs across tasks. This hierarchical structure enables MTHS to simultaneously conduct broad strategic exploration at the shared metaheuristic level and specialized, fine-grained program optimization at the individual task level. We introduce each phase as follows. We expand the subalgorithms and present the detailed specific prompts in Appendix B.

## 2.3 HIGH-LEVEL EVOLUTION

The high-level evolution maintains a population of individuals,  $\mathcal{P}_H$ . In each generation, we employ the LLM as an evolutionary operator to generate a new candidate metaheuristic. The process begins by selecting a set of k parent individuals,  $\{I_1,\ldots,I_k\}$ , from  $\mathcal{P}_H$ . The corresponding metaheuristic descriptions of these parents,  $\{MH_1,\ldots,MH_k\}$ , are then formatted into a carefully designed prompt. This prompt instructs the LLM,  $\mathcal{L}$ , to synthesize a novel and potentially superior metaheuristic, adhering to a predefined description template. The LLM's textual output constitutes the metaheuristic description,  $MH_{new}$ , for the new offspring.

```
162
           Algorithm 1 Multi-Task Heuristic Search (MTHS)
163
164
             1: \mathcal{T} = \{T_1, \dots, T_m\}: Set of m tasks, each with description D_t, template Temp_t, and evaluator
166
                L: Large Language Model
167
             3: N_{eval}: total evaluation limits
             4: N_H, k, N_L: High-level population size and number of parents, Low-level evaluation budget
169
           Output: The final population of high-performing individuals \mathcal{P}_H
             5: procedure MTHS(\mathcal{T}, \mathcal{L}, N_{eval}, N_H, k)
170
             6:
                      \mathcal{P}_H \leftarrow \emptyset
171
             7:
                      InitialMHs \leftarrow \mathcal{L}(BuildInitialPrompt(\{D_t\}_{t=1}^m))
172
             8:
                      for each MH_{init} in InitialMHs do
173
                            I_{new} \leftarrow \text{LowLevelEvolution}(MH_{init}, \mathcal{T}, N_L, \mathcal{L})
             9:
174
                            \mathcal{P}_H \leftarrow \mathcal{P}_H \cup \{I_{new}\}
           10:
175
                      while evaluation count \leq N_{eval} do
           11:
176
           12:
                            \{I_1,\ldots,I_k\} \leftarrow \text{SelectParents}(\mathcal{P}_H,k)
177
           13:
                            \{MH_1,\ldots,MH_k\}\leftarrow\{I_1,\ldots,I_k\}
178
                            prompt \leftarrow BuildEvolutionPrompt(\{D_t\}_{t=1}^m, \{MH_i\}_{i=1}^k)
           14:
179
           15:
                            MH_{new} \leftarrow \mathcal{L}(prompt)
                            I_{new} \leftarrow \text{LowLevelEvolution}(MH_{new}, \mathcal{T}, N_L, \mathcal{L})
                                                                                                                                      ⊳ Sec. 2.4
           16:
181
           17:
                            if I_{new} is valid then
                                 I_{new} \leftarrow \text{KnowledgeTransfer}(I_{new}, \mathcal{T}, \mathcal{L})
                                                                                                                                      ⊳ Sec. 2.5
            18:
183
           19:
                                 \mathcal{P}_H \leftarrow \text{UpdatePopulation}(\mathcal{P}_H \cup \{I_{new}\}, N_H)
                                                                                                                                      ⊳ Sec. 2.6
184
           20:
                      return \mathcal{P}_H
185
186
           21: procedure LOWLEVELEVOLUTION(MH_{new}, \mathcal{T}, N_L, \mathcal{L})
                      I_{new} \leftarrow \text{new Individual with } MH_{new}
187
           22:
           23:
                      for t \leftarrow 1 to m do
188
                            X_{new,t} \leftarrow \mathcal{L}(\text{BuildProgramPrompt}(MH_{new}, T_t.Temp, \mathcal{L}))
           24:
189
                            F_{new,t} \leftarrow \mathcal{L}(\text{BuildKeyFuncPrompt}(T_t.D, X_{new,t}, \mathcal{L}))
           25:
190
           26:
                            (X_{new,t}^*, S_{new,t}^*) \leftarrow \text{EvolveKeyFunction}(X_{new,t}, F_{new,t}, T_t.E, \mathcal{L})
191
           27:
                            I_{new}.X_{new,t} \leftarrow X^*_{new,t}
192
           28:
                            I_{new}.S_{new,t} \leftarrow S^*_{new,t}
193
           29:
                      return I_{new}
194
```

## 2.4 Low-Level Evolution

196 197 198

199 200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

Each newly generated metaheuristic  $MH_{new}$  must be instantiated and optimized for all m tasks to determine its fitness. This evaluation is a multi-step procedure executed for each task  $T_t$ . First, in i) Task-Specific Program Generation, the LLM generates a full, compilable program  $X_{new,t}$  by integrating the logic of  $MH_{new}$  with the task-specific template  $Temp_t$ . Second, during ii) Key Function Identification, the LLM is prompted to analyze the generated code  $X_{new,t}$  and identify its most performance-critical component, which we designate as the key function  $F_{new,t}$ . Third, a dedicated low-level evolutionary search is performed to refine the key function in a process of iii) Key Function Refinement. An ephemeral population  $\mathcal{P}_{L,t}$  is initialized with variants of  $F_{new,t}$  generated by the LLM. This population then undergoes a short evolutionary process for a fixed budget of  $N_L$  evaluations. The LLM acts as a mutation operator, creating new function variants from existing high-performing ones. Each new variant is injected back into the base program  $X_{new,t}$  and evaluated using  $E_t(\cdot)$ . Finally, for iv) Fitness Assignment, after the low-level search concludes, the best-performing key function variant,  $F_{new,t}^*$ , is identified. The program incorporating this optimized function,  $X_{new,t}^*$ , yields the fitness score  $S_{new,t}$  for the metaheuristic  $MH_{new}$  on task  $T_t$ .

Once this process is completed for all m tasks, a new individual  $I_{new}$  is formed, comprising the metaheuristic  $MH_{new}$ , its vector of scores  $\{S_{new,1},\ldots,S_{new,m}\}$ , and the set of optimized programs  $\{X_{new,1}^*,\ldots,X_{new,m}^*\}$ . This individual is then added to the high-level population  $\mathcal{P}_H$ .

## 2.5 KNOWLEDGE TRANSFER

To facilitate explicit cross-task learning, we introduce a knowledge transfer phase. For the new individual  $I_{new}$ , we identify its best-performing program,  $X^*_{new,src}$ , on some source task  $T_{src}$ . The LLM is then prompted to adapt the logic of this program to every other target task  $T_{tgt}$  (where  $tgt \neq src$ ). This adaptation creates a new set of candidate programs. If an adapted program for  $T_{tgt}$  achieves a better score than the incumbent program  $X^*_{new,tgt}$ , it replaces it. This process directly transfers successful algorithmic patterns discovered on one task to others within the context of the same metaheuristic,  $MH_{new}$ .

## 2.6 PARETO-BASED POPULATION MANAGEMENT

MTHS uses a Pareto-based survival strategy to manage the high-level population  $\mathcal{P}_H$ . Since each metaheuristic is evaluated on m tasks, its performance is represented by a score vector  $\mathbf{S}_i = (S_{i,1}, \dots, S_{i,m})$ , framing the search as a multi-objective optimization problem. The most straightforward way to tackle multi-objective search is to transfer the multiple objectives into a single objective using some scalarization method, such as weighted-sum. However, it is hard to determine proper weights because the task scores are on different scales.

Therefore, we adopt a Pareto-based approach that works as follows: i): Task Champions (Elitism): For each task  $t \in \{1, \dots, m\}$ , the individual with the highest score  $S_{i,t}$  on that task is automatically preserved for the next generation. This ensures that the best-known performance on any single task is never lost. ii): Pareto Dominance Ranking: All remaining individuals in the candidate pool are ranked based on Pareto dominance. An individual  $I_i$  is said to dominate  $I_j$  if:  $(\forall t, S_{i,t} \geq S_{j,t}) \land (\exists t', S_{i,t'} > S_{j,t'})$ , where missing scores are treated as  $-\infty$ . Individuals are sorted into non-dominated fronts. iii): Selection and Truncation: The next generation is populated by adding individuals from the first non-dominated front, then the second, and so on, until the population size  $N_H$  is reached. If adding an entire front would exceed the population size, individuals from that front are selected based on their average scores on all tasks. An illustration of populations in objective space is presented in Appendix E.

## 3 EXPERIMENTAL STUDIES

## 3.1 TASKS AND DATASETS

We evaluate our method on four combinatorial optimization problems: the Traveling Salesman Problem (TSP), Capacitated Vehicle Routing Problem (CVRP), Flow Shop Scheduling Problem (FSSP), and Bin Packing Problem (BPP). For each problem, we generate a set of 64 diverse instances for the heuristic evolution phase. The final performance of the evolved heuristics is then validated on established, standard benchmark datasets. Further details on instance generation and benchmark specifics are provided in Appendix C.

- Traveling Salesman Problem (TSP): We aim to find the shortest tour visiting a set of locations. Our evolution set consists of 100node instances with locations uniformly sampled in  $[0, 1]^2$ . Fitness is the average optimality gap relative to the Concorde solver (Applegate et al., 2006). For final evaluation, we use standard instances from TSPLib (Reinelt, 1991).
- Capacitated Vehicle Routing Problem (CVRP): The goal is to design minimum-cost routes for a fleet of capacitated vehicles to serve a set of customers. The evolution set contains 100-customer instances. Fitness is measured as the average gap to solutions found by the LKH3 solver (Helsgaun, 2017). We test the final heuristics on CVRPLib benchmarks (Uchoa et al., 2017).
- Flow Shop Scheduling Problem (FSSP): We seek to schedule n jobs on m machines to minimize the makespan (total completion time). Our evolution instances feature 50 jobs and a variable number of machines ( $m \in [2, 20]$ ). Fitness is the average makespan. Final validation is performed on the Taillard benchmark suite (Taillard, 1993).
- Bin Packing Problem (BPP): The objective is to pack items of various sizes into the minimum number of fixed-capacity bins. Following prior work (Ye et al., 2024; Zheng et al., 2025), our evolution instances feature a bin capacity of 150 and item sizes sampled from [20, 100]. Fitness is the average number of bins used.

