MINIOPT: REASONING TO MODEL AND SOLVE GENERAL OPTIMIZATION PROBLEMS WITH LIMITED RESOURCES

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

034

037 038

039

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Modeling and solving optimization problems via large language models (LLMs) has attracted increasing attention recently. Although both prompt-based and learning-based methods have achieved progress, they remain limited by their reliance on large data volumes, high-quality annotations, expensive intermediate step verification, and huge computational overhead. From a data privacy perspective, a low-cost localized deployment of small-scale LLMs is of significant value. To train a small-scale LLM with excellent optimization generalization under limited resources, this paper proposes a reasoning to model and solve paradigm called MiniOpt based on reinforcement learning (RL) with verifiable reward. To reduce the demand for training data, MiniOpt adopts two-stage RL training. In the first stage the model quickly learns the model-and-solve paradigm and in the second stage it acquires strong optimization generalization ability. To reduce the cost of verifying the response of LLMs, OptReward in MiniOpt verifies the completeness of problem modeling and avoids the need for content validation. The above techniques enable the training of small-scale LLMs with strong optimization generalization ability under limited resources, thereby resulting in low inference cost for localized deployment and usage. Extensive experiments show that MiniOpt-3B exhibits strong optimization generalization across various optimization types and scenarios. For models with parameters fewer than 10B, MiniOpt-3B achieves the highest average solving accuracy (SA). For models with more than 10B parameters, MiniOpt-3B still shows competitive performance. Notably, MiniOpt-3B indicates superior SA on the hard OptMATH-Bench while only consuming 37.64% of the average output tokens required by DeepSeek-R1. The code is available at https://anonymous.4open.science/r/MiniOpt-6194.

1 Introduction

Optimization problems are ubiquitous in real-world scenarios, profoundly affecting diverse domains, including industrial production and transportation planning (Song et al., 2023; Gui et al., 2025). While traditional optimization solvers are efficient, their application heavily relies on expert knowledge, requiring the manual conversion of problems described in natural language into precise mathematical formulations or code, which is a process that both time-consuming and non-generalizable. The rise of LLMs has opened new pathways for the automated modeling and solving of optimization problems using natural language descriptions (Liu et al., 2024; Sun et al., 2024; Li et al., 2025a;b), bringing them closer to practical application scenarios. Representative work like LLMOPT (Jiang et al., 2025) or Text2Zinc (Singirikonda et al., 2025) significantly advances this field by parsing natural language descriptions into structured formulation, which are a unified general expression for optimization problems, and subsequently generating solution codes efficiently.

However, deploying such LLM-based approaches faces three critical bottlenecks. First, ensuring accurate text generations on a smaller parameter scale requires Supervised Fine-Tuning (SFT) with a large amount of training data (Lu et al., 2025). However, obtaining such data is very difficult, often requires a lot of time and effort, and is prone to errors (Xiao et al., 2025). Second, it is challenging to verify whether the result generated meets the requirements (Swamy et al., 2025). The non-verifiability of this task also exposes the limitations of the learning-based solving paradigm. It

is a barrier to the accuracy of the final result. Finally, to ensure the reliability of the result, common methods include reflection or self-correction mechanisms (Jiang et al., 2025; Li et al., 2025c), or multi-agent chains (AhmadiTeshnizi et al., 2024). However, these methods significantly multiply the computational overhead. These disadvantages pose significant barriers to implementing LLMs with powerful optimization-solving capabilities in the business scenarios of small and medium-sized enterprises or even on mobile devices. At the same time, considering data privacy issues, it is significant to locally deploy small-scale LLMs with strong performance while reducing training costs, including the volume and quality of training data.

To address these problems, this paper proposes a framework called MiniOpt to train a small-scale LLM with excellent optimization generalization under limited resources. To reduce the demand for data volume, MiniOpt adopts a two-stage RL training pipeline. Before conducting RL training, an SFT warm-up is applied to enable the model to get effective rewards in the early stages of RL. Subsequently, in the first RL stage, the model rapidly acquires the model-and-solve paradigm. In the second RL stage, training is conducted on a high-quality data subset obtained through type-scenario guided data selection, yielding strong optimization generalization (Jiang et al., 2025). The RL process is training using the proposed OptGRPO algorithm, which enhances data utilization efficiency and improves the model's ability to solve complex problems. To reduce the validation cost, we introduce OptReward in MiniOpt, which will perform three tasks: validating data format, ensuring the completeness of optimization problem modeling, and verifying solution accuracy. This encourages the model to learn in accordance with the model-and-solve paradigm. The framework enables the training of small-parameter LLMs with strong optimization generalization capabilities under limited resources, making localized deployment of LLMs for optimization feasible, thereby resulting in low inference cost for localized deployment and usage. Moreover, due to the high solving efficiency of MiniOpt-3B, reducing token usage also leads to considerable economic benefits.

Building upon the aforementioned methodology, this paper conducts extensive experiments with MiniOpt-3B on 9 benchmarks across different optimization types and problem scenarios. The results demonstrate its high optimization generalization capability. Compared to baselines below 10B parameters, MiniOpt-3B achieves the best performance. When evaluated against baselines exceeding 10B parameters, MiniOpt-3B exhibits only a 1.57% lower average solving accuracy (SA) than DeepSeek-R1, 0.37% higher than GPT-5 and outperforms LLMOPT-14B by 2.13% in average SA. Notably, MiniOpt Pareto dominates other baselines in terms of both parameter scale and the SA metric. On the hard OPTMATH-Bench, MiniOpt-3B has a higher SA than DeepSeek-R1 and only consumes about 37.64% of DeepSeek-R1's token count. Furthermore, results from ablation studies indicate that RL training contributes most significantly to MiniOpt-3B.

The subsequent sections review the related work, introduce the MiniOpt framework, present experimental results and analysis, provide an in-depth discussion, and finally conclude the paper.

2 RELATED WORK

LLMs for Modeling and Solving Optimization Problems. For modeling and solving optimization problems with LLMs, there are already a variety of benchmarks (Huang et al., 2025b; AhmadiTeshnizi et al., 2024; Yang et al., 2024b). Challenging benchmarks like OptMATH (Lu et al., 2025) and NL4Opt (Huang et al., 2025a) have led to numerous studies utilizing LLMs to solve optimization problems. Prompt-based approaches such as OptiMUS (AhmadiTeshnizi et al., 2024), CoE (Xiao et al., 2024) and Reflection (Shinn et al., 2023) utilizes the powerful generation capability of LLMs to generate the solver code of the optimization problem through multi-stage pipeline, without performing any post-training. Learning-based methods enhance LLMs' capabilities in modeling and solving mathematical problems. For example, LLaMoCo (Ma et al., 2024) proposes an SFT-based framework comprising a meticulously designed instruction set and a two-stage training methodology that incorporates contrastive learning warm-up followed by SFT. LLMOPT (Jiang et al., 2025) and NER4OPT (Dakle et al., 2023; Singirikonda et al., 2025) adopts a two-stage training process of modeling the optimization problems first and then solving them by generating solution code.

Reinforcement Learning with Verifiable Reward. While Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) plays a crucial role in post-training alignment, it suffers from high annotation costs and inherent human bias (Xiao et al., 2025). Reinforcement Learning with Verifiable Reward (RLVR) (Lambert et al., 2024) leverages externally grounded, easily verifi-

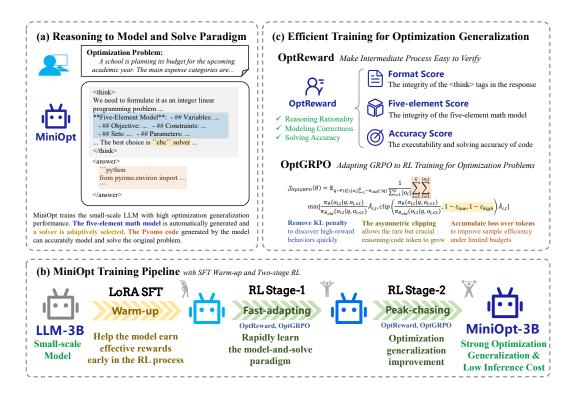


Figure 1: The framework of MiniOpt. Sub-figure (a) demonstrates the reasoning to model and solve paradigm of MiniOpt, encompassing problem modeling and solver adaptation during the thinking (i.e., reasoning) process of RL, and solution code generation in the response. Sub-figure (b) illustrates the training pipeline of MiniOpt, which involves sequential execution of SFT warm-up followed by two-stage RL. Sub-figure (c) presents the reward function OptReward and training algorithm OptGRPO used in MiniOpt's RL training.

able rewards (e.g., rule-based reward) to provide dense and structurally simple supervision (Wang et al., 2024; Xie et al., 2025; Jia et al., 2025; Gao et al., 2024). Its practicality is especially valuable in real-world black-box systems (Zhang et al., 2025; Xin et al., 2025; Pan et al., 2025), where verification is typically feasible only at the output stage, making RLVR a broadly applicable framework for aligning LLMs. In the field LLMs for solving optimization problems, SIRL (Chen et al., 2025) significantly improves the model's performance through the RLVR training paradigm.

3 METHODOLOGY: THE PROPOSED MINIOPT

3.1 OVERVIEW

This paper studies how to endow small-scale LLMs with strong optimization generalization under tight data and compute budgets. We introduce MiniOpt, whose framework is shown in Figure 1. MiniOpt introduces the reasoning process within the model-and-solve paradigm for optimization problems, the reward function and algorithm used for RL training, and the training pipeline for small-scale models. Through this framework, MiniOpt formulates the path from a natural language problem to an executable solver code as a single verifiable end-to-end task.

3.2 REASONING TO MODEL AND SOLVE PARADIGM

As shown in subfigure (a) in Figure 1, we introduce a reasoning to model and solve paradigm that turns a natural language optimization problem into a single verifiable objective. The paradigm is enforced by two compulsory and machine-parsable segments <think>...
think> and <answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...</answer>...