Table 1: Results on standard benchmark instances from TSPLib and CVRPLib (Sets A, B, E, F, M, P, and X). The table reports the average percentage gap to the best-known solutions for four distinct groups of heuristics: constructive heuristics, metaheuristics, LLM-designed heuristics, and our methods. The best result for each instance set is highlighted in **bold**, and the second-best one is underlined.

		T	SPLib			CVRPLib							
	50-99	100-199	200-499	500-1000	A	В	E	F	M	P	X		
NN	27.07	23.76	24.79	26.57	39.40	42.32	41.51	60.01	52.88	36.18	27.63		
Insert	13.99	16.15	20.00	26.30	33.86	33.07	32.51	65.45	44.49	25.96	31.19		
Or-tools SA	3.01	3.74	4.62	10.08	6.58	5.40	7.48	9.44	15.54	5.60	7.82		
Or-tools TS	1.81	3.20	4.62	10.11	1.05	1.12	1.57	4.56	5.74	1.11	6.06		
Or-tools GLS	0.63	1.62	3.34	6.84	1.24	1.14	1.30	3.49	7.74	1.07	6.29		
MS	1.82	2.92	4.15	6.58	8.19	11.60	10.25	9.65	42.66	7.52	42.85		
ALNS	1.62	1.90	5.24	8.28	6.39	5.80	3.93	3.56	14.94	4.61	11.33		
TS	4.10	5.54	7.52	12.61	5.06	4.00	5.83	4.51	6.40	5.56	5.77		
ACO_EoH	7.95	8.09	14.71	22.36	20.11	16.05	18.15	34.48	31.20	12.90	19.72		
ACO_MCTS	3.68	3.40	9.13	22.64	15.70	10.90	17.80	35.53	29.34	12.82	18.77		
GLS_EoH	0.67	0.63	1.62	2.67	2.69	3.89	3.99	6.56	4.43	5.23	<u>5.17</u>		
GLS_ReEvo	0.79	0.68	1.71	2.72	2.60	3.72	4.00	6.96	2.45	5.61	5.62		
GLS_MCTS	0.75	0.64	1.53	2.93	3.07	3.97	4.79	6.89	4.23	5.02	6.22		
STHS	0.87	0.60	1.47	3.59	3.48	3.88	4.41	3.41	6.80	5.64	5.36		
MTHS	0.72	0.49	1.03	2.64	1.08	1.50	0.94	1.23	3.51	1.06	4.29		

#### 3.2 METHODS AND SETTINGS

We compare our method, MTHS, against a diverse set of baselines representing the state of the art in both conventional and LLM-assisted heuristic design.

- Conventional Heuristics: We include widely-used constructive: Nearest Neighbor (NN) (Rosenkrantz et al., 1977) and a standard Insertion heuristic (Insert) (Rosenkrantz et al., 1977) and metaheuristic: Tabu Search (TS) (Glover & Laguna, 1998), Adaptive Large Neighborhood Search (ALNS) (Pisinger & Ropke, 2018), Memetic Search (MS) (Neri et al., 2011), and Guided Local Search (GLS) (Voudouris et al., 2010). For FSSP, we investgate GUPTA (Gupta, 1971), CDS (Campbell et al., 1970), NEH (Nawaz et al., 1983) and NEHFF (Fernandez-Viagas & Framinan, 2014), where NEH (Nawaz et al., 1983) and NEHFF (Fernandez-Viagas & Framinan, 2014) are widely recognized heuristics for this problem.
- Google OR-Tools: A high-performance, unified solver for CO problems. We utilize its standard metaheuristic solvers: Guided Local Searach (OR-Tools GLS), Simulated Annealing (OR-Tools SA) (Van Laarhoven & Aarts, 1987), and Tabu Search (OR-Tools TS) with their default parameter configurations.
- **LLM-driven Methods:** We compare against three recent LLM-based AHD methods: **EoH** (Liu et al., 2024b), **ReEvo** (Ye et al., 2024), and **MCTS-AHD** (Zheng et al., 2025). As these methods operate on a base heuristic framework, we test them with Ant Colony Optimization (ACO) and GLS, consistent with their original papers.
- MTHS (Ours): We evaluate our method in two configurations: MTHS (Multi-Task): The full proposed method, and STHS (Single-Task): An ablation where knowledge transfer and Pareto-based population management are disabled to assess the single-task performance of our hierarchical search.

**Experimental Setup for LLM-driven AHD** For MTHS, we conduct AHD on three tasks (i.e., TSP, CVRP and FSSP) with a budget of 1,000 program evaluations (i.e.,  $N_{eval}=1,000$ ). The high-level population size is  $N_H=8$  and the low-level search budget is  $N_L=4$ . For STHS and all compared LLM-driven AHD methods, including EoH, ReEvo, and MCTS-AHD, we conduct one search run per task with a budget of 1,000 program evaluations with their default settings. To prevent excessively long evaluations from stalling the search process, we impose a 20-minute time limit on

Table 2: Results on benchmark FSSP instances. The average gap (%) to the upper bounds from Taillard's FSSP benchmarks (Taillard, 1993), calculated over the 10 instances in each problem set. A set with n jobs and m machines is denoted as  $n_{-}m$ . The best result in each row is shown in **bold**, and the second-best is <u>underlined</u>.

	20_5	20_10	20_20	50_5	50_10	50_20	100_5	100_10	100_20	Average
GUPTA	12.89	23.42	21.79	12.23	20.11	22.78	5.98	15.03	21.00	17.25
CDS	9.03	12.87	10.35	6.98	12.72	15.03	5.10	9.36	13.55	10.55
NEH	3.24	4.05	3.06	0.57	3.47	5.48	0.39	2.07	3.58	2.88
NEHFF	2.30	4.15	2.72	0.40	3.62	5.10	0.31	1.88	3.73	2.69
LS	1.91	2.77	2.60	0.32	3.33	4.67	0.28	1.38	3.51	2.31
ILS	0.18	0.59	0.45	0.03	1.27	1.99	-0.03	0.34	1.29	0.68
MTHS	-0.01	0.03	0.03	0.00	0.22	0.45	<u>-0.02</u>	0.52	0.98	0.24

each individual heuristic evaluation. We use GPT-5-mini as the underlying LLM for all LLM-driven methods, except for ACO framework, where we directly adopt the best heuristics reported by Zheng et al. (2025) for EoH, ReEvo, and MCTS-AHD, rather than re-running the search. A summary of settings and running times is listed in Appendix D.

**Implementation and Execution Environment** All heuristic algorithms were implemented in Python, with the exception of Google OR-Tools, which uses a C++ library with a Python interface. Following standard practice in LLM-driven AHD research, we use the Numba JIT compiler to accelerate computationally intensive components, such as local search operators, for all metaheuristics, including those designed by the LLM-based methods.

Establishing a perfectly fair comparison based on a fixed evaluation budget is challenging due to the diverse frameworks and iterative components of different metaheuristics. Therefore, we adopted a time-based comparison protocol. We carefully configured the parameters of all baseline methods to commonly accepted values, ensuring that all algorithms had a comparable average wall-clock time for their execution. Detailed parameter settings are provided in the Appendix D.

The LLM-driven AHD experiments were conducted on a workstation equipped with two Intel Xeon 6248R CPUs and 128 GB of RAM. A single multi-task AHD using MTHS, utilizing 8-core parallel evaluations, took approximately 1.5 days to complete. The final heuristic evaluations were performed on a machine with an Intel Core Ultra 7 CPU and 32 GB of RAM.

## 3.3 MAIN RESULTS

We present our main experimental results in Table 1 and Table 2. A key contribution of our work is the ability of MTHS to discover a general-purpose metaheuristic. To demonstrate this, we selected a single metaheuristic from the final MTHS population and applied its associated three programs to the three distinct problem domains. This approach highlights the task-agnostic nature and generalization capabilities of the designed metaheuristic across different tasks. In contrast, existing LLM-driven AHD frameworks, including EoH, ReEvo, and MCTS-AHD, must execute a separate search to design a specialized heuristic for each task. The performance metric is the percentage gap relative to best-known solutions (for TSP and CVRP) or established upper bounds (for FSSP), with lower values signifying superior performance. For clarity, the best-performing heuristic is marked **in bold**, while the second-best is underlined.

Table 1 summarizes the results on the TSPLib and CVRPLib benchmarks. Our proposed method, MTHS, demonstrates superior performance, consistently outperforming all baseline methods across nearly all instance sets. For the TSP, MTHS achieves the lowest average optimality gaps on three size categories, from smaller instances to the largest ones (500-1000 nodes). The performance advantage of MTHS is consistent on the CVRP benchmarks. It secures the best results on four out of the seven CVRPLib sets (E, F, P, X) and is highly competitive on the remaining three. In contrast, while highly optimized solvers like Google OR-Tools perform well, especially on smaller CVRP instances, their performance degrades on larger TSP instances compared to the best LLM-evolved heuristics. Furthermore, comparing MTHS to its single-task ablation, STHS, reveals the clear benefit of multi-

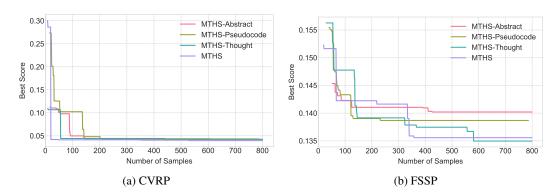


Figure 2: A comparison of different metaheuristic representations on CVRP and FSSP.

task learning; MTHS consistently outperforms STHS, underscoring the effectiveness of knowledge transfer in discovering more robust and powerful heuristics.

Table 2 shows the results on the Taillard benchmark for FSSP. The heuristic discovered by MTHS establishes a new state of the art, substantially outperforming all conventional constructive heuristics and local search methods. It achieves an average gap of just 0.24%. On several instance sets (20\_5 and 100\_5), the MTHS-designed heuristic finds solutions that are slightly better than or close to the existing upper bounds provided in Taillard (1993).

#### 3.4 METAHEURISTIC REPRESENTATION

We now analyze the representation MTHS uses to design metaheuristics, a key part of its success. The representation determines the LLM's level of abstraction, which in turn affects search efficiency and the quality of the resulting algorithms.

We compare four distinct metaheuristic representation strategies. Examples of different metaheuristic representations are provided in Appendix E.

- **Abstract:** The LLM is prompted to design a task-agnostic code structure directly, without a predefined template. This offers maximum flexibility but minimal structural guidance.
- Pseudocode: The LLM is prompted to design a task-agnostic pseudocode, which is then translated into executable code.
- Thought: The LLM describes the high-level strategy or "thought proces" of a metaheuristic.
- MTHS (Template): Our proposed method, which uses a structured, task-agnostic template to define the metaheuristic's components and control flow.

We use different metaheuristic representation in MTHS and perform the cross-task AHD on the three tasks with the same settings. Figure 2 illustrates the convergence behaviour of the automated search process for each representation on the CVRP and FSSP. It depicts the current best score (related gap to baseline on training instances) with respect to the number of program samples. The results clearly demonstrate the superiority of the template-based metaheuristic representation used in MTHS. For both problems, MTHS achieves a faster convergence compared to the other representations. This suggests that providing the LLM with a well-defined, modular structure is helpful for efficiently navigating the vast search space of possible metaheuristics.

## 3.5 GENERALIZATION TO NEW TASKS AND LLMS

A key hypothesis of our work is that a well-designed, task-agnostic metaheuristic can serve as a powerful and generalizable scaffold for solving new problems. To test this, we evaluate the generalization of the metaheuristic discovered by MTHS to an unseen task and across different LLMs.

Specifically, we prompt LLMs to generate code for solving BPP without any evolution (i.e., repeated sampling). We evaluate three models, including GPT-5-mini, Gemini-2.5-pro, and Claude-3.7-Sonnet, under two conditions: i) the model writes a solver from scratch, and ii) the model is

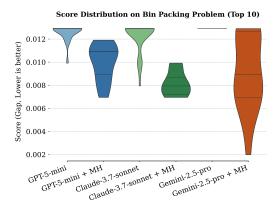


Figure 3: A comparison of results on BPP instances. + MH represents that we inform LLM to create programs using a given metaheuristic automatically designed by MTHS.

Table 3: Results on two sets of BPP instances. The results of three baseline methods are from (Zheng et al., 2025). + MH represents that we inform LLM to create programs using a given metaheuristic automatically designed by MTHS.

Method	n500, c150	n1000, c150
ЕоН	0.75%	0.85%
ReEvo	1.76%	2.06%
MCTS-AHD	0.48%	0.53%
GPT-5-mini	0.98%	0.98%
GPT-5-mini + MH	0.82%	0.65%
Gemini-2.5-pro	1.32%	1.25%
Gemini-2.5-pro + MH	0.34%	0.25%

explicitly instructed to implement a solver based on the metaheuristic template designed by MTHS (+MH). For each model and condition, we generate 100 programs and evaluate their performance on five BPP instances.

Figure 3 presents the performance distribution of the top 10 programs from each setting. The results show a notable and consistent improvement when the LLMs are guided by the MTHS-designed metaheuristic. For all three models, the + MH setting yields programs with significantly lower optimality gaps. Notably, Gemini-2.5-pro, when guided by the metaheuristic, produced a solver achieving a near-optimal gap of 0.002%, while when no metaheuristic is given, it struggled in designing high-quality BPP solvers. Results demonstrate that the task-agnostic metaheuristic designed by MTHS provides a general problem-solving logic that effectively transfers to new related tasks and can be leveraged by different LLMs.

To further validate the effectiveness of this generalization, we evaluate the top-performing LLM-generated solvers against established baselines on two BPP test sets with 500 and 1000 items. Each set contains 64 instances, and the average gap to the lower bound is reported. The solver generated by Gemini-2.5-pro with our metaheuristic guidance (+ MH) achieves state-of-the-art performance, recording optimality gaps of just 0.34% and 0.25% on the n=500 and n=1000 instances, respectively. This significantly outperforms not only the scratch-generated LLM solvers but also existing task-specific approaches like MCTS-AHD (0.48% and 0.53%). This demonstrates that the MTHS-discovered metaheuristic can generalized to other related tasks and enables LLMs to create programs that are not only conceptually sound but also highly competitive.

## 4 Conclusion

This paper addresses the limited cross-task generalization of current task-specific LLM-dirven AHD. We introduced Multi Task Hierarchical Search (MTHS), a framework that shifts the focus from crafting monolithic solvers to co-designing task-agnostic metaheuristics together with their task-specific realizations. Through a hierarchical representation and evolution, the method creates high-level metaheuristics that are reusable across tasks. Experiments on four problems show that metaheuristics produced by our approach outperform strong classical baselines, specialized metaheuristic solvers, and existing LLM-driven AHD methods. More importantly, the learned metaheuristic exhibits strong out-of-distribution behaviour. Used as a template on an unseen BPP, it enabled different LLMs to instantiate high-quality solvers without iterative search. These results indicate that designing at the metaheuristic level within a hierarchical representation offers a viable path to cross-task generalization in LLM-driven automated algorithm design.

In future work, we plan to expand the task suite, refine transfer mechanisms, and incorporate resource and reliability constraints directly into the search process. A deeper analysis of when and why transfer succeeds could further amplify the benefits of this paradigm.

## REFERENCES

- David Applegate, Robert Bixby, Vašek Chvátal, and William Cook. Concorde tsp solver, 2006.
  - Federico Berto, Chuanbo Hua, Nayeli Gast Zepeda, André Hottung, Niels Wouda, Leon Lan, Kevin Tierney, and Jinkyoo Park. Routefinder: Towards foundation models for vehicle routing problems. In *ICML 2024 Workshop on Foundation Models in the Wild*, 2024.
  - Leonardo CT Bezerra, Manuel López-Ibánez, and Thomas Stützle. Automatic component-wise design of multiobjective evolutionary algorithms. *IEEE Transactions on Evolutionary Computation*, 20(3):403–417, 2015.
  - Edmund K Burke, Matthew R Hyde, Graham Kendall, Gabriela Ochoa, Ender Özcan, and John R Woodward. A classification of hyper-heuristic approaches: revisited. In *Handbook of metaheuristics*, pp. 453–477. Springer, 2018.
  - Herbert G Campbell, Richard A Dudek, and Milton L Smith. A heuristic algorithm for the n job, m machine sequencing problem. *Management science*, 16(10):B–630, 1970.
  - 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 Ilms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 26931–26938, 2025.
  - Victor Fernandez-Viagas and Jose M Framinan. On insertion tie-breaking rules in heuristics for the permutation flowshop scheduling problem. *Computers & Operations Research*, 45:60–67, 2014.
  - Michel Gendreau, Jean-Yves Potvin, et al. *Handbook of metaheuristics*, volume 2. Springer, 2010.
  - Fred Glover and Manuel Laguna. Tabu search. handbook of combinatorial optimization. *Handbook of Combinatorial Optimization*, pp. 2093–2229, 1998.
  - Jatinder ND Gupta. A functional heuristic algorithm for the flowshop scheduling problem. *Journal of the Operational Research Society*, 22:39–47, 1971.
  - Can Gurkan, Narasimha Karthik Jwalapuram, Kevin Wang, Rudy Danda, Leif Rasmussen, John Chen, and Uri Wilensky. Lear: Llm-driven evolution of agent-based rules. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pp. 2309–2326, 2025.
  - Keld Helsgaun. An extension of the lin-kernighan-helsgaun tsp solver for constrained traveling salesman and vehicle routing problems. *Roskilde: Roskilde University*, 12:966–980, 2017.
  - Zhuoyi Lin, Yaoxin Wu, Bangjian Zhou, Zhiguang Cao, Wen Song, Yingqian Zhang, and Jayavelu Senthilnath. Cross-problem learning for solving vehicle routing problems. In *The 33rd International Joint Conference on Artificial Intelligence (IJCAI-24)*, 2024.
  - Fei Liu, Xi Lin, Zhenkun Wang, Qingfu Zhang, Tong Xialiang, and Mingxuan Yuan. Multi-task learning for routing problem with cross-problem zero-shot generalization. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1898–1908, 2024a.
  - 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 International Conference on Machine Learning*, pp. 32201–32223, 2024b.
  - Rafael Martí, Marc Sevaux, and Kenneth Sörensen. Fifty years of metaheuristics. *European Journal of Operational Research*, 321(2):345–362, 2025.
  - Muhammad Nawaz, E Emory Enscore Jr, and Inyong Ham. A heuristic algorithm for the m-machine, n-job flow-shop sequencing problem. *Omega*, 11(1):91–95, 1983.
  - Ferrante Neri, Carlos Cotta, and Pablo Moscato. *Handbook of memetic algorithms*, volume 379. Springer, 2011.

- Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025.
- Michael O'Neill and Conor Ryan. Grammatical evolution. *IEEE Transactions on Evolutionary Computation*, 5(4):349–358, 2001.
  - Nelishia Pillay and Rong Qu. Hyper-heuristics: theory and applications. Springer, 2018.
  - David Pisinger and Stefan Ropke. Large neighborhood search. In *Handbook of metaheuristics*, pp. 99–127. Springer, 2018.
  - Rong Qu, Graham Kendall, and Nelishia Pillay. The general combinatorial optimization problem: Towards automated algorithm design. *IEEE Computational Intelligence Magazine*, 15(2):14–23, 2020.
  - Gerhard Reinelt. Tsplib—a traveling salesman problem library. *ORSA Journal on Computing*, 3(4): 376–384, 1991.
    - Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
  - Daniel J Rosenkrantz, Richard E Stearns, and Philip M Lewis, II. An analysis of several heuristics for the traveling salesman problem. *SIAM journal on computing*, 6(3):563–581, 1977.
    - Yiding Shi, Jianan Zhou, Wen Song, Jieyi Bi, Yaoxin Wu, and Jie Zhang. Generalizable heuristic generation through large language models with meta-optimization. *arXiv* preprint *arXiv*:2505.20881, 2025.
    - Thomas Stützle and Manuel López-Ibáñez. Automated design of metaheuristic algorithms. In *Handbook of metaheuristics*, pp. 541–579. Springer, 2018.
    - Eric Taillard. Benchmarks for basic scheduling problems. *european journal of operational research*, 64(2):278–285, 1993.
    - Eduardo Uchoa, Diego Pecin, Artur Pessoa, Marcus Poggi, Thibaut Vidal, and Anand Subramanian. New benchmark instances for the capacitated vehicle routing problem. *European Journal of Operational Research*, 257(3):845–858, 2017.
    - Peter JM Van Laarhoven and Emile HL Aarts. Simulated annealing. In *Simulated annealing: Theory and applications*, pp. 7–15. Springer, 1987.
    - Niki van Stein and Thomas Bäck. Llamea: A large language model evolutionary algorithm for automatically generating metaheuristics. *IEEE Transactions on Evolutionary Computation*, 2024.
    - Christos Voudouris, Edward PK Tsang, and Abdullah Alsheddy. Guided local search. In *Handbook of metaheuristics*, pp. 321–361. Springer, 2010.
    - Lindong Xie, Yang Zhang, Zhixian Tang, Edward Chung, Genghui Li, and Zhenkun Wang. Coevolution of large language models and configuration strategies to enhance surrogate-assisted evolutionary algorithm. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 3321–3332, 2025a.
    - Zhuoliang Xie, Fei Liu, Zhenkun Wang, and Qingfu Zhang. Llm-driven neighborhood search for efficient heuristic design. In *Proceedings of the IEEE Congress on Evolutionary Computation*, pp. 1–8. IEEE, 2025b.
    - Shunyu Yao, Fei Liu, Xi Lin, Zhichao Lu, Zhenkun Wang, and Qingfu Zhang. Multi-objective evolution of heuristic using large language model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 27144–27152, 2025.