The first segment, enclosed by <think>...</think>, carries all modeling contents. It performs an analysis of the problem statement, specifies a complete five-element formulation inspired by LLMOPT (Jiang et al., 2025) and determines the appropriate open-source solver. Specifically, the optimization problem can be described as minimizing the **objective** function f(x) subject to the **constraints** $G(x) \leq c$, where the $x \in \mathcal{X} \subseteq \mathbb{R}^D$ is the D-dimensional decision **variables**, and $I = \{1, 2, \ldots, D\}$ is the index **sets**. \mathcal{X} is the feasible region, and $c \in \mathbb{R}^m$ provides the upper bounds. The constants in f(x) and G(x) form the **parameters** set, which also includes c. The five-element formulation, $\mathcal{M} = (\mathbf{Variables}, \mathbf{Objective}, \mathbf{Constraints}, \mathbf{Sets}, \mathbf{Parameters})$ map one-to-one to the components of an optimization problem. More details about the five-element formulation are provided in Appendix C.

During the reasoning process, to enhance the capability of LLM for solving optimization problems across diverse optimization types, it will analyze the problem during the thinking (i.e., reasoning) phase and select different open-source solvers for the given problem depending on the optimization types and the characteristics of the solvers, thereby ensuring an optimized match between the problem and the solver. Solver selection is guided by a prompt, which is introduced in Appendix N.

The second segment, enclosed by <answer>...</answer>, converts the five-element formulation into an executable Pyomo program that models the problem, invokes the solver, solves the instance, and prints the numerical answer. By constraining responses in this manner, we collapse the action space from free-form language to a programmable artifact whose intermediate structure and final result can be deterministically parsed and verified. The <answer> segment must implement the blueprint as a single python code fence containing a complete Pyomo script. Because these outputs are easy to verify, the OptReward in Section 3.4.1 can score format, five-element formulation, and accuracy based on rules via an automated procedure, providing low-cost supervision. During the training process, the prompt corresponding to this paradigm is shown in Appendix O.

3.3 TRAINING PIPELINE OF MINIOPT

Based on the paradigm described in Section 3.2 and the key techniques introduced in Section 3.4, we propose a training pipeline that enables MiniOpt to learn from limited resource constraints while achieving powerful solving performance and strong optimization generalization. This pipeline begins with a lightweight SFT warm-up followed by two-stage RL under OptReward. The illustration of the training pipeline is shown in subfigure (b) in Figure 1.

3.3.1 WARM-UP BASED ON LIGHTWEIGHT SFT

Small scale LLMs struggle to produce an end-to-end, executable trajectory for complex optimization tasks, and direct RL thus suffers from sparse rewards and unstable exploration. We therefore apply an SFT warm-up whose purpose is to provide a run-through capability, thereby allowing the model to obtain effective rewards in the early stages of RL training. This SFT warm-up does not adopt the reasoning to model and solve paradigm introduced in Section 3.2. Instead, it establishes a starting point upon which our two-stage RL can subsequently focus on paradigm acquisition and optimization generalization. The training data for warm-up are constructed from OptMATH-Train. For each instance, we prompt Qwen2.5-Coder-32B-Instruct (Hui et al., 2024) using the prompt in Appendix M to rewrite the original GurobiPy program into an equivalent Pyomo implementation and to select an appropriate open-source solver according to the detected structure. Every rewritten instance is executed and only those that compile, solve, and print the correct structure are retained. The remaining data serves as a candidate pool for the next stage of the two-stage RL training.

3.3.2 THE TWO-STAGE RL

After undergoing SFT warm-up, the models can conduct RL training more effectively. We employ a two-stage RL training under the same reward function (OptReward) for paradigm acquisition and optimization generalization, respectively. Both stages share the pipeline of parsing, executing, and scoring in Section 3.2 and the OptReward of Sections 3.4.1, they differ in training data and hyperparameter setting. Algorithmically, both stages use the same OptGRPO (cf. paragraphs in Section 3.4.2).

Stage-1 aims to enable the model to acquire the reasoning to model and solve Paradigm: the model must generate a valid <think>/<answer> pair, produce the executable Pyomo code, and make a coherent solver choice. To this end we train on 1,585 relatively easy problems, so that most signal arises from the formatting and structural components of OptReward, rapidly improving executability and the ability to solve the optimization problem of natural language description.

Stage-2 focuses on optimization generalization once the paradigm is established. The training distribution shifts to the problems with diverse optimization types and problem scenario, and the optimization emphasis moves to the accuracy score, encouraging refined modeling and solving behaviors (e.g., variable/constraint formulation and solver selection).

Based on the annotations of optimization types and problem scenario tags this paper assigned to each instance in OptMATH-Train, we sample a data subset from the candidate pool mentioned in Section 3.3.1 subject to two constraints: (i) type-uniform coverage, with exactly 600 instances per type, and (ii) within each type, the scenario frequencies match the distribution in the full dataset. The resulting data serves as the training data for the second stage of RL. While the dataset of the first stage is the union of the NL4Opt (AhmadiTeshnizi et al., 2024) and ICML Competition (Yang et al., 2024b) training splits. Detailed information on the construction of the training set is provided in the Appendix F.

Such a training strategy makes efficient use of limited data and reduces training costs, and ultimately allows for strong optimization generalization of LLMs with even small parameters and limited computational resources.

3.4 EFFICIENT TRAINING FOR OPTIMIZATION GENERALIZATION

Building upon the reasoning to model and solve paradigm mentioned in Section 3.2, we propose two key components for the RL training of MiniOpt as shown in subfigure(c) in Figure 1. An informative and easily verifiable reward function OptReward, and an improved algorithm OptGRPO builds upon GRPO (Shao et al., 2024).

3.4.1 OPTREWARD: VERIFIABLE REWARDS DESIGNED FOR MINIOPT

Based on the reasoning to model and solve paradigm described in section 3.2, we propose OptReward, including three automatically computed components: formatting correctness, structural sufficiency, and numerical accuracy. Each component is derived from deterministic parsing or execution of the output, It ensures the completeness of the modeling of problems but also enables verification to scale with machine time, yielding substantially lower verification costs.

Format Score: The format score $S_{\rm fmt}$ validates the response format. A response must contain exactly one <think>...</think> and one <answer>...</answer> in the correct order. If all conditions hold, we assign $S_{\rm fmt}=+1$, otherwise $S_{\rm fmt}=-1$. If the specified format is not present in the response, we deterministically set the remaining components to their error defaults, $S_{\rm five}=-1$ and $S_{\rm acc}=-2$, so that the total reward immediately reaches the global minimum. This forces the model to adopt the correct response format early in training and prevents expensive evaluation of malformed samples.

Five-element Score: To avoid ground-truth labeling of five-element content and the bias it may introduce, we use a presence-based rule aligned with the paradigm in Section 3.2. Conditional on valid formatting, we compute the five-element score $S_{\rm five}$. In the <think> segment, the model response is expected to include five labeled summaries starting with "## Sets:", "## Parameters:", "## Variables:", "## Objective:", and "## Constraints:". Each present summary adds 0.2 scores; if none is present we assign $S_{\rm five} = -1$. This structure shaping keeps the modeling blueprint parsable.

Accuracy Score: The accuracy score is obtained by executing the Pyomo program contained in <answer>. Failure to extract a program or to complete execution yields $S_{\rm acc}=-2$. When execution succeeds, we retrieve the optimal objective value \hat{f} from the model output and compare it with the ground-truth value f^* ; if they are equal, we assign $S_{\rm acc}=2$, otherwise $S_{\rm acc}=-1.5$.

The mathematical formulas of the three components of OptReward are provided in Appendix D. Combining the components with the formatting gate yields the total OptReward $R = S_{\rm fmt} + S_{\rm five} + S_{\rm acc}$ if $S_{\rm fmt} = 1$. When $S_{\rm fmt} = -1$, R = -4. By OptReward, the format score enforces the

strict <think>/<answer> format, the five-element score shapes a complete modeling blueprint in the think phase, and the accuracy score certifies correctness through equality of optimal objective values, enabling low-cost verifiable RL for problems in the field of optimization.

3.4.2 OPTGRPO: TRAINING SMALL-SCALE LLMs WITH LIMITED RESOURCES

GRPO (Shao et al., 2024) is an efficient RL algorithm, which replaces the critic model with a group baseline and updates the policy at the group level so that improves training stability and efficiency. The details of the algorithm are introduced in Appendix E. Based on the previous work on the improvement of GRPO like DAPO (Yu et al., 2025), we introduces three modifications, with the goal of achieving stronger optimization generalization within a limited budget.

First, building upon the strict gating already provided by OptReward through formatting score and five-element score, we set the coefficient on the KL penalty $\beta=0$ to remove the KL penalty to encourage greater exploration by the model. This approach does not jeopardize stability while facilitating faster discovery of high-reward behaviors. Second, to prevent entropy collapse and allow rare but crucial reasoning/code tokens to grow, we replace symmetric clipping with an asymmetric interval $[1-\varepsilon_{\rm low},\,1+\varepsilon_{\rm high}]$ with a higher upper clipping threshold $\varepsilon_{\rm high}$ than the lower threshold $\varepsilon_{\rm low}$. This relaxes the trust region on probability increases while keeping a firm lower bound on decreases, which empirically improves executability in stage-1 RL training and supports generalization in stage-2 RL training. Following Specifically, we raise $\varepsilon_{\rm high}$ to 0.28 during the training. Finally, to improve sample efficiency, critical for small-scale LLMs under limited data and verification budgets, instead of optimizing a sequence-averaged loss, we accumulate the loss over tokens and normalize by the total number of tokens in the group $\sum_i |o_i|$, which can yield denser learning signals for long outputs. The final loss of OptGRPO is as follows:

$$\begin{split} \mathcal{J}_{\text{OptGRPO}}(\theta) &= \mathbb{E}_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \\ &\min \left[\frac{\pi_{\theta}(o_{i,t}|q, o_{i, < t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i, < t})} \hat{A}_{i,t}, \text{ clip}(\frac{\pi_{\theta}(o_{i,t}|q, o_{i, < t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i, < t})}, 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}}) \hat{A}_{i,t} \right]. \end{split}$$

4 EXPERIMENT

We evaluate MiniOpt models on diverse optimization benchmarks spanning multiple types and scenarios to assess whether small scale parameter LLMs (3B/7B/14B) can achieve strong optimization generalization ability. We report Solving Accuracy (SA) as primary metrics, SA presents the proportion of output solutions of executed code that are numerically equal to the optimal solution provided from the labels of benchmarks. We use Execution Rate (ER) as a supplement that shows the proportion of generated code samples that run successfully without errors. Comparisons cover general LLMs, general reasoning LLMs, prompting-based baselines and learning-based baselines. The experiments aim to answer four key questions below.