Yiming Yao, Fei Liu, Ji Cheng, and Qingfu Zhang. Evolve cost-aware acquisition functions using large language models. In *Proceedings of the International Conference on Parallel Problem Solving from Nature*, pp. 374–390. Springer, 2024.

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. *Advances in Neural Information Processing Systems*, 37:43571–43608, 2024.

Huigen Ye, Hua Xu, An Yan, and Yaoyang Cheng. Large language model-driven large neighborhood search for large-scale milp problems. In *Proceedings of the International Conference on Machine Learning*, 2025.

Rui Zhang, Fei Liu, Xi Lin, Zhenkun Wang, Zhichao Lu, and Qingfu Zhang. Understanding the importance of evolutionary search in automated heuristic design with large language models. In International Conference on Parallel Problem Solving from Nature, pp. 185–202. Springer, 2024.

Zhi Zheng, Zhuoliang Xie, Zhenkun Wang, and Bryan Hooi. Monte carlo tree search for comprehensive exploration in LLM-based automatic heuristic design. In *Proceedings of the International Conference on Machine Learning*, 2025.

## A RELATED WORKS

## A.1 AUTOMATED HEURISTIC DESIGN (AHD)

Automated Heuristic Design (AHD), often discussed under the umbrella of hyper-heuristics (Burke et al., 2018; Stützle & López-Ibáñez, 2018), aims to automate the process of selecting, combining, or generating simpler heuristics to solve complex computational search problems (Pillay & Qu, 2018). AHD methods are broadly categorized into selection and generation approaches.

Genetic programming or grammatical evolution (O'Neill & Ryan, 2001) are commonly used in generating new algorithms from fundamental building blocks. Recent advances in this area include component-based frameworks that assemble novel algorithms by integrating diverse operators and algorithmic stages (Bezerra et al., 2015; Qu et al., 2020). While powerful, these approaches often rely on hand-crafted components and require significant domain-specific knowledge, which can limit their flexibility and ease of application.

#### A.2 LLM-DRIVEN AHD

The advent of LLMs has introduced a new paradigm for AHD. A prominent strategy employs an evolutionary framework where LLMs iteratively propose and refine algorithms (Zhang et al., 2024; van Stein & Bäck, 2024). For example, Evolution of Heuristics (EoH) (Liu et al., 2024b) evolves both high-level thoughts and executable code using distinct prompt strategies to guide the search. FunSearch (Romera-Paredes et al., 2024) uses a multi-island evolutionary approach with a focused prompt strategy for refinement, while ReEvo (Ye et al., 2024) integrates reflection mechanisms to provide LLMs with more structured guidance. Other search strategies, such as Monte Carlo Tree Search (MCTS) (Zheng et al., 2025) and neighborhood search (Xie et al., 2025b), have also been explored to steer the design process.

Despite their success, a common limitation of these methods is their focus on discovering a single heuristic optimized for average performance on a specific task. The heuristic and knowledge can hardly be generalized to solving other tasks.

#### A.3 MULTI-TASK LEARNING FOR AHD

In the adjacent field of neural combinatorial optimization, multi-task learning has emerged as a key strategy for improving cross-problem generalization (Liu et al., 2024a; Berto et al., 2024). Researchers have developed single neural solvers trained across multiple problem types. Others demonstrate that models pre-trained on one problem (e.g., TSP) can be efficiently fine-tuned for related tasks (e.g., VRPs) using techniques like LoRA (Lin et al., 2024). Moverover, recent work (Shi

et al., 2025) has explored using LLMs to extract symbolic features that enhance the generalization of a backbone neural solver. However, these neural solvers are often black-box models that lack interpretability and typically require large datasets and substantial computational resources for training.

#### A.4 HEURISTIC REPRESENTATION IN LLM-DRIVEN AHD

The representation of the heuristic itself is a critical design choice in LLM-driven AHD. A common approach, popularized by EoH (Liu et al., 2024b), is a dual "thought-and-code" representation, where a high-level idea guides the generation of executable code. This or single code-based representations have been adopted by many subsequent works (Ye et al., 2024; Zheng et al., 2025; van Stein & Bäck, 2024). Recent explorations have introduced intermediate representations like pseudocode (Gurkan et al., 2025) or more flexible code structures (Novikov et al., 2025).

Closer to our work, some methods have used high-level algorithmic templates to enable meta-learning across different distributions of the same problem (Shi et al., 2025). However, to our knowledge, the challenge of learning generalizable metaheuristic structures that can be applied across entirely different tasks has not yet been addressed. Our hierarchical representation is designed specifically to fill this gap.

## 

## B MORE METHOD DETAILS

## 

## B.1 HIGH-LEVEL POPULATION INITIALIZATION

The InitializePopulation procedure (Algorithm 2) is responsible for seeding the initial high-level population,  $\mathcal{P}_H$ , which serves as the starting point for the main evolutionary search. The goal is to generate a diverse and competent set of initial individuals, where each individual represents a complete multi-task problem-solving strategy.

The procedure begins by constructing a single, comprehensive prompt using the BuildInitialPrompt function. This prompt aggregates the descriptions,  $\{D_t\}_{t=1}^m$ , of all m tasks in the set  $\mathcal{T}$ . This contextual information guides the LLMs ( $\mathcal{L}$ ) to generate a set of initial metaheuristics, denoted as InitialMHs. Each metaheuristic,  $MH_{init}$ , is a high-level textual description of a problem-solving approach.

For each generated  $MH_{init}$ , the procedure invokes LowLevelEvolution (as defined in the main MTHS algorithm). This critical step translates the abstract metaheuristic into a concrete, executable individual,  $I_{new}$ . The LowLevelEvolution procedure instantiates the metaheuristic into task-specific programs, refines them, and evaluates their performance, consuming a low-level evaluation budget of  $N_L$ . The resulting individual,  $I_{new}$ , contains a collection of optimized programs  $\{X_{new,t}^*\}_{t=1}^m$  and their corresponding scores  $\{S_{new,t}^*\}_{t=1}^m$ .

If the newly created individual  $I_{new}$  is deemed valid (e.g., it compiles and runs without fatal errors), it is added to the high-level population  $\mathcal{P}_H$ . This process repeats until the population reaches its target size,  $|\mathcal{P}_H|$ . The final, fully populated  $\mathcal{P}_H$  is then returned, ready for the main evolutionary loop of the MTHS algorithm.

## 

## B.2 Low-Level Key Function Evolution

The EvolveKeyFunction procedure (Algorithm 3) implements a fine-grained, task-specific optimization process. It is a core component of the LowLevelEvolution routine and is responsible for refining a single program by iteratively improving its most critical component: the key function. Its inputs are the initial program code  $X_t$  for a task  $T_t$ , the identified key function  $F_t$  within that code, the task object  $T_t$  (which provides the description  $D_t$  and evaluator  $E_t(\cdot)$ ), the LLM  $\mathcal{L}$ , and the low-level evaluation budget  $N_L$ .

The procedure operates as a micro-evolutionary search. It first initializes a local, low-level population,  $\mathcal{P}_{L,t}$ , by seeding it with the initial program  $X_t$ , its key function  $F_t$ , and its evaluated score. The main loop then commences, running until the evaluation budget  $N_L$  is exhausted.

In each iteration, a parent program,  $p_{parent}$ , is selected from  $\mathcal{P}_{L,t}$  using a selection strategy (e.g., tournament selection). A mutation prompt is then constructed via BuildMutationPrompt, providing the LLM with the task description  $T_t.D$  and the body of the parent's key function,  $p_{parent}$ .function. The LLM acts as a sophisticated mutation operator, generating a new function body,  $F'_{body}$ , that represents a plausible variation of the original.

This new function body is integrated back into the parent's base code to create a new program candidate,  $X'_{new}$ . This candidate is then executed and evaluated using the task-specific evaluator  $T_t.E(\cdot)$ , yielding a new score,  $S'_{new}$ . The new program, its function, and its score are registered as a new member of the low-level population  $\mathcal{P}_{L,t}$ . After the loop terminates, the procedure identifies the best-performing program in  $\mathcal{P}_{L,t}$  and returns its optimized code,  $X^*_t$ , and final score,  $S^*_t$ .

#### B.3 KNOWLEDGE TRANSFER

The KnowledgeTransfer procedure (Algorithm 4) is designed to enhance the multi-task proficiency of a newly generated individual,  $I_{new}$ , before it is integrated into the main population. This is achieved by systematically attempting to adapt its successful solutions from one task to another, leveraging the inherent relationships between tasks. The procedure takes the new individual  $I_{new}$ , the set of all tasks  $\mathcal{T}$ , and the LLM  $\mathcal{L}$  as input.

The process operates through a series of pairwise comparisons across all tasks. It iterates through every possible source task,  $t_{src}$ , and target task,  $t_{tgt}$ , within the individual's repertoire. For each pair where  $t_{src} \neq t_{tgt}$ , the procedure attempts to transfer knowledge.

Specifically, it constructs a transfer-oriented prompt using BuildTransferPrompt. This prompt provides the LLM with the description of the target task  $(T_{tgt}.D)$ , the full program code of the successful solution for the source task  $(I_{new}.X_{new,t_{src}})$ , and the code template for the target task  $(T_{tgt}.Temp)$ . The LLM's objective is to synthesize this information and generate a new program,  $X'_{transfer}$ , that is a plausible adaptation of the source solution for the target context.

This newly generated program is immediately evaluated on the target task using its evaluator,  $T_{tgt}.E(\cdot)$ , to obtain a transfer score,  $S'_{transfer}$ . This score is then compared against the individual's existing score for the target task,  $I_{new}.S_{new,t_{tgt}}$ . If the transfer results in a performance improvement ( $S'_{transfer} > I_{new}.S_{new,t_{tgt}}$ ), the individual is updated: its program and score for the target task,  $t_{tgt}$ , are replaced with the superior transferred versions,  $X'_{transfer}$  and  $S'_{transfer}$ . After all possible transfers have been attempted, the potentially improved individual  $I_{new}$  is returned.

## Prompt for Metaheuristic Generation

You are an expert algorithm designer. Your task is to create one novel algorithm for the following tasks:

{tasks\_formatted}

Design and present the high-level task-agnostic pseudocode for your new algorithm refer to the following template.

```
ALGORITHM < Algorithm_Name >
/\star PURPOSE: Brief description of the algorithm's purpose \star/
INPUT: <Description of input parameters/data>
OUTPUT:  OUTPUT:  Construction of expected results/return values>
/* Initialization Phase */
Initialize necessary data structures, variables, or state
Set up initial conditions or constraints
/* Main Processing Loop (if applicable) */
WHILE termination criteria not satisfied DO
    Perform core algorithm operations
    Update algorithm state
    Evaluate progress or intermediate results
    Adjust parameters if needed
END WHILE
/* Post-Processing Phase (if applicable) */
Finalize results
Perform any cleanup or final transformations
RETURN output
```

- The pseudocode must describe the core strategy and logical flow of the algorithm at a conceptual level
- Crucially, avoid low-level task-specific implementation details. Do not include specific variable names, data structures, or numerical constants.
- Ensure the pseudocode has a consistent shape (10–20 lines).

Enclose the entire pseudocode block within a single code block marked by '' pseudocode and ''.

810 Prompt for Program Generation 811 812 You are an expert algorithm implementer. Given a pseudocode algorithm, convert it to an efficient 813 Python implementation. 814 PSEUDOCODE: 815 {pseudocode} 816 817 IMPLEMENTATION REQUIREMENTS: 818 1. Use the template structure provided below 2. Ensure the implementation runs in acceptable time complexity 819 3. Maintain the core logic of the pseudocode 820 4. Use appropriate Python data structures and libraries 821 822 TEMPLATE: 823 {template\_program\_str} 824 RESPONSE FORMAT: 825 Return ONLY the Python code without explanations or examples, enclosed between "'python and " 826 markers as shown: 827 "python 828 # Your program here 829 830 831 832 833 Prompt for Key Function Identification 834 835 You are given a program. Please identify the most important function in this program that would benefit most from optimization. 836 Program: 837 {program\_str} 838 839 Task Description: 840 {task\_description} 841 Return only the key function. It should be enclosed between "'python and" markers exactly as shown 842 below: 843 "python 844 # Your key function here 845 846 847 848 849 Prompt for Key Function Generation 850 851 You are given a function. Please create a variation of this function that with the same inputs and 852 outputs but might be more effective or use a different approach. The function is part of a larger program solving the following task: 853 854 Task Description: 855 {task\_description} 856

Original function body:

{original\_function}

857

858 859

861

862

863

Return only the modified function. It should be enclosed between "python and" markers exactly as shown below:

"python # Your key function here

## Prompt for Knowledge Transfer

You are an expert algorithm designer specialized in translating knowledge between related problems. Your task is to implement an algorithm for the specified target problem by drawing inspiration from a reference algorithm that solves a related task.

## Target Problem:

{Task Description}

## Reference Algorithm:

The following is a high-quality implementation for a related problem that can inform your approach: {refer\_program}

## Implementation Template:

Your implementation should follow this template structure:

{template\_program\_str}

Return ONLY the Python code without explanations or examples, enclosed between "python and" markers as shown:
"python
"X

# Your program here

## Algorithm 2 Initialization of High-Level Population

```
Input:
```

```
    T = {T<sub>1</sub>,..., T<sub>m</sub>}: Set of tasks, each with description D<sub>t</sub>
    Large Language Model
    N<sub>H</sub>: Target size for the high-level population
    N<sub>L</sub>: Low-level evaluation budget per individual
    Output: Initialized high-level population P<sub>H</sub>
    procedure InitializePopulation(T, L, N<sub>H</sub>, N<sub>L</sub>)
    P<sub>H</sub> ← ∅
```

7: prompt  $\leftarrow$  BuildInitialPrompt( $\{D_t\}_{t=1}^m$ )
8: InitialMHs  $\leftarrow$   $\mathcal{L}(\text{prompt})$   $\triangleright$  Generate a set of initial metaheuristics
9: **for** each  $MH_{init}$  in InitialMHs **do**10: **if**  $N_H \geq N_H$  **then break**11:  $I_{new} \leftarrow \text{LowLevelEvolution}(MH_{init}, \mathcal{T}, N_L, \mathcal{L})$   $\triangleright$  Use main evaluation procedure

12: **if**  $I_{new}$  is valid **then**13:  $\mathcal{P}_H \leftarrow \mathcal{P}_H \cup \{I_{new}\}$ 14: **return**  $\mathcal{P}_H$ 

## C PROBLEM AND EXPERIMENTAL SETTINGS

## C.1 TRAVELING SALESMAN PROBLEM (TSP)

**Problem Definition:** Let G = (V, E) be a complete graph where:  $V = \{v_1, \dots, v_n\}$  represents n cities with coordinates  $\mathbf{x}_i \in [0, 1]^2$  and E contains edges with costs  $c_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$  (we consider Euclidean distance).

The objective is to find a Hamiltonian cycle  $\pi = (\pi_1, \dots, \pi_n, \pi_1)$  minimizing:

$$\mathcal{L}_{TSP} = \sum_{k=1}^{n-1} c_{\pi_k \pi_{k+1}} + c_{\pi_n \pi_1}.$$

**Task Description and Template: TSP Task Description:** Develop an algorithm to address the Traveling Salesman Problem. The objective is to determine the shortest route that visits each city in a given list exactly once and then returns to the starting city, thereby minimizing the total distance traveled.

945

946

947

948

949

950

951

952

953

954

955

956

957958959

960

961

962

963 964

965

966 967

968 969

970

971

#### 918 Algorithm 3 Low-Level Key Function Evolution 919 920 1: $X_t$ : Program code for task t921 2: $F_t$ : Identified key function for task t922 3: $T_t$ : Task object, containing description $D_t$ and evaluator $E_t(\cdot)$ 923 4: L: Large Language Model 924 5: $N_L$ : Low-level evaluation budget (can be used to limit iterations) 925 **Output:** Optimized program code $X_t^*$ and its score $S_t^*$ 6: **procedure** EVOLVEKEYFUNCTION( $X_t, F_t, T_t, \mathcal{L}, N_L$ ) 926 $\mathcal{P}_{L,t} \leftarrow \text{Initialize with } (X_t, F_t, E_t(X_t))$ ⊳ Seed low-level population 7: 927 $eval\_count\_L \leftarrow 1$ 8: 928 9: while $eval\_count\_L < N_L$ do 929 10: $p_{parent} \leftarrow \mathcal{P}_{L,t}.Selection()$ ▶ Select a program from the low-level pool 930 11: $prompt \leftarrow BuildMutationPrompt(T_t.D, p_{parent}.function)$ 931 12: $F'_{body} \leftarrow \mathcal{L}(\text{prompt})$ 932 $X'_{new} \leftarrow \text{IntegrateFunction}(p_{parent}.\text{code}, F'_{bodu}) \Rightarrow \text{Insert new function into base code}$ 13: 933 $S'_{new} \leftarrow T_t.E(X'_{new})$ 14: 934 $eval\_count\_L \leftarrow eval\_count\_L + 1$ 15: 935 Register new program $(X'_{new}, F'_{body}, S'_{new})$ in $\mathcal{P}_{L,t}$ 16: 936 $(X_t^*, S_t^*) \leftarrow \text{GetBest}(\mathcal{P}_{L,t})$ ▶ Get code and score of the best program 17: 937 return $(X_t^*, S_t^*)$ 18: 938 939 940 **Algorithm 4** Knowledge Transfer 941 **Input:** 942 1: $I_{new}$ : A new high-level individual with programs $\{X_{new,t}\}_{t=1}^m$ and scores $\{S_{new,t}\}_{t=1}^m$ 943

```
2: \mathcal{T} = \{T_1, \dots, T_m\}: Set of tasks
 3: L: Large Language Model
Output: Updated individual I_{new}
 4: procedure KNOWLEDGETRANSFER(I_{new}, \mathcal{T}, \mathcal{L})
 5:
          for t_{src} \leftarrow 1 to m do
 6:
               for t_{tqt} \leftarrow 1 to m do
 7:
                    if t_{tqt} = t_{src} then continue
                    \mathsf{prompt} \leftarrow \mathsf{BuildTransferPrompt}(T_{tgt}.D, I_{new}.X_{new,t_{src}}, T_{tgt}.Temp)
 8:
 9:
                    X'_{transfer} \leftarrow \mathcal{L}(prompt)

    Adapt source solution to target task

                    S'_{transfer} \leftarrow T_{tgt}.E(X'_{transfer})
10:
                    if S'_{transfer} < I_{new}.S_{new,t_{tqt}} then
11:
                        I_{new}.S_{new,t_{tqt}} \leftarrow S'_{transfer}
12:
                                                                               ▶ Update score if transfer is successful
                         I_{new}.X_{new,t_{tqt}} \leftarrow X'_{transfer}
                                                                                                   ▶ Update program code
13:
14:
          return I_{new}
```

**Training Instances:** For the heuristic evolution process, we use a set of 64 TSP instances, each with 100 locations randomly sampled from a uniform distribution over  $[0,1]^2$ . The fitness of a candidate heuristic is measured by its average optimality gap, calculated against the optimal solutions found by the Concorde solver.

**Testing Instances:** For testing, we select commonly used 49 symmetric Euclidean TSPLib instances (Reinelt, 1991), with problem sizes ranging from 52 to 1,000 nodes.