- (Q1) Optimization Generalization Ability at Small-Scale LLMs: To what extent can MiniOpt at 3B/7B/14B achieve high SA across types and scenarios, and how does it compare with larger reasoning LLMs and prior learning-based approaches?
- (Q2) Pareto Front of Performance vs. Cost: What is the limit of the scale of model parameters for achieving strong optimization generalization?
- **(Q3) Importance of the Training Pipeline of MiniOpt**: How critical are the lightweight SFT warm-up and the two-stage RL to the performance of MiniOpt?
- **(Q4) Importance of the Reasoning to Model and Solve Paradigm and OptReward**: How do the proposed paradigm and the corresponding Opt Reward contribute to boost SA and ER in solving optimization problems?

The four questions are answered sequentially in the following sections, which first provide a detailed introduction to the experiments and then present an analysis of the results.

Table 1: Comparison of the SA metric across 9 benchmarks with rankings (NL4Opt, ICML Competition, Mamo Easy, Mamo Complex, NLP4LP, ComplexOR, IndustryOR, OptiBench, OptMATH-Bench). **Bold** indicates 1st, wavy underline indicates 2nd, <u>underline</u> indicates 3rd. "Rank*" represents the result of sorting methods among parameter sizes below 10B.

Category	Models / Methods	Avg.	Rank	Rank*	NL4Opt	ICML.C	Mamo.E	Mamo.C	NLP4LP	Com.OR	Indus.OR	OptiBench	OptMATH- Bench
	Qwen2.5-3B-Instruct	9.98	18	7	19.13	18.78	17.18	2.37	18.60	0.00	2.00	11.74	0.00
General Models	Qwen2.5-7B-Instruct	30.05	13	4	53.48	51.71	35.58	4.27	55.37	16.67	13.00	35.54	4.82
General Wodels	Qwen2.5-14B-Instruct	43.59	10	-	67.39	63.17	80.21	15.64	67.36	22.22	22.00	41.65	12.65
	DeepSeek-V3	57.81	4	-	78.26	77.56	84.82	26.54	79.34	44.44	26.00	64.13	39.16
	Qwen3-4B	10.25	17	6	16.52	17.56	13.80	6.64	15.29	5.56	2.00	11.90	3.01
	Qwen3-8B	19.77	15	5	30.87	29.51	23.93	9.95	34.71	11.11	6.00	28.26	3.61
General Models	Qwen3-14B	22.54	14	-	24.35	22.68	36.20	11.37	19.42	38.89	15.00	22.31	12.65
(Thinking)	DeepSeek-R1	58.51	3	-	83.91	75.37	74.54	39.81	69.83	44.44	32.00	66.94 ~~~~	39.76
	Gemini-2.5-Pro	57.04	5	-	78.26	71.22	65.95	30.81	73.55	50.00	28.00	61.32	54.21
	GPT-5	56.57	7	-	80.43	73.66	58.12	23.22	73.14	61.11	26.00	64.63	48.80
Prompt-based	Chain-of-Experts	41.03	12	-	66.52	56.59	63.65	22.75	59.09	33.33	19.00	45.29	3.01
Methods	OptiMUS	18.76	16	-	13.48	33.17	37.27	11.85	18.18	16.67	8.00	26.61	3.61
Methods	Reflexion	41.28	11	-	56.52	52.20	84.82	18.01	53.72	38.89	19.00	41.16	7.23
Learning-based	OptMATH-7B	52.37	9	3	78.70	66.83	84.20	34.12	68.60	33.33	19.00	52.23	34.34
Models	LLMOPT-14B	54.81	8	-	80.28	75.35	89.53	44.08	73.42	35.29	29.00	53.83	12.50
	MiniOpt-3B	56.94	6	2	83.04	68.05	85.43	35.07	73.55	50.00	21.00	53.55	42.77
Ours	MiniOpt-7B	62.76	2	ì	89.13	77.56	88.34	38.39	79.34	55.56	26.00	59.34	51.20
	MiniOpt-14B	66.10	î	-	92.17	86.34	90.80	33.65	79.75	61.11	27.00	67.44	56.63

4.1 EXPERIMENTAL SETUP

The evaluation encompasses nine benchmarks about operations research optimization: NL4Opt (Ramamonjison et al., 2022), Mamo (Easy and Complex subsets) (Huang et al., 2025b), IndustryOR (Huang et al., 2025a), NLP4LP (AhmadiTeshnizi et al., 2024), ComplexOR (Xiao et al., 2024), OptMATH (Lu et al., 2025), OptiBench (Yang et al., 2025b), and ICML Competition (Yang et al., 2024b). We follow the same setting in LLMOPT (Jiang et al., 2025) to ensure consistency and comparability. For the newly included datasets, OptMATH-Bench and OptiBench, we adhere to their original data splits provided by the authors. For ICML Competition dataset, we use it exclusively as out-of-distribution test data to evaluate the generalization capability of our model.

To validate the correctness of the solutions obtained by MiniOpt and all baselines, this paper compares their log files generated during the execution process against the ground-truth solutions provided by the benchmarks. A solution is deemed correct if the log file yields a non-empty match with the ground-truth; otherwise, it is considered unsuccessful.

4.2 Analysis of Optimization Generalization

In this section, we compare MiniOpt with general LLMs (Qwen2.5 Series, DeepSeek-V3), general reasoning LLMs (Qwen3 Series, DeepSeek-R1, Gemini-2.5-Pro, GPT-5), prompt-based methods (Chain-of-Experts, OptiMUS, Reflexion), learning-based methods (OptMATH-7B, LLMOPT-14B), demonstrating the optimization generalization capability of MiniOpt. The information on these methods can be found in Appendix B. Tables 1 summarizes SA on nine benchmarks that span 7 optimization types and 22 scenarios, and the statistics on problem categories and scenarios are provided in Appendix A.2. The ER metrics for all the methods on the nine benchmarks and its analysis can be found in the Appendix G.

Overall Performance (**Answer to Q1**). Across all nine benchmarks, MiniOpt-7B achieved the strongest average performance among all baselines. Notably, MiniOpt-3B surpassed all prompt-based and learning-based methods. For instance, compared to LLMOPT-14B, the state-of-the-art learning-based model, MiniOpt-3B achieved an average SA that is 2.13% higher. When compared to DeepSeek-V3 and DeepSeek-R1, MiniOpt-3B scored 0.87% and 1.57% lower, respectively. Furthermore, as a higher-parameter variant of MiniOpt, MiniOpt-14B achieves the highest average SA of 66.10%, significantly raising the performance ceiling of MiniOpt in practical applications.

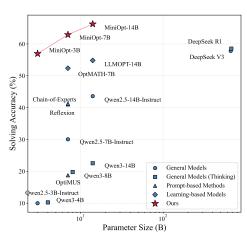
Besides, SIRL (Chen et al., 2025) is not included in Table 1 because the paper was accepted only shortly before our submission deadline. Although it proposes a reasoning model training pipeline, the important issue of improving the optimization generalization ability of small-scale models with limited data and computing resources has not been studied. According to the results presented in the

SIRL paper (Chen et al., 2025), MiniOpt-7B outperforms the SIRL-Qwen2.5-7B model of similar size by 22.20% solving accuracy on the hard OptMATH-Bench benchmark.

Competitiveness of Small-Scale Models (Answer to Q1). MiniOpt remains competitive even at smaller scales. MiniOpt-3B reaches 56.94% SA on average, this performance already matching or exceeding several much larger reasoning models (e.g., the average SA higher than GPT-5 at 56.57%) and clearly outperforming general-purpose 14B post-trained models (e.g., +13.35% over Qwen2.5-14B-Instruct on average). Performance grows smoothly with scale under the same training pipeline, The average SA metric grows by +5.82% when parameters increase from 3B to 7B, and by +9.16% when further scaled to 14B.

Challenging Benchmarks (Answer to Q1). On the most demanding sets that require faithful modeling and solver usage, MiniOpt shows clear advantages. On the latest challenging benchmark OptMATH-Bench, MiniOpt-14B achieves 56.63% SA, outperforming Gemini-2.5-Pro (54.21%) and GPT-5 (48.80%). Even on extremely high-dimensional test sets such as Indus.OR, where small-parameter MiniOpt does not attain the highest SA due to limitations in instruction following, it still delivers competitive performance levels in both metrics.

Breadth across Types and Scenarios (Answer to Q1). To rigorously evaluate the optimization generalization, we analyze the SA of MiniOpt-7B and MiniOpt-3B across three optimization types and three application scenarios under two difficulty levels. At the medium-difficulty OptiBench, MiniOpt-7B achieved 59% (M), 74% (T), 68% (L), and 65% (LP), 63% (IP), 45% in Mixed-Integer Linear Programming (MILP). MiniOpt-3B scored 55% (M), 71% (T), 59% (L), and 60% (LP), 58% (IP), 38% (MILP). On the hard-difficulty OptMATH-Bench, MiniOpt-7B achieved 58% (M), 56% (T), 57% (L), and 55% (LP), 100% (IP), 56% (MILP). MiniOpt-3B attained 43% (M), 56% (T), 76% (L), and 73% (LP), 67% (IP), 47% (MILP). These results reflect consistent and scalable generalization across optimization types, problem scenarios, and difficulty levels. In order to investigate whether the pipeline training in this paper would impair other abilities of the LLMs, we also discussed the generalization performance of MiniOpt-3B and 7B under different tasks except optimization problems. Detailed information is presented in Appendix K.