## C.2 CAPACITATED VEHICLE ROUTING PROBLEM (CVRP)

**Problem Definition:** CVRP aims to minimize the total traveling distances of a fleet of vehicles given a depot and a set of customers with coordinates and demands. Given: 1) Depot  $v_0$  and customers  $\{v_1,...,v_n\}$  with coordinates  $\mathbf{x}_i \in [0,1]^2$ , 2) Demands  $d_i \in \mathbb{Z}^+$  ( $d_0 = 0$ ), 3) Vehicle capacity  $Q \in \mathbb{Z}^+$ , 4) Distance metric  $c_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$ , find routes  $\mathcal{R} = \{r_1,...,r_m\}$ . Each route  $r_k$ 

```
972
         Template for TSP
973
974
975
        1 import numpy as np
976
        3 class TSPSolver:
977
             def __init__(self, coordinates: np.ndarray, distance_matrix: np.
978
                 ndarray):
979
980
        6 Initialize the TSP solver.
981
        8 Aras:
982
             coordinates: Numpy array of shape (n, 2) containing the (x, y)
983
                 coordinates of each city.
984
             distance_matrix: Numpy array of shape (n, n) containing pairwise
985
                 distances between cities.
       11 """
986
       12 self.coordinates = coordinates
987
       13 self.distance_matrix = distance_matrix
988
989
              \# --- your code here ---
       15
990
       16
991
       17
             def solve(self) -> np.ndarray:
       18 """
992
       19 Solve the Traveling Salesman Problem (TSP).
993
       20
994
       21 Returns:
995
       22
             A numpy array of shape (n,) containing a permutation of integers
             [0, 1, \ldots, n-1] representing the order in which the cities are
996
       23
                 visited.
997
       24
998
              The tour must:
999
              - Start and end at the same city (implicitly, since it's a loop)
       26
1000
              - Visit each city exactly once
       27
       28 """
1001
       29 n = len(self.coordinates)
1002
1003
       31 \# --- your code here ---
1004
       32
1005
       33 \ # Example (naive ordered tour replace with your algorithm):
       34 tour = np.arange(n)
1007
       36 return tour
1008
1009
```

starts/ends at  $v_0$ . Capacity constraints are satisfied  $\sum_{v_i \in r_k} d_i \leq Q$  and all customers served exactly once. The objective is to minimize total distance.

**Task Description and Template: CVRP Task Description:** Develop an algorithm to solve the Capacitated Vehicle Routing Problem (CVRP). The objective is to determine the optimal set of routes for a fleet of vehicles that all start and end at a central depot. Each vehicle has a maximum capacity, and the routes must collectively serve all customer nodes exactly once without exceeding the vehicle's capacity. The goal is to minimize the total distance traveled across all routes.

**Training Instances:** The heuristic evolution is conducted on 64 randomly generated CVRP instances, each with 100 customers. Customer and depot locations are randomly sampled from  $[0,1]^2$ . Each vehicle has a capacity of 50, and customer demands are integers sampled uniformly from  $\{1,\ldots,9\}$ . The fitness value is the average gap to LKH solver (Helsgaun, 2017).

**Testing Instances:** We select 7 commonly used benchmark sets, including A, B, E, F, M, P, and X, from CVRPLib (Uchoa et al., 2017). The chosen sets and their characteristics are summarized in Table 4. Due to the time limit, we do not test on all instances from the X set.

Table 4: CVRPLib benchmark sets

Benchmark Set	Number of Instances	Instance Size
Set A	27	31-79
Set B	23	30-77
Set E	11	22-101
Set F	3	44-134
Set M	5	100-199
Set P	23	15-100
Set X	43	100-500

## C.3 FLOW SHOP SCHEDULING PROBLEM (FSSP)

**Problem Definition:** The Flow Shop Scheduling Problem (FSSP) aims to minimize the makespan (total time to complete all jobs) for a set of jobs that must be processed on a series of machines in a fixed order. Given: 1) A set of n jobs  $\mathcal{J}=\{J_1,...,J_n\}$ , 2) A set of m machines  $\mathcal{M}=\{M_1,...,M_m\}$ , 3) The processing time  $p_{ij}\in\mathbb{Z}^+$  for each job  $J_i$  on each machine  $M_j$ . The problem is to find a permutation (sequence)  $\pi$  of the jobs. This sequence dictates the order in which jobs are processed on the first machine, and this same order is maintained for all subsequent machines. The objective is to find the sequence  $\pi$  that minimizes the makespan,  $C_{max}(\pi)$ , which is the completion time of the last job on the last machine.

**Task Description and Template: FSSP Task Description:** Develop an algorithm to solve the Flow Shop Scheduling Problem (FSSP) by determining the optimal sequence of jobs to minimize makespan. In FSSP, all jobs must be processed on all machines in the same order (machine 0, then machine 1, then machine 2, etc.). The goal is to find the job sequence that minimizes the makespan (total completion time) while ensuring that: (1) all jobs follow the same machine processing order, (2) each machine processes only one job at a time, and (3) each job can only be processed on one machine at a time. The algorithm should return a permutation of job indices representing the order in which jobs should be processed.

**Training Instances:** For heuristic evolution, we use 64 randomly generated instances, each comprising 50 jobs and a number of machines varying between 2 and 20. The processing times for each job are sampled from a uniform distribution over  $[0,1]^2$ . The average makespan (gap to lower bound) is used as the fitness value.

**Testing Instances:** We evaluate the algorithms on the widely-used Taillard instances (Taillard, 1993). We test 9 different test sets. The number of jobs in these instances ranges from 20 to 100, and the number of machines ranges from 5 to 20.

## C.4 BIN PACKING PROBLEM (BPP)

**Problem Definition:** We consider one-dimensional bin packing problem. The primary goal is to pack a set of n items, each with a specific size or weight  $w_j \in \mathbb{Z}^+$ , into the minimum number of identical bins, each having a uniform capacity  $C \in \mathbb{Z}^+$ . The core challenge is to find a partition of the items into a set of bins  $B = \{B_1, \ldots, B_m\}$  such that the sum of item sizes in any single bin does not exceed the capacity C, and the total number of bins used, m, is minimized.

**Task Description and Template: BPP Task Description:** You are given a set of items, each with a specific weight, and a number of identical bins, each with a fixed capacity. The goal is to pack all items into the minimum number of bins possible, such that the sum of the weights of the items in each bin does not exceed the bin's capacity.

**Instances:** Following the setup of Ye et al. (2024) and Zheng et al. (2025), we generate instances where bins have a capacity of 150 and item sizes are uniformly sampled from the range [20, 100].

## D More Details on Metaheuristics

Table 5 details the configuration and average running times of several metaheuristic solvers on standard TSPLib and CVRPLib benchmarks. We compare our Python implementations of Tabu Search, ALNS, and Memetic Search—accelerated using Numba, against Google's C++ OR-Tools solvers, an existing LLM-based AHD approach, and our proposed MTHS.

Table 5: Average Running Time and Configuration of Metaheuristic Solvers on Benchmark Datasets. Our Python-based implementations are compared against the highly optimized C++ solvers from Google OR-Tools, existing LLM-driven AHD approach and our proposed MTHS. Average running times are reported for standard TSPLib and CVRPLib instances.

Category	Metaheuristic	Key Parameters	Key Accelerated Functions (Numba/C++)	Average Running Time		
Cutegory	Wetaneur istic	ncy runnecers	recy receivated 1 unctions (runnba/C11)	TSPLib	CVRPLib	
Our Python	Tabu Search (TS)	max_iterations: 100 tabu_tenure: 20	_calculate_tour_distance_numba _find_best_neighbor_numba	58s	8s	
Implementations	Adaptive Large Neighborhood Search (ALNS)	max_iterations: 1000 removal_rate: [0.1, 0.4] reaction_factor: 0.5	_calculate_tour_cost _greedy_insertion _shaw_removal	155s	138s	
	Memetic Search (MS)	population_size: 30 generations: 50 tournament_size: 5 patience: 40	_calculate_tour_distance _two_opt_local_search _generate_nearest_neighbor_tour	190s	156s	
OR-Tools Solvers	Tabu Search (TS)  Is Solvers  Simulated Annealing (SA)  Guided Local Search (GLS)  Default  Default		C++ C++ C++	60s 60s 60s	60s 60s 60s	
Existing AHD Approach	GLS (EoH, ReEvo, MCTS-AHD)	iter_limit: 100 perturbation_moves: 30	_two_opt_once _relocate_once	60–100s	25-30s	
MTHS (Ours)		time_limit: 100 population_size: 10	_two_opt _insert _swap	100s	60s	

## E MORE EXPERIMENTAL RESULTS AND ANALYSES

## E.1 PARETO FRONT

Figure 4 illustrates the search trajectory and final population of our proposed method, MTHS, in the three-dimensional objective space for the multi-task AHD. The background points represent all candidate metaheuristics in the population throughout the evolutionary process, colored by their generation index from early (dark purple) to late (bright yellow). This visualization demonstrates the algorithm's progression, showing how it initially explores a broad region of the objective space before intensifying its search and converging towards the Pareto front. The final, non-dominated population is highlighted in red, showcasing a well-distributed set of high-quality trade-off solutions across the three conflicting objectives: TSP, CVRP, and FSSP. The distinct separation and advancement of the final front from the historical samples underscore the effectiveness of our approach in achieving both convergence and diversity.

## Search Trajectory & Final Population (MTHS)

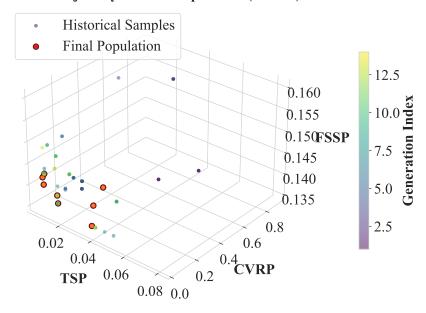


Figure 4: Metaheuristics generated by MTHS colored by generation and the final non-dominated front highlighted in red. The three objectives are fitness on three tasks (i.e., TSP, CVRP and FSSP). These metaheuristics that are removed during population management are not included.

## E.2 METAHEURISTIC REPRESENTATION

We identify and illustrate four distinct levels of abstraction for describing a metaheuristic algorithm: i) a high-level metaheuristic in MTHS, ii) an algorithmic pseudocode, iii) a code-level abstraction, and iv) a natural language thought description. The conceptual design outlines the overarching strategy, while pseudocode and code-level abstractions provide structured, implementation-oriented views. The thought description captures the core inventive idea in a dense, human-readable format.