Solving Accuracy vs. Parameter Size

Figure 2: Comparison of average SA against model parameter scales for various methods. MiniOpt is the Pareto optimal among compared methods.

4.3 PARETO FRONT OF PERFORMANCE VS. COST

Analysis of the Pareto Front (Answer to Q2). Figure 2 and Figure 5 in Appendix I indicate that MiniOpt series (represented by the solid red line establishes a new and superior Pareto front in the performance-versus-cost trade-off. Since the parameter size of GPT-5 and Gemini-2.5-Pro have not been disclosed, we do not label these two models in the figures. MiniOpt demonstrates strong competitiveness even at smaller scales. For instance, the MiniOpt-3B model, with only 3 billion parameters, achieves an average SA of 56.94%. This performance already surpasses that of larger-scale general reasoning models like GPT-5 (cf. Table 1) and significantly outperforms Qwen2.5-14B-Instruct. As the scale of the model increases, the average SA performance of MiniOpt also grows steadily. The MiniOpt-14B variant achieves the highest average SA of 66.10% among all models. It achieves a comprehensive performance lead while having substantially fewer parameters than top-tier general reasoning models such as DeepSeek-R1.

4.4 ABLATION STUDY

In this section we ablate core components of MiniOpt-3B and report results of the SA metric in Table 2 below and the ER metric in Table 6 in Appendix H. As evidenced in Table 2, each module

Table 2: Ablation study (MiniOpt-3B) on the SA metric across 9 benchmarks.

Category	Model / Method	Avg.	NL4Opt	ICML.C	Mamo.E	Mamo.C	NLP4LP	Com.OR	Indus.OR	OptiBench	OptMATH-Bench
	MiniOpt-3B	56.94	83.04	68.05	85.43	35.07	73.55	50.00	21.00	53.55	42.77
	MiniOpt-3B w/o SFT Warm-up	53.59	83.91	70.00	83.28	33.65	68.18	38.89	15.00	53.88	35.54
	MiniOpt-3B w/o RL	39.25	52.61	48.29	80.06	25.12	60.33	16.67	15.00	34.71	20.48
Ablations	MiniOpt-3B w/o Two-stage RL	54.12	83.48	69.02	82.67	33.65	69.01	44.44	16.00	53.88	34.94
Adiations	MiniOpt-3B w/o Data Selection	53.50	78.70	68.05	84.82	30.81	72.31	38.89	19.00	52.89	36.75
	MiniOpt-3B w/ Random Selection	50.94	83.48	69.76	82.98	29.38	70.66	27.78	14.00	52.73	27.71
	MiniOpt-3B w/o OptReward	52.44	78.70	66.83	82.52	33.18	70.25	33.33	18.00	51.24	37.95
	MiniOpt-3B w/o OptGRPO	48.33	70.87	60.00	81.29	32.23	65.70	27.78	17.00	45.79	34.34

of the proposed reasoning to model and solve paradigm demonstrates substantial contributions to modeling and solving optimization problems with smaller-scale models under limited training resources. Among these, RL provides the most significant improvement, highlighting the importance of the proposed paradigm. For a detailed analysis of the proposed model-and-solve paradigm, the training process, and the significance of OptReward and OptGRPO, please see Appendix H.

5 DISCUSSION

Exploration of Models with Smaller Parameter Scale. To further explore the performance of MiniOpt with smaller parameter scales, thereby better balancing the trade-off between parameters and performance, we deploy MiniOpt-1.5B with the same reasoning to model and solve paradigm (Section 3.2), OptReward (Section 3.4.1), and training pipeline (Section 3.3). The results shown in Table 7 and 8 in Appendix J demonstrate that when the model scale is reduced from 3B to 1.5B, MiniOpt achieves an average SA of 46.15% and an average ER of 80% across the 9 benchmarks, still surpassing all baselines except for learning-based methods, maintaining a strong competitiveness.

Dimensionality Scalability Issue. Our empirical results show that MiniOpt-3B can correctly model optimization problems with dimensionalities as high as 72 on the OptMATH-Bench task and 80 on the Mamo Complex task. This constitutes a remarkably strong performance for models of such limited scale, which highlights that the MiniOpt pipeline effectively enhances the model's optimization modeling capacity and retains robust generalization even in high-dimensional settings. For problems with even higher dimensionality, the descriptions of numerical parameters typically become substantially longer. Such extended input sequences inherently pose a significant challenge to the comprehension capabilities of small-scale models. As corroborating evidence, larger models such as MiniOpt-7B successfully solve problems up to 109 dimensions on the OptMATH-Bench task.

Cost Savings of MiniOpt and Generality of Base Models. We find that MiniOpt-3B achieves a 3.01% higher SA than DeepSeek-R1 on the OptMATH-Bench while using 62.4% fewer average output tokens. For MiniOpt-7B, the corresponding improvements are 11.44% higher SA and a 39.6% reduction in average output tokens. More details and discussions are provided in Appendix L. This demonstrates that MiniOpt constitutes a general framework capable of enabling small-scale models to achieve strong optimization generalization under constrained data and computational resources. Furthermore, although Qwen2.5 series models are adopted as the base model in this paper, we have also conducted experiments with different base models and obtained closely aligned results.

6 CONCLUSION

This paper proposes a novel reasoning to model and solve paradigm and the small scale model, MiniOpt-3B, achieves higher performance with a small-scale parameter and limited resources. We explore the optimization generalization of the model in various types of optimization, problem scenarios, and high-variable dimensions. Empirical results demonstrate that MiniOpt exhibits strong generalization performance under these varying conditions. Furthermore, this study explores the minimum parameter scale required for MiniOpt to maintain competitive performance. Future work includes exploring efficient modeling and solving methods for optimization problems with high-dimensional variables or a large number of constraints.

ETHICS AND REPRODUCIBILITY STATEMENT

Ethics. This work does not involve any human subjects, personal data, or sensitive information. All the test datasets used are publicly available, and no proprietary or confidential information is used.

Reproducibility. Experimental settings are described in Section 4.1 and datasets included in Appendix F. The code is available at https://anonymous.4open.science/r/MiniOpt-6194.

LLM USAGE STATEMENT

No LLMs were used in the research ideation and paper writing of this work.

REFERENCES

486

487 488

489

490

491

492

493 494

495 496

497 498 499

500

501

502

504

505

506

507

509

510

511

512

513 514

515

516

517

518

519 520

521

522

523

524

525

527

528

529

530

531

532

534

535

536

538

Ali AhmadiTeshnizi, Wenzhi Gao, and Madeleine Udell. OptiMUS: Scalable optimization modeling with (MI)LP solvers and large language models. In *Advances in Forty-first International Conference on Machine Learning*, Vienna, Austria, 2024.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021.

Yitian Chen, Jingfan Xia, Siyu Shao, Dongdong Ge, and Yinyu Ye. Solver-Informed RL: Grounding large language models for authentic optimization modeling. *CoRR*, abs/2505.11792, 2025.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. CoRR, abs/2110.14168, 2021.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit S. Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornraphop Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilaï Deutel, Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller, Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero Rissa, Dan Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves, Cindy Wang, Nitesh Bharadwaj Gundavarapu, Ilia Shumailov, Bo Wang, Mantas Pajarskas, Joe Heyward, Martin Nikoltchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik, Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kavukcuoglu, Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov, Bo-Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas Leichner, Haichuan Yang, Zelda Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios Vytiniotis, Jieru Mei, and Mu Cai. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. CoRR, abs/2507.06261, 2025.

Parag Pravin Dakle, Serdar Kadioglu, Karthik Uppuluri, Regina Politi, Preethi Raghavan, SaiKrishna Rallabandi, and Ravisutha Srinivasamurthy. Ner4Opt: Named entity recognition for optimization modelling from natural language. In *Integration of Constraint Programming, Artificial*

Intelligence, and Operations Research - 20th International Conference, volume 13884, pp. 299–319, Nice, France, 2023.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, and Wangding Zeng. DeepSeek-V3 technical report. CoRR, abs/2412.19437, 2024.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *CoRR*, abs/2501.12948, 2025.

- Jiaxuan Gao, Shusheng Xu, Wenjie Ye, Weilin Liu, Chuyi He, Wei Fu, Zhiyu Mei, Guangju Wang, and Yi Wu. On designing effective RL reward at training time for LLM reasoning. *CoRR*, abs/2410.15115, 2024.
- Runquan Gui, Zhihai Wang, Jie Wang, Chi Ma, Huiling Zhen, Mingxuan Yuan, Jianye HAO, Defu Lian, Enhong Chen, and Feng Wu. Hypertree planning: Enhancing LLM reasoning via hierarchical thinking. In *Forty-second International Conference on Machine Learning*, 2025.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *9th International Conference on Learning Representations*, 2021a.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021b.
- Chenyu Huang, Zhengyang Tang, Shixi Hu, Ruoqing Jiang, Xin Zheng, Dongdong Ge, Benyou Wang, and Zizhuo Wang. ORLM: A customizable framework in training large models for automated optimization modeling. *Operations Research*, May 2025a. ISSN 1526-5463. doi: 10.1287/opre.2024.1233.
- Xuhan Huang, Qingning Shen, Yan Hu, Anningzhe Gao, and Benyou Wang. LLMs for mathematical modeling: Towards bridging the gap between natural and mathematical languages. In *Findings of the Association for Computational Linguistics* 2025, pp. 2678–2710, Albuquerque, New Mexico, 2025b.

- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,
 Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang, Xingzhang Ren, Xuancheng Ren, Jingren
 Zhou, and Junyang Lin. Qwen2.5-Coder technical report. *CoRR*, abs/2409.12186, 2024.
 - Zeyu Jia, Alexander Rakhlin, and Tengyang Xie. Do we need to verify step by step? rethinking process supervision from a theoretical perspective. In *Forty-second International Conference on Machine Learning*, 2025.
 - Caigao Jiang, Xiang Shu, Hong Qian, Xingyu Lu, Jun Zhou, Aimin Zhou, and Yang Yu. LLMOPT: learning to define and solve general optimization problems from scratch. In *Advances in The Thirteenth International Conference on Learning Representations*, Singapore, 2025.
 - Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, 2017.
 - Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard H. Hovy. RACE: large-scale ReAding comprehension dataset from examinations. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017.
 - Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. TÜLU 3: Pushing frontiers in open language model post-training. *CoRR*, abs/2411.15124, 2024.
 - Jinyuan Li, Yi Chu, Yiwen Sun, Mengchuan Zou, and Shaowei Cai. AutoPBO: LLM-powered optimization for local search PBO solvers. *CoRR*, abs/2509.04007, 2025a.
 - Kai Li, Fei Liu, Zhenkun Wang, Xialiang Tong, Xiongwei Han, Mingxuan Yuan, and Qingfu Zhang. Ars: Automatic routing solver with large language models. *CoRR*, abs/2502.15359, 2025b.
 - Xiaozhe Li, Jixuan Chen, Xinyu Fang, Shengyuan Ding, Haodong Duan, Qingwen Liu, and Kai Chen. OPT-BENCH: evaluating LLM agent on large-scale search spaces optimization problems. *CoRR*, abs/2506.10764, 2025c.
 - 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 *Advances in Forty-first International Conference on Machine Learning*, Vienna, Austria, 2024.
 - Hongliang Lu, Zhonglin Xie, Yaoyu Wu, Can Ren, Yuxuan Chen, and Zaiwen Wen. OptMATH: A scalable bidirectional data synthesis framework for optimization modeling. In *Forty-second International Conference on Machine Learning*, 2025.
 - Zeyuan Ma, Hongshu Guo, Jiacheng Chen, Guojun Peng, Zhiguang Cao, Yining Ma, and Yue-Jiao Gong. LLaMoCo: Instruction tuning of large language models for optimization code generation. *CoRR*, abs/2403.01131, 2024.
 - OpenAI. GPT-5 system card. https://cdn.openai.com/gpt-5-system-card.pdf, 2025.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, New Orleans, LA, 2022.
 - Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with SWE-gym. In *ICLR 2025 Third Workshop on Deep Learning for Code*, 2025.

Rindranirina Ramamonjison, Timothy Yu, Raymond Li, Haley Li, Giuseppe Carenini, Bissan Ghaddar, Shiqi He, Mahdi Mostajabdaveh, Amin Banitalebi-Dehkordi, Zirui Zhou, and Yong Zhang. NL4Opt competition: Formulating optimization problems based on their natural language descriptions. In *Proceedings of the NeurIPS 2022 Competitions Track*, volume 220, pp. 189–203, 2022.

- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. DeepSeekMath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023*, New Orleans, LA, 2023.
- Akash Singirikonda, Serdar Kadioglu, and Karthik Uppuluri. Text2Zinc: A cross-domain dataset for modeling optimization and satisfaction problems in MiniZinc. *CoRR*, abs/2503.10642, 2025.
- Wen Song, Xinyang Chen, Qiqiang Li, and Zhiguang Cao. Flexible job-shop scheduling via graph neural network and deep reinforcement learning. *IEEE Transactions on Industrial Informatics*, 19(2):1600–1610, 2023.
- Yiwen Sun, Furong Ye, Xianyin Zhang, Shiyu Huang, Bingzhen Zhang, Ke Wei, and Shaowei Cai. AutoSAT: Automatically optimize SAT solvers via large language models. *CoRR*, abs/2402.10705, 2024.
- Gokul Swamy, Sanjiban Choudhury, Wen Sun, Zhiwei Steven Wu, and J. Andrew Bagnell. All roads lead to likelihood: The value of reinforcement learning in fine-tuning. *CoRR*, abs/2503.01067, 2025.
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-Shepherd: Verify and reinforce LLMs step-by-step without human annotations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9426–9439, Bangkok, Thailand, 2024.
- Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin Fu, Tao Zhong, Jia Zeng, Mingli Song, and Gang Chen. Chain-of-Experts: When LLMs meet complex operations research problems. In *The Twelfth International Conference on Learning Representations*, Vienna, Austria, 2024.
- Ziyang Xiao, Jingrong Xie, Lilin Xu, Shisi Guan, Jingyan Zhu, Xiongwei Han, Xiaojin Fu, Wing Yin Yu, Han Wu, Wei Shi, Qingcan Kang, Jiahui Duan, Tao Zhong, Mingxuan Yuan, Jia Zeng, Yuan Wang, Gang Chen, and Dongxiang Zhang. A survey of optimization modeling meets LLMs: Progress and future directions. In James Kwok (ed.), *Proceedings of the Thirty-Fourth International Joint Conference on Artificial Intelligence, IJCAI-25*, pp. 10742–10750. International Joint Conferences on Artificial Intelligence Organization, 8 2025. doi: 10.24963/ijcai.2025/1192. Survey Track.
- Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. Logic-RL: Unleashing LLM reasoning with rule-based reinforcement learning. CoRR, abs/2502.14768, 2025.
- Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong Shao, Wanjia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Haowei Zhang, Qihao Zhu, Dejian Yang, Zhibin Gou, Z. F. Wu, Fuli Luo, and Chong Ruan. DeepSeek-Prover-V1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. In *Advances in the Thirteenth International Conference on Learning Representations*, Singapore, 2025.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu

Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *CoRR*, abs/2412.15115, 2024a.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jian Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *CoRR*, abs/2505.09388, 2025a.

Zhicheng Yang, Yinya Huang, Wei Shi, Liang Feng, Linqi Song, Yiwei Wang, Xiaodan Liang, and Jing Tang. Benchmarking LLMs for optimization modeling and enhancing reasoning via reverse socratic synthesis. *CoRR*, abs/2407.09887, 2024b.

Zhicheng Yang, Yiwei Wang, Yinya Huang, Zhijiang Guo, Wei Shi, Xiongwei Han, Liang Feng, Linqi Song, Xiaodan Liang, and Jing Tang. OptiBench meets ReSocratic: Measure and improve LLMs for optimization modeling. In *The Thirteenth International Conference on Learning Representations*, Singapore, 2025b.

Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. DAPO: an open-source LLM reinforcement learning system at scale. *CoRR*, abs/2503.14476, 2025.

Kechi Zhang, Ge Li, Yihong Dong, Jingjing Xu, Jun Zhang, Jing Su, Yongfei Liu, and Zhi Jin. CodeDPO: Aligning code models with self generated and verified source code. In *Advances in the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15854–15871, Vienna, Austria, 2025.

A DATASETS

A.1 THE INTRODUCTION OF EVALUATION DATASETS

In this section, we provide an overview of the datasets used for performance evaluation in our experiments. These datasets cover a wide range of optimization types and scenarios, ensuring the robustness and generalization of our proposed method. In our practice, we use the version of the benchmark datasets above from https://github.com/antgroup/LLMOPT.

Table 3: Statistics of the optimization problem datasets

Dataset Name	# of Data
NL4Opt	230
ICML.C	410
Mamo.E	652
Mamo.C	211
NLP4LP	242
Com.OR	18
Indus.OR	100
OptiBench	605
OptMATH-Bench	166

NL4Opt (Ramamonjison et al., 2022) dataset is curated from the NL4Opt Competition. For this benchmark, we used the test split containing 230 annotated linear programming word problems after manually removing 15 unsolvable problems from the original 245 problems. Each problem is sourced from domains such as sales, advertising, and investment, ensuring a balanced representation.

Mamo (Huang et al., 2025b) dataset (optimization split of the original Mamo dataset) consists of two parts: Easy_LP and Complex_LP. These two subsets provide 652 high-school-level and 211 undergraduate-level Linear Programming (LP) and Mixed-Integer Linear Programming (MILP) problems, respectively.

IndustryOR (Huang et al., 2025a) is the first industrial dataset specifically designed for optimization modeling. It incorporates data from 13 different industries and covers a variety of real-world scenarios. The dataset includes real operations research problems from eight different industries, covering five types of optimization problems, and divided into three difficulty levels. The test dataset contains 100 instances with optimal solutions.

NLP4LP (AhmadiTeshnizi et al., 2024) dataset includes 242 feasible samples sourced from optimization textbooks and lecture notes. These problems cover areas such as facility location, network flow, scheduling, and portfolio management. Each instance in NLP4LP includes a description, sample parameter data file, and optimal value derived from textbook solutions or manual solving, offering a range of complex optimization challenges of varying difficulty levels.

ComplexOR (Xiao et al., 2024) dataset is developed in collaboration with three specialists in operations research. It contains 18 samples sourced from diverse references such as academic papers, textbooks, and real-world industrial scenarios. These problems encompass a broad spectrum of topics, including supply chain optimization, scheduling problems, and warehousing optimization, providing comprehensive and complex optimization challenges.

OptMATH-Bench (Lu et al., 2025) is a large-scale, challenging benchmark specifically designed to evaluate the optimization modeling capabilities of LLMs, encompassing diverse optimization problem types across 10+ real-world application domains such as logistics, manufacturing, transportation, and finance. The benchmark features significantly more complex problem descriptions with an average length 2.9× longer than Mamo Easy, containing extended natural language contexts and intricate constraints that pose greater challenges.

OptiBench (Yang et al., 2025b) is a comprehensive benchmark for evaluating large language models' end-to-end optimization problem-solving capabilities. The dataset contains 605 carefully curated optimization problems that span multiple optimization types and formats. OptiBench includes problems of Linear Programming (LP), Integer Programming (IP), and Mixed-Integer Linear Programming (MILP), encompassing a wide range of optimization complexities.

ICML Competition (Yang et al., 2024b) dataset comprises data from the ICML 2024 Challenges on Automated Math Reasoning - Track 3: Automated Optimization curated from the competition's test split. Since the original ground truth is not released by the organizers, all solutions in this dataset are manually labeled. The dataset serves as a challenging benchmark for evaluating end-to-end optimization reasoning and problem-solving capabilities of language models.