For brevity and due to their structural similarity, we present a single example for the pseudocode and code-level abstraction formats. Each format is demonstrated below with a representative metaheuristic.

1188 1189 1190 Adaptive Cooperative Substructure Search (ACSS) 1191 **Purpose:** A unified, task-agnostic metaheuristic to find high-quality feasible solutions for routing and 1192 scheduling problems by combining constructive heuristics, cooperative memory of useful substruc-1193 tures, adaptive perturbation, and constraint-aware repair. 1194 **INPUT:** A problem instance with a solution representation, an objective evaluator, and a constraint 1195 checker **OUTPUT:** A feasible solution (permutation or set of routes) with a near-optimal objective value 1196 1197 ▷ Initialization Phase 1198 Construct a diverse set of initial candidate solutions using problem-aware constructive methods 1199 and randomization Extract and record promising substructures from initial candidates into a cooperative memory 1201 ▶ Main Processing Loop 1202 while stopping condition not met do 1203 Select one or more candidates for improvement based on quality and diversity Intensify: apply local improvement operators guided by cooperative memory to reduce objec-1205 tive while preserving feasibility Diversify: apply adaptive, constraint-aware perturbations to escape local optima and generate varied neighborhood proposals 1207 Repair: enforce feasibility by applying generic constraint-handling procedures that adapt to 1208 problem specifics 1209 **Recombine:** optionally merge complementary substructures from cooperative memory into 1210 candidates to create new high-quality solutions Evaluate updated candidates with objective evaluator and constraint checker 1211 Update cooperative memory with newly discovered high-quality substructures and adjust op-1212 erator selection probabilities based on recent success 1213 1214 *⊳ Post-Processing Phase* 1215 Polish the best feasible solution using targeted local refinement and a final constraint-aware repair 1216 if needed 1217 **RETURN** best feasible solution found 1218 1219 Metaheuristic as Code Abstract/Pseudocode 1220 1: **procedure** AMOCGS(problem, params) 1222 population ← multi\_construct(problem, params.heuristics) ▷ task-specific constructive seeds 3: population  $\leftarrow$  map(lambda s: repair\_and\_evaluate(s, problem), population) ▷ enforce 1224 constraints and score 4: operators ← init\_operator\_pool(problem) ▷ problem-aware neighborhood & crossover 1225 5: op\_scores ← init\_scores(operators); memory ← init\_elite\_memory(population) 1226 6:  $best \leftarrow argmin(population)$ 1227 7: **while** not termination\_condition(params) **do** 1228 8: parents ← select\_parents(population, op\_scores, params) ⊳ biased by quality and diversity 1229 9: op ← adaptive\_select(operators, op\_scores, params) 10: offspring  $\leftarrow$  apply\_operator(op, parents, problem) 1230 11: offspring ← local\_search\_and\_repair(offspring, problem, params) ⊳ e.g., tabu/SA/LNS 1231 respecting constraints 1232 12: offspring.score  $\leftarrow$  evaluate(offspring, problem) 1233 13: population ← replace\_population(population, offspring, params) ⊳ elite preservation + 1234 diversity maintenance 14: op\_scores \( \to \) update\_op\_scores(op\_scores, op, offspring, improvement\_metric(best, offspring)) 15:  $best \leftarrow select\_best(best, offspring)$ 1237 adapt\_parameters(params, op\_scores, population, memory) beta temperature, operator
 constant of temperature.
 const 16: weights, restart triggers 1239 17: if intensify\_trigger(params) then 18: path\_relink\_and\_intensify(population, memory, problem) 1240 19: return best 1241

## Metaheuristic as Thought Description

Recursively partition the instance into manageable clusters (spatial for routing, temporal for scheduling), stochastically generate diverse candidate partial sequences within each cluster using lightweight local cost models, and iteratively merge clusters with a capacity- and precedence-aware repair operator that enforces feasibility; concurrently adapt sampling biases via online learning of high-value move patterns to concentrate search.

During merges apply a multi-objective adaptive acceptance criterion that balances global cost reduction and constraint satisfaction, allowing focused local search and occasional exploratory perturbations to rapidly converge to high-quality feasible permutations and route sets.

## E.3 TOKENS AND COST

The computational cost of MTHS is dominated by LLM API calls, measured in latency and token usage. We analyze the cost to generate and evaluate a single new metaheuristic individual,  $I_{new}$ , across m tasks.

**High-Level Evolution** This step generates one new metaheuristic,  $MH_{new}$ , from k parents.

• LLM Calls: 1

• Token Cost: The input includes a prompt and k parent metaheuristics; the output is  $MH_{new}$ .

**Low-Level Evolution** This is the most expensive phase, executed for each of the m tasks to evaluate  $MH_{new}$ . The cost for a single task  $T_t$  includes:

• i) Program Generation: 1 LLM call to combine  $MH_{new}$  and a task template  $Temp_t$  into a program  $X_{new,t}$ .

• ii) **Key Function Identification:** 1 LLM call to analyze  $X_{new,t}$  and extract the key function  $F_{new,t}$ .

 • iii) Key Function Refinement: An evolutionary loop with a budget of  $N_L$  evaluations, where each step uses the LLM as a mutation operator. This requires  $N_L$  LLM calls.

The total cost for this stage scales linearly with the number of tasks (m) and the refinement budget  $(N_L)$ .

**Knowledge Transfer** After evaluation, this optional step adapts the best-performing program from a source task,  $X_{new\ src}^*$ , to the other m-1 target tasks.

• LLM Calls: m-1

 - Token Cost: Each call prompts the LLM with  $X^*_{new,src}$  and a target task template.

**Summary of Costs** The total number of LLM calls required to evaluate one new metaheuristic individual is:

$$Calls_{total} = \underbrace{1}_{\text{High-Level}} + \underbrace{m \times (2 + N_L)}_{\text{Low-Level}} + \underbrace{(m - 1)}_{\text{Knowledge Transfer}} \tag{1}$$

The dominant cost factor is the Low-Level Evolution, particularly the Key Function Refinement loop ( $m \times N_L$  calls), making it the primary bottleneck in terms of time and expense.

When compared to existing LLM-driven AHD methods that target a single task, evaluating one MTHS individual requires a larger number of LLM requests due to the per-task evaluations. However, because MTHS simultaneously designs heuristics for multiple tasks within a single evolutionary run, the total computational budget required to find effective heuristics for an entire set of tasks is lower than running a single-task AHD method independently for each task. Table 6 lists the tokens used for each components in one run of MTHS on three tasks. It costs around 10 dollars when using GPT-5-mini.

1296 1297

Table 6: Breakdown of token cost and number of LLM requests in MTHS.

1316 1317

1318 1319 1320 1321 1322

1334

1341 1342 1343

1340

1344 1345

1346 1347 1348 1349

Tokens/ Total Type Sub-type No. Percentage Sample **Tokens** High-level Metaheuristic generation 80 80k 1% 1k Evolution Lower-level Initial program 240 1440k 18% 6k **Evolution** generation Key function 240 4k 960k 12% identification Key function 0.6k720 432k 5% generation Program generation 6k 720 4320k 52% using new key function Knowledge Program generation with 6k 160 960k 12% Transfer knowledge transfer 8192k **Total**  $100\,\%$ 

Table 7: Detailed results for selected TSPLib instances (first seven instances in alphabetical order with different sizes and distributions): Gap Performance and Runtimes.

Mada al	a280 berlin52 bier127		r127	ch	ch130 ch1			0 d198		d493				
Method	Gap	Time	Gap	Time	Gap	Time	Gap	Time	Gap	Time	Gap	Time	Gap	Time
NN	27.43	0.23	19.08	0.01	14.86	0.19	23.98	0.19	25.53	0.21	19.33	0.21	24.02	0.34
Insert	20.13	0.02	4.55	0.00	12.02	0.00	6.31	0.00	9.45	0.01	12.48	0.01	24.48	0.06
Or-tools SA	4.39	60.13	4.80	60.04	2.25	60.18	1.73	60.39	1.67	60.11	1.23	60.01	3.29	60.68
Or-tools TS	4.67	60.23	0.03	60.04	1.43	60.18	1.73	60.42	1.67	60.14	1.23	60.25	3.58	60.17
Or-tools GLS	5.27	60.07	0.03	60.06	1.43	60.05	0.49	60.13	0.73	60.19	2.86	60.22	4.65	60.22
MS	6.24	35.30	0.03	1.98	1.03	7.08	3.53	6.59	2.85	3.64	2.24	13.78	4.76	409.47
ALNS	7.04	40.82	0.03	1.03	2.14	10.20	1.57	9.44	3.24	9.19	1.85	19.06	6.32	300.42
TS	11.32	124.54	3.42	70.16	6.37	83.09	4.23	87.29	5.81	92.50	3.73	116.51	9.46	113.10
ACO_EoH	27.77	39.56	1.79	4.59	N	I/A	10.39	22.45	3.97	31.06	12.03	23.46	18.52	862.68
ACO_MCTS	9.50	83.46	0.03	28.92	3.23	18.79	3.40	27.90	1.63	29.78	2.61	20.38	18.88	909.84
GLS_EoH	1.78	351.41	0.03	2.72	0.04	2.02	1.15	4.49	0.84	4.10	0.95	259.85	1.66	563.71
GLS_ReEvo	2.94	349.71	0.03	3.09	0.62	2.51	0.64	5.62	0.97	4.39	1.16	400.06	2.72	563.09
GLS_MCTS	3.22	340.82	0.03	1.43	0.59	2.71	0.32	4.18	0.97	4.21	1.06	265.00	1.78	559.99
STHS	2.07	49.46	0.03	37.28	0.39	42.05	1.05	41.56	0.45	38.89	0.59	43.50	2.16	57.56
MTHS	1.34	100.17	0.03	81.26	0.39	100.00	0.64	100.00	0.37	100.01	0.39	100.00	1.63	100.01

## DETAILED RESULTS ON BENCHMARK INSTANCES

#### F REPRODUCTION

We are committed to making our research fully reproducible and accessible to the broader community. We have made our code for metaheuristics and data publicly available. Our resources are hosted on an anonymous link https://anonymous.4open.science/r/MTHS-E80B.

The following components are provided:

Table 8: Detailed results for selected CVRPLib instances (middle-size X instances): Gap Performance and Runtimes.