A.2 THE DISTRIBUTION OF OPTIMIZATION TYPES AND PROBLEM SCENARIOS OF BENCHMARKS

To evaluate the generalization ability of the MiniOpt across different problem scenarios through experiments, this paper has counted the number of optimization types and scenarios in 9 benchmarks. The distribution histogram of optimization types in the nine benchmarks used in this paper is shown in Figure 3, and the distribution histogram of problem scenarios is shown in Figure 4.

A.3 TRAINING DATASETS FOR SFT WARM-UP AND TWO-STAGE RL

The data volume and data sources for each stage of MiniOpt's training process are illustrated in Table 4.

B BASELINES

B.1 GENERAL MODELS

Qwen2.5-3B/7B/14B (Yang et al., 2024a). The models of the Qwen2.5 series are widely adopted as base models or baselines. The series showcases significant enhancements such as substantially

Type	Com.OR	ICML.C	Indus.OR	Mamo.C	Mamo.E	NL4Opt	NLP4LP	OptMAT H-Bench	Optibench
CO	4	8	12	39	0	0	0	19	11
IP	4	211	23	11	2	161	167	7	291
LP	9	175	31	123	650	69	75	16	214
MILP	1	16	30	38	0	0	0	91	44
MOP	0	0	4	0	0	0	0	0	0
NLP	0	0	0	0	0	0	0	0	45
SOCP	0	0	0	0	0	0	0	33	0
# Data	18	410	100	211	652	230	242	166	605

Figure 3: Histogram showing the distribution of optimization types across 9 benchmarks. We categorize the problems in the benchmarks into these types: Combinatorial Optimization (CO), Integer Programming (IP), Linear Programming (LP), Mixed-Integer Linear Programming (MILP), Multi-Objective Optimization Problems (MOP), Nonlinear Programming (NLP), Second-Order Cone Programming (SOCP).

Scenario	Com.OR	ICML.C	Indus.OR	Mamo.C	Mamo.E	NL4Opt	NLP4LP	OptMAT H-Bench	Optibench
Agriculture	0	37	5	4	31	13	12	0	46
Aviation	3	2	2	2	0	2	1	12	2
Construction	0	1	1	1	42	1	1	0	8
Education	0	5	1	0	30	0	0	0	7
Energy	0	0	0	5	25	1	2	7	15
Environment	0	0	0	0	31	0	1	0	0
Finance	0	16	7	5	85	5	5	6	22
Healthcare	0	15	1	5	31	23	29	0	20
Logistics	1	27	8	32	25	21	23	22	42
Manufacturing	5	169	37	12	34	77	83	56	254
Marketing	1	9	0	0	43	2	2	0	11
Military	0	0	2	0	36	0	0	0	3
Public Utilities	0	9	1	4	7	2	2	2	19
Resources	0	15	0	12	9	11	8	3	17
Retail	0	16	4	2	31	8	9	1	19
Services	4	38	10	6	73	29	29	7	40
Sports	0	0	0	0	31	0	0	0	0
Supply Chain	2	12	7	48	31	5	3	11	11
Technology	0	0	1	3	0	0	0	1	3
Telecommunications	1	1	0	8	27	0	0	9	1
Transportation	1	31	8	57	45	23	26	9	40
Other	- 0	7	5	5	15	7	6	20	25
# Data	18	410	100	211	682	230	242	166	605

Figure 4: Histogram showing the distribution of optimization problem scenarios across 9 benchmarks. We categorize the problems in the benchmarks into these scenarios: Agriculture, Aviation, Construction, Education, Energy, Environment, Finance, Healthcare, Logistics, Manufacturing, Marketing, Military, Public Utilities, Resources, Retail, Services, Sports, Supply Chain, Technology, Telecommunications, Transportation, Other.

improved knowledge, coding, and mathematical capabilities. Key features excel at instruction following, processing long contexts up to 128K tokens, and robustly handling structured data like JSON.

DeepSeek-V3 (DeepSeek-AI et al., 2024). DeepSeek-V3 introduces a sparse Mixture-of-Experts (MoE) model with 671B total parameters. It achieves training efficiency through Multi-head Latent Attention (MLA) architecture and an auxiliary-loss-free load balancing strategy. Pretrained on 14.8T tokens with an SFT and reinforcement learning (RL).

B.2 GENERAL REASONING MODELS

Qwen3-4B/8B/14B (Yang et al., 2025a). Qwen3 pioneers a unified architecture (0.6B to 235B) integrating thinking mode (complex reasoning) and non-thinking mode (rapid responses) with dynamic switching. Its thinking budget mechanism enables adaptive computational allocation. The series outperforms larger MoE models in tasks such as coding, mathematics, and agent application.

Table 4: The number of samples in the dataset for training. Note that all data mentioned in the table comes entirely from the training split of the corresponding dataset.

Training Stage	Dataset Size	Data Source
SFT Warm-up	140K	OptMATH-Train
RL-Stage 1	1585	NL4Opt (Train) & ICML.C (Train)
RL-Stage 2	3000	OptMATH-Train

DeepSeek-R1 (DeepSeek-AI et al., 2025). The DeepSeek-R1 is an enhanced model based on DeepSeek-R1-Zero presented in this work. It's a purely RL-driven reasoning model requiring no SFT pretraining. As the most representative model with thinking ability, DeepSeek-R1 is an important baseline for reasoning models.

Gemini-2.5-Pro (Comanici et al., 2025). Gemini-2.5-Pro is a powerful multimodal agent that has excellent programming / reasoning performance and enables the processing of long video content. The Gemini-2.5 family spans the full Pareto front of capability-cost optimization. Its integration of long-context understanding, multimodality, and reasoning unlocks novel agentic applications.

GPT-5 (OpenAI, 2025). GPT-5 is the latest unified, router-mediated system that instantiates a spectrum of language-model instances ranging from a high-throughput, low-latency model (gpt-5-main) to a deliberative, compute-intensive reasoning model (gpt-5-thinking). The router selects the appropriate instantiation by conditioning on conversation type, task complexity, tool requirements, and explicit user directives, thereby optimizing both instruction adherence and inference efficiency.

B.3 PROMPT-BASED METHODS

Reflection (Shinn et al., 2023). Reflexion is an enhanced language agent framework utilizing feedback mechanisms. It enables agents to excel at sequential decision-making tasks through task feedback analysis and memory buffering without requiring weight updates. This framework accommodates diverse feedback signals and demonstrates effectiveness across programming, math problems and language reasoning domains.

OptiMUS (AhmadiTeshnizi et al., 2024). OptiMUS is a highly modular solver that leverages the text understanding and generating capabilities of LLMs. It constructs specialized agents for entity extraction, mathematical modeling, and code generation using concise prompts, while incorporating a reflection mechanism for iterative improvement.

Chain-of-Experts (Xiao et al., 2024). Chain-of-Experts is a multi-agent framework specifically designed for operations research optimization problems. The system features a central controller that coordinates an interaction sequence among specialized agents, including a term interpreter, modeling agent, and programming expert. Thus solving optimization problems through precise coordination of multiple modules.

B.4 LEARNING-BASED MODELS

LLMOPT-14B (Jiang et al., 2025). LLMOPT is a novel framework for optimization problem solving that leverages LLMs. It begins by formulating a unified representation of optimization problems, thereby enhancing the model's ability to generalize across diverse types of scenarios. Based on this unified description of five-element formulation, the framework generates the solving code. LLMOPT uses multi-instruction SFT and KTO alignment during training to enhance modeling accuracy and reduce model hallucinations.

OptMATH-Qwen2.5-7B (Lu et al., 2025). OptMATH-Qwen2.5-7B is trained end-to-end on the OptMATH-Train dataset, it generates both mathematical formulations and solver code from problem descriptions. The input consists of textual problem specifications, while the target output comprises concatenated sequences. Optimization follows the standard sequence-to-sequence loss function, enabling single-stage joint optimization of formulation and code generation.

C THE DETAILS OF THE FIVE-ELEMENT FORMULATION

The five-element modeling formulation is a universal mathematical model for optimization problems, which consists of five parts. We start from the following formulation:

$$\min_{\boldsymbol{x} \in \mathcal{X} \subset \mathbb{R}^D} f(\boldsymbol{x}), \quad \text{s.t. } G(\boldsymbol{x}) \le \boldsymbol{c} , \qquad (2)$$

where $\boldsymbol{x}=(x_1,x_2,\ldots,x_D)^{\top}$ is the D-dimensional decision variable, \mathcal{X} is the feasible region, $f:\mathcal{X}\to\mathbb{R}$ is the objective, $G(\boldsymbol{x}):\mathbb{R}^D\to\mathbb{R}^m$ collects the constraints, and $\boldsymbol{c}\in\mathbb{R}^m$ provides the upper bounds. Among them, Variables, Objective, and Constraints correspond to $\boldsymbol{x}, f(\boldsymbol{x})$, and $G(\boldsymbol{x})$, while Sets and Parameters provide indices and numerical tables that instantiate and vectorize f and G. Sets determine the dimensions and naming of decision and constraint families. Parameters supply exogenous constants such as costs, coefficients, budgets, and demands. Variables specify domains and bounds(e.g., continuous, nonnegative, integer, or binary), which jointly define the feasible region \mathcal{X} ; domain-type restrictions such as "positive integers" may equivalently be encoded as explicit constraints, and our parser maps both styles to \mathcal{X} . Objective gives the minimization or maximization expression, and Constraints provide named families of linear or nonlinear equalities/inequalities composing G and the bound vector G. This representation naturally spans LP/IP/MILP: integrality arises through \mathcal{X} , linearity or nonlinearity is captured by the form of G, and multi-objective problems can be accommodated by extending f(x) to a vector F(x) with a scalarization scheme. The think output ends with the five-element formulation, which serves as a modeling blueprint.

D FORMULAS OF OPTREWARD

As described in Section 3.4.1, OptReward consists of three parts: format score S_{fmt} , five-element score S_{fmt} , and accuracy score S_{fmt} . Their detailed mathematical formulas are expressed as follows:

$$S_{fmt} = \begin{cases} 1, & \text{if format is fully correct,} \\ -1, & \text{otherwise.} \end{cases}$$
 (3)

$$S_{\text{five}} = \begin{cases} 0.2 \sum_{k=1}^{5} I_k, & \text{if } \sum_{k=1}^{5} I_k \ge 1, \\ -1, & \text{if } \sum_{k=1}^{5} I_k = 0, \end{cases}$$

$$(4)$$

where $I_k = 1$ if the k-th required element e_k is present in the <think> segment, and $I_k = 0$ otherwise. $(e_1, \ldots, e_5) = (Sets, Parameters, Variables, Objective, Constraints).$

$$S_{\rm acc} = \begin{cases} +2, & \text{if execution succeeds and } \hat{f} = f^{\star}, \\ -1.5, & \text{if execution succeeds but } \hat{f} \neq f^{\star}, \\ -2, & \text{if no executable code or execution fails .} \end{cases}$$
 (5)

$$R = \begin{cases} -4, & \text{if } S_{\text{fmt}} = -1, \\ S_{\text{fmt}} + S_{\text{five}} + S_{\text{acc}}, & \text{if } S_{\text{fmt}} = 1. \end{cases}$$

$$(6)$$

E DETAILS OF GRPO ALGORITHM

This paper proposes OptGRPO based on the improvement of the GRPO algorithm (Shao et al., 2024). For each query, GRPO sample G responses, compute the mean and standard deviation of their scalar rewards, and form a group-normalized advantage $\hat{A}_i = \frac{(r_i - \mu)}{\sigma}$. Specifically, GRPO optimizes the policy π_{θ} as follows:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \\
\left[\sum_{i=1}^G \sum_{t=1}^{|o_i|} \frac{1}{|o_i|} \left(\min \left[\frac{\pi_{\theta}(o_{i,t}|q, o_{i, < t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i, < t})} \hat{A}_{i,t}, \operatorname{clip}\left(\frac{\pi_{\theta}(o_{i,t}|q, o_{i, < t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i, < t})}, 1 - \varepsilon, 1 + \varepsilon\right) \hat{A}_{i,t} \right] \\
-\beta D_{\text{KL}}\left[\pi_{\theta}(\cdot|q) \| \pi_{\text{ref}}(\cdot|q) \right] \right) \right], \tag{7}$$

where q denotes queries sampled from the input dataset and o denotes the model's outputs. ε is the clipping threshold, β is the coefficient on the KL penalty, and $D_{\rm KL}$ is the KL divergence between the current policy π_{θ} and the reference policy $\pi_{\rm ref}$.

F THE PROCESSING PIPELINE OF TRAINING DATA

The processing pipeline for the SFT training set is detailed in Section 3.3.1. The resulting collection of this pipeline is named OptMATH-Train-Pyomo, which contains approximately 140K samples. The prompt template used in code conversion is introduced in Appendix M, and the prompt of the solver adapter is introduced in Appendix N.

In the first stage of MiniOpt's RL training, we employ a set of relatively easy and small-scale optimization problems, for which the accuracy score is more readily maximized. This, in turn, incentivizes the model to attain higher format score and five-element score, thereby accelerating mastery of the reasoning to model and solve paradigm. Concretely, the dataset of the first stage is the union of the NL4Opt (AhmadiTeshnizi et al., 2024) and ICML Competition (Yang et al., 2024b) training splits, comprising 1585 problems. Each instance is presented as a natural-language prompt paired with a reference answer and is fully compatible with our pipeline of parsing, executing, and scoring, enabling straightforward computation of the OptReward.

The second stage of RL targets optimization generalization (Jiang et al., 2025) under a limited training budget. The objective is to construct a training set that simultaneously covers diverse optimization types and scenarios, while preserving the scenario proportions observed in the real distribution. Starting from the OptMATH-Train pool containing 201K problems, we label each instance with types and scenarios using the DeepSeek-V3 (DeepSeek-AI et al., 2024), with prompt templates and data distributions provided in Appendix P. We then sample a 3000-instance stage-2 training set subject to two constraints: (i) type-uniform coverage, with exactly 600 instances per type; and (ii) within each type, the scenario frequencies match the distribution in the full pool. Formally, letting p(s) denote the overall scenario distribution in the complete pool, the target number of samples for each scenario s under type t is allocated as $n_t(s) = \text{round} \left(600 \cdot p(s)\right)$, where $\sum_s n_t(s) = 600$.

To prevent leakage, problems used during the SFT warm-up (Section 3.3.1) are excluded from consideration. The resulting 3K training set is uniform across types and aligned with the scenario distribution of the data pool. Compared with random or non-selected baselines, this selection improves overall performance of MiniOpt under the same compute budget.

G COMPARISON OF EXECUTION RATE ACROSS 9 BENCHMARKS

This section we use Execution Rate (ER), the proportion of generated code samples that run successfully without errors. As shown in the Table 5, MiniOpt also exhibits superior text generation capabilities compared to baseline methods, which suggests its excellent code generation performance given the problem modeling.

H ABLATION STUDY OF MINIOPT ON THE ER METRIC ACROSS 9 BENCHMARKS

Importance of the Training Pipeline of MiniOpt (Answer to Q3). In the training pipeline of MiniOpt, each module plays a distinct role in improving SA and ER. First, the lightweight SFT

Table 5: Comparison of the ER metric across 9 benchmarks with rankings (NL4Opt, ICML Competition, Mamo Easy, Mamo Complex, NLP4LP, ComplexOR, IndustryOR, OptiBench, OptMATH-Bench). **Bold** indicates 1st, wavy underline indicates 2nd, <u>underline</u> indicates 3rd. "Rank*" represents the result of sorting methods among parameter sizes below 10B.

Category	Model / Method	Avg.	Rank	Rank*	NL4Opt	ICML.C	Mamo.E	Mamo.C	NLP4LP	Com.OR	Indus.OR	OptiBench	OptMATH- Bench
	Qwen2.5-3B-Instruct	17.11	17	6	31.30	28.54	21.47	7.11	28.93	5.56	10.00	20.50	0.60
General Models	Qwen2.5-7B-Instruct	41.69	14	4	66.09	68.54	40.34	19.43	71.90	16.67	30.00	53.22	9.04
General Wodels	Qwen2.5-14B-Instruct	63.53	11	-	82.17	80.73	86.66	59.24	84.71	38.89	52.00	61.49	25.90
	DeepSeek-V3	83.50	7	-	97.83	97.32	96.63	71.56	97.93	72.22	70.00	84.79	63.25
	Qwen3-4B	14.15	18	7	19.57	20.00	15.18	18.48	18.18	11.11	5.00	16.20	3.61
	Qwen3-8B	25.55	16	5	36.96	38.29	26.53	18.01	39.26	22.22	10.00	35.04	3.61
General Models	Qwen3-14B	31.17	15	-	28.70	25.12	40.18	30.81	22.73	61.11	27.00	27.44	17.47
(Thinking)	DeepSeek-R1	83.07	9	-	96.09	94.15	89.72	84.83	91.32	61.11	82.00	88.92	53.01
	Gemini-2.5-Pro	89.65	4	-	94.35	96.10	95.86	86.26	96.28	77.78	87.00	91.90	81.32
	GPT-5	83.26	8	-	98.08	97.32	71.47	55.92	97.52	88.89	78.00	92.07	69.88
Prompt-based	Chain-of-Experts	61.72	12	-	79.57	73.41	72.85	56.87	77.27	55.56	54.00	65.45	20.48
Methods	OptiMUS	52.13	13	-	45.22	75.85	73.77	43.60	44.63	44.44	48.00	65.95	27.71
Wichiods	Reflexion	80.42	10	-	91.74	91.46	97.55	67.77	95.45	83.33	66.00	81.65	48.80
Learning-based	OptMATH-7B	85.07	6	3	99.13	95.85	98.47	90.05	99.17	66.67	69.00	82.81	64.46
Models	LLMOPT-14B	90.03	<u>3</u>	- 7	97.42	93.90	92.29	77.73	97.93	88.89	61.00	73.22	31.93
	MiniOpt-3B	88.04	5	2	99.57	95.85	98.47	87.68	99.59	88.89	70.00	83.64	68.67
Ours	MiniOpt-7B	90.61	2	$\widetilde{1}$	99.57	98.05	98.93	95.26	100.00	88.89	74.00	84.30	76.51
	MiniOpt-14B	92.35	ĩ	-	99.57	98.54	98.93	97.16	100.00	88.89	77.00	89.09	81.93

Table 6: Ablation study (MiniOpt-3B) on the ER metric across 9 benchmarks.