	OR-T	ools TS	OR-T	ools SA	GLS	БеоН	GLS	ReEvo	GLS MCTS		MTHS	
Instance	Gap	Time	Gap	Time	Gap	Time	Gap	Time	Gap	Time	Gap	Time
X-n303-k21	7.5	60.0	6.6	60.0	3.7	43.3	5.6	33.1	8.9	49.2	4.7	60.0
X-n308-k13	8.9	60.0	10.4	60.0	5.4	28.9	5.9	29.7	8.4	47.2	6.8	60.0
X-n313-k71	8.1	60.0	9.3	60.0	7.8	72.5	10.7	57.6	10.9	89.5	3.4	60.0
X-n317-k53	1.3	60.0	1.3	60.0	1.5	39.9	1.4	40.5	1.4	70.2	1.4	60.0
X-n322-k28	8.1	60.0	8.4	60.0	4.7	42.8	10.4	33.3	7.8	81.4	5.1	60.0
X-n327-k20	7.6	60.0	7.4	60.0	3.4	38.3	5.9	31.8	6.5	51.7	6.2	60.0
X-n331-k15	7.6	60.0	6.4	60.0	4.6	34.0	5.5	28.8	5.0	46.8	5.5	60.0
X-n336-k84	4.0	60.0	4.1	60.0	4.8	93.0	4.4	71.1	4.8	95.2	3.8	60.0
X-n344-k43	5.1	60.0	5.1	60.0	6.3	45.8	6.2	42.6	7.3	65.3	4.7	60.0
X-n351-k40	9.7	60.0	9.1	60.0	6.0	65.6	8.4	52.6	9.2	73.8	4.4	60.0
X-n359-k29	7.1	60.0	6.9	60.0	4.2	64.9	4.8	42.0	5.7	64.3	3.0	60.0
X-n367-k17	10.0	60.0	6.8	60.0	10.0	86.5	9.5	108.6	8.6	182.0	10.6	60.0
X-n376-k94	0.7	60.0	0.7	60.0	0.8	106.6	0.8	114.7	0.8	156.9	1.0	60.0
X-n384-k52	5.6	60.0	5.3	60.0	4.9	135.4	5.7	101.4	5.0	160.9	3.8	60.0
X-n393-k38	8.6	60.0	8.2	60.0	9.1	108.5	7.8	111.5	8.7	164.1	4.4	60.0
X-n401-k29	3.7	60.0	3.7	60.0	3.2	155.4	5.3	162.6	3.7	233.3	2.5	60.0
X-n411-k19	13.4	60.0	13.4	60.0	9.3	155.3	9.1	139.8	10.0	218.7	12.9	60.0
X-n420-k130	6.4	60.0	6.9	60.2	4.9	180.1	4.7	159.5	5.3	221.2	5.0	60.0
X-n429-k61	5.4	60.0	5.8	60.0	5.5	141.8	5.7	128.2	8.0	184.1	4.1	60.0
X-n439-k37	4.7	60.0	4.9	60.1	3.0	140.2	2.9	120.5	3.3	180.6	5.2	60.0
X-n449-k29	11.0	60.0	10.4	60.0	7.4	145.9	7.6	156.8	8.2	201.4	4.3	60.0
X-n459-k26	12.6	60.0	10.6	60.0	12.8	161.8	10.4	168.9	10.4	207.4	7.9	60.0
X-n469-k138	4.2	60.0	4.5	60.0	6.9	154.9	7.2	158.1	8.0	180.1	6.1	60.0
X-n480-k70	4.1	60.0	4.0	60.0	4.7	139.2	5.3	137.5	5.6	169.4	3.7	60.0
X-n491-k59	10.4	60.0	8.5	60.0	8.3	257.5	7.1	232.8	8.6	193.4	4.0	60.0

- 1. **Detailed Experimental Results:** In the sections of this appendix, we present detailed tables and figures that elaborate on the results discussed in the main text. This includes per-instance performance and running times.
- Open-Sourced Algorithms: The core contribution of our work, the generated metaheuristics, is available in our public repository. The code is commented to facilitate understanding and extension.
- 3. **Open-Sourced Evaluation Datasets and Scripts:** To ensure fair and consistent comparison, we have released the complete set of evaluation datasets, including TSP, CVRP, FSSP and BPP, used in our experiments. The repository also contains the exact scripts used to run the evaluations.

## G USE OF LLMS

First, for manuscript preparation, the LLM was employed as a writing assistant to check grammar and refine phrasing, particularly in the introduction section. Second, the LLM was integrated as a core component of our proposed method to design and generate heuristics and programs.

```
1458
1459
         Template for CVRP
1460
1461
1462
       1 import numpy as np
       2 class CVRPSolver:
1463
             def __init__(self, coordinates: np.ndarray, distance_matrix: np.
1464
                 ndarray, demands: list, vehicle_capacity: int):
1465
1466
       5 Initialize the CVRP solver.
1467
       7 Args:
1468
             coordinates: Numpy array of shape (n, 2) containing the (x, y)
1469
                 coordinates of each node, including the depot.
1470
             distance_matrix: Numpy array of shape (n, n) containing pairwise
1471
                distances between nodes.
1472
             demands: List of integers representing the demand of each node (
                 first node is typically the depot with zero demand).
1473
             vehicle_capacity: Integer representing the maximum capacity of
       11
1474
                 each vehicle.
1475
       12 """
1476
       13 self.coordinates = coordinates
1477
       14 self.distance_matrix = distance_matrix
       15 self.demands = demands
1478
       16 self.vehicle_capacity = vehicle_capacity
1479
1480
             \# --- your code here ---
1481
1482
             def solve(self) -> list:
       20
1483
       22 Solve the Capacitated Vehicle Routing Problem (CVRP).
1484
       23
1485
       24 Returns:
1486
          A one-dimensional list of integers representing the sequence of
1487
                nodes visited by all vehicles.
             The depot (node 0) is used to separate different vehicle routes
1488
                and appears at the start and end
1489
            of each route. For example: [0, 1, 4, 0, 2, 3, 0] represents:
1490
       28 - Route 1: 0 - 1 - 4 - 0
1491
       29 - Route 2: 0 - 2 - 3 - 0
       30 """
1492
       31 n = len(self.coordinates)
1493
       32
1494
       33 \# --- your code here ---
1495
1496
       35 \# Example (naive solution replace with your algorithm):
       36 solution = [0] \# Start at the depot
1497
       37 current_capacity = 0
1498
1499
       39 for i in range(1, n):
1500
             if current_capacity + self.demands[i] > self.vehicle_capacity:
       40
1501
            solution.append(0) \# return to depot and start a new route
       41
1502
            current\_capacity = 0
1503
       43
             solution.append(i)
       44
1504
       45
             current_capacity += self.demands[i]
1505
1506
       47 if solution[-1] != 0:
1507
             solution.append(0) \# end the last route at the depot
       48
       49
1508
       50 return solution
1509
1510
```

```
1516
1517
1518
1519
1520
         Template for FSSP
1521
1522
1523
       1 import numpy as np
1524
       3 class FSSPSolver:
1525
             def __init__(self, num_jobs: int, num_machines: int,
1526
                 processing_times: list):
1527
1528
       6 Initialize the FSSP solver.
1529
       8 Args:
1530
            num_jobs: Number of jobs in the problem
1531
            num_machines: Number of machines in the problem
1532
            processing_times: List of lists where processing_times[j][m] is
1533
                 the processing time of job j on machine m
       12 """
1534
       13 self.num_jobs = num_jobs
1535
       14 self.num_machines = num_machines
1536
       15 self.processing_times = processing_times
1537
1538
             \# --- your code here ---
1539
       18
       19
            def solve(self) -> list:
1540
1541
       21 Solve the Flow Shop Scheduling Problem (FSSP).
1542
       22
1543
       23 Returns:
1544
           A list representing the sequence of jobs to be processed.
       24
             For example, [0, 2, 1] means job 0 is processed first, then job
1545
                 2, then job 1.
1546
             All jobs must be processed on all machines in the same order.
1547
       27
1548
             The sequence must include all jobs exactly once.
       29 """
1549
       30
1550
       31 \# --- your code here ---
1551
1552
       33 \# Simple solution: process jobs in their original order (0, 1, 2,
1553
1554
       34 job_sequence = list(range(self.num_jobs))
1555
       36 return job_sequence
1556
1557
1558
```

```
1567
1568
         Template for BPP
1569
1570
        1 import numpy as np
1571
1572
        3 class BPPSolver:
1573
             def __init__(self, capacity: int, weights: list[int | float]):
1574
1575
        6 Initialize the BPP solver.
1576
        8 Args:
1577
            capacity (int): The capacity of each bin.
1578
             weights (list[int | float]): A list of item weights.
1579
1580
       12 self.capacity = capacity
       13 self.weights = weights
1581
       14 self.num_items = len(weights)
1582
1583
             \# --- your code here ---
1584
1585
            def solve(self) -> list[list[int]]:
       19 """
1586
       20 Solve the Bin Packing Problem.
1587
       21
1588
       22 Returns:
1589
       23 A list of lists, where each inner list represents a bin and
1590
                 contains the
             original indices of the items packed into it.
1591
           e.g., [[0, 2], [1, 3]] means item 0 and 2 are in the first bin,
1592
              and item 1 and 3 are in the second.
       26
1593
       27 """
1594
1595
       29 \# --- your code here ---
1596
       31 bins = [] \# Stores the content (indices) of each bin
1597
       32 bin_loads = [] \# Stores the current load of each bin
1598
1599
       34 \setminus# Store items as tuples of (index, weight) to keep track of original
1600
              indices
       35 items = sorted([(i, w) for i, w in enumerate(self.weights)], key=
1601
             lambda x: x[1], reverse=True)
1602
1603
       37 for item_index, item_weight in items:
1604
             placed = False
1605
             \# Try to place the item in an existing bin
       39
             for i in range(len(bins)):
1606
            if bin_loads[i] + item_weight <= self.capacity:</pre>
       41
1607
                bins[i].append(item_index)
       42
1608
                bin_loads[i] += item_weight
       43
1609
                placed = True
       44
1610
       45
                break
1611
       46
             \# If not placed, open a new bin
       47
1612
             if not placed:
       48
1613
            bins.append([item_index])
1614
       50
            bin_loads.append(item_weight)
1615
       51
       52 return bins
1616
1617
1618
```