Category	Model / Method	Avg.	NL4Opt	ICML.C	Mamo.E	Mamo.C	NLP4LP	Com.OR	Indus.OR	OptiBench	OptMATH-Bench
	MiniOpt-3B	88.04	99.57	95.85	98.47	87.68	99.59	88.89	70.00	83.64	68.67
	MiniOpt-3B w/o SFT Warm-up	83.66	99.13	96.59	98.16	84.36	98.35	66.67	68.00	84.46	57.23
	MiniOpt-3B w/o RL	75.36	96.09	92.68	97.09	69.19	97.93	55.56	62.00	73.39	34.34
Ablations	MiniOpt-3B w/o Two-stage RL	84.38	98.26	96.10	97.85	84.36	98.35	72.22	70.00	84.46	57.83
Adiations	MiniOpt-3B w/o Data Selection	82.04	99.57	95.61	96.78	86.73	98.76	55.56	66.00	82.15	57.23
	MiniOpt-3B w/ Random Selection	82.01	99.13	96.10	97.09	77.73	97.93	67.67	69.00	84.63	48.80
	MiniOpt-3B w/o OptReward	83.39	98.26	95.61	97.70	83.89	97.93	72.22	64.00	82.48	58.43
	MiniOpt-3B w/o OptGRPO	80.00	97.39	94.63	96.32	83.41	95.87	55.56	64.00	81.65	51.20

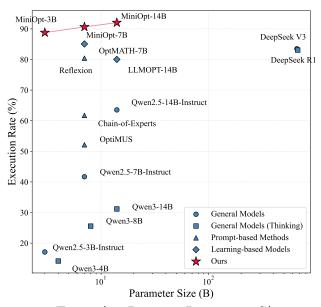
warm-up provides a better starting point for RL training. Without it, averages for SA and ER fall to 53.59% and 83.66%, respectively. The decreases are $\Delta SA=-3.35$ and $\Delta ER=-4.38$. Second, removing all RL training collapses the average performance from SA of 56.94% and ER of 88.04% to 39.25% and 75.36%, respectively. With the sharpest drops on challenging datasets like Com.OR $(\Delta SA = -33.33\%, \Delta ER = -33.33\%)$ and OptMATH-Bench $(\Delta SA = -22.29\%, \Delta ER = -34.33\%)$ where accurate modeling and solver selection are indispensable. Third, collapsing the two-stage RL removes the progressive training that first consolidates the paradigm (stage-1) and then targets generalization (stage-2), averages drop to the SA of 54.12% and ER of 84.38% (Δ SA=-2.82%, Δ ER=-3.66%). Moreover, the data selection of stage-2 is crucial for sample efficiency that training on the full pool (w/o Data Selection) yield the average SA of 53.58% and the average ER of 82.04% ($\Delta SA = -3.36\%$, $\Delta ER = -6.00\%$), while the random selected training data (w/ Random Selection) yield SA of 50.94% and ER of 82.01% on average (Δ SA = -6.00%, Δ ER = -6.03%), indicating that type-uniform, globally scenario-aligned sampling concentrates updates where they best improve cross-type, cross-scenario behavior. Finally, reverting our OptGRPO to the original GRPO (w/o OptGRPO) further decreases both metrics, averages SA of 48.33% and ER of 80.00%. The drop aligns with our algorithmic choices: removing KL part frees exploration for small-scale LLMs; Clip-Higher part prevents entropy collapse by allowing probability increases on rare but crucial reasoning/code tokens; Token-Level Loss part can enhance the impact of long output, which is conducive to the training of the reasoning model. Together these changes improve sample efficiency and training stability, which is critical for eliciting strong optimization generalization at small parameter scales.

Importance of the Reasoning to Model and Solve Paradigm and OptReward (Answer to Q4). The reasoning to model and solve paradigm and OptReward together turn free-form generation into a verifiable formulation, which are easy to verify. Removing them and keeping only a final-answer signal (w/o OptReward) drops averages to 52.44% and 83.39% (Δ SA=-4.50%, Δ ER=-4.65%), respectively. The largest losses appear where the problems are challenging: Com.OR (Δ SA = -16.67%, Δ ER = -16.67%), Indus.OR (Δ SA = -3.00%, Δ ER = -6.00%), and OptMATH-Bench (Δ SA = -4.82%, Δ ER = -10.24%). These patterns align with the roles of the three reward

components: the format score enforces the labels <think> and <answer> in the responses are complete; the five-element score shapes the intermediate blueprint, so the models learns to extract problem structure before coding; the Accuracy Score certifies numerical correctness by executing the Pyomo code and comparing the returned optimum with the reference. In combination, this reasoning to model and solve paradigm, together with verifiable reward, steers learning toward structurally consistent behaviors, yielding higher SA and ER across various optimization problems.

I COMPARISON OF AVERAGE ER AGAINST MODEL PARAMETER SCALES FOR VARIOUS METHODS

In this chapter, we present a comparative plot of the model parameters scale and the average ER across multiple methods, as illustrated in Figure 5.



Execution Rate vs. Parameter Size

Figure 5: Comparison of average ER against model parameter scales for various methods. MiniOpt is the Pareto optimal among compared methods.

J COMPARISON OF RESULTS BETWEEN THE MINIMUM MODEL AND MODELS OF OTHER SCALES

As shown in Tables 7 and 8, we compared the SA and ER of four different sizes of MiniOpt (1.5B, 3B, 7B, and 14B) across 9 benchmarks, respectively.

Table 7: Comparison of the SA metric between MiniOpt-1.5B and larger scale counterparts across 9 benchmarks.

Solvi	ng Accuracy (SA)	Avg.	NL4Opt	ICML.C	Mamo.E	Mamo.C	NLP4LP	Com.OR	Indus.OR	OptiBench	OptMATH-Bench
Datas	et size		230	410	652	211	242	18	100	605	166
Ours	MiniOpt-1.5B MiniOpt-3B MiniOpt-7B MiniOpt-14B	46.15 56.94 62.76 66.10	63.48 83.04 89.13 92.17	55.37 68.05 77.56 86.34	77.15 85.43 88.34 90.80	27.01 35.07 38.89 33.65	58.68 73.55 79.34 79.75	38.89 50.00 55.56 61.11	16.00 21.00 26.00 27.00	41.98 53.55 59.34 67.44	36.75 42.77 51.20 56.63

Table 8: Comparison of the ER metric between MiniOpt-1.5B and larger scale counterparts across 9 benchmarks.

Execu	ition Rate (ER)	Avg.	NL4Opt	ICML.C	Mamo.E	Mamo.C	NLP4LP	Com.OR	Indus.OR	OptiBench	OptMATH-Bench
Datas	et size	[230	410	652	211	242	18	100	605	166
Ours	MiniOpt-1.5B MiniOpt-3B MiniOpt-7B MiniOpt-14B	80.00 88.04 90.61 92.35	93.91 99.57 99.57 99.57	91.71 95.85 98.05 98.54	97.09 98.47 98.93 98.93	79.15 87.68 95.26 97.16	91.74 99.59 100.00 100.00	66.67 88.89 88.89 88.89	63.00 70.00 74.00 77.00	75.87 83.64 84.30 89.09	60.84 68.67 76.51 81.93

K THE SEESAW ISSUE OF LLMS

We assess whether adapting our models to optimization modeling introduces a seesaw issue for small-parameter models. After evaluating pre-training and post-training scores on six widely used general-purpose benchmarks: MMLU (Hendrycks et al., 2021a), MATH (Hendrycks et al., 2021b), HumanEval (Chen et al., 2021), TriviaQA (Joshi et al., 2017), RACE (Lai et al., 2017), GSM8K (Cobbe et al., 2021), we observe that the average scores of MiniOpt-7B and MiniOpt-3B decrease by only 1.83% and 3.05%, respectively. Notably, the 3B model even exhibits a 0.7% improvement on GSM8K, a proxy for mathematical reasoning. These results indicate that our training paradigm for optimization modeling does not induce a pronounced seesaw effect. We attribute this to two factors. First, the two-stage RL framework coupled with a verifiable OptReward constrains learning to structural modeling correctness and executable solving efficacy, mitigating overfitting to superficial linguistic style or lengthy chain-of-thought and thereby substantially reducing the risks of catastrophic forgetting and cross-task seesaw effects. Second, the foundational competencies required for optimization modeling, such as mathematical understanding, symbolic reasoning, program synthesis, and execution, which highly overlap with those assessed by general benchmarks such as MMLU, MATH, GSM8K, and HumanEval. Consequently, targeted reinforcement in this domain does not overwrite existing representations, instead, it yields small positive transfer on tasks closely aligned with modeling and solving, as exemplified by GSM8K.

L DETAILED DISCUSSION ON THE EFFICIENCY OF MINIOPT TO EODEL AND SOLVE OPTIMIZATION PROBLEMS

To validate the high efficiency of the MiniOpt model during inference, we compared the average number of tokens generated in full responses on the OptMATH dataset by two small-parameter models of MiniOpt and DeepSeek-R1. As shown in Table 9, MiniOpt achieves higher ER and SA while using fewer average output tokens. This advantage stems from MiniOpt's internalized unified reasoning to model and solve paradigm tailored for optimization modeling problems, which guides the model along a more efficient reasoning path during inference.

Table 9: Comparison of average output token consumption in OptMATH-Bench.

Model	Avg. Token Count	SA (%)	ER (%)
MiniOpt-3B	1911.72	42.77	68.67
MiniOpt-7B	3068.99	51.20	76.51
DeepSeek-R1	5078.90	39.76	53.01

M SYSTEM PROMPT FOR CODE CONVERSION FROM GUROBIPY TO PYOMO

This section provides the system prompt for the large language model to convert GurobiPy code in OptMATH-Train into Pyomo code, where the content within $[\cdot]$ and $\{\cdot\}$ will be replaced with the corresponding parts.

PROMPT TEMPLATE FOR CODE CONVERSION You are an expert in optimization problems. Your task is to convert the given gurobipy code into pyomo code. **Instructions:** 1. Don't give any explanation, just provide the converted pyomo code in the following format: '''python [pyomo code here] 2. Please note that the following solvers are available for use: ' glpk', 'cbc', 'ipopt', 'scip'. Other solvers should not be utilized. 3. Please add 'from pyomo.environ import *' at the beginning of your code. 4. Please print the optimal objective value at the end of the code. **Gurobipy code: ** {gurobipy}

N PROMPT FOR THE SOLVER SELECTION

This section provides the system prompt to guide MiniOpt models in autonomously selecting solvers after modeling optimization problems.

Solver Selection Guide: - ''glpk'': Best for small-to-medium linear problems (LP). - ''cbc'': Recommended for mixed-integer linear programming (MILP) and larger linear problems. Handles binary/integer variables well. - ''ipopt'': Use for nonlinear problems (NLP) with continuous variables. Does NOT support discrete variables. - ''scip'': Most versatile - handles mixed-integer nonlinear problems (MINLP), large-scale problems, and complex constraints. **Select solver based on:** 1. Variable types (continuous vs integer/binary) 2. Linearity of objective/constraints 3. Problem scale (small: glpk/cbc, large: scip/ipopt) 4. Nonlinearity presence (use ipopt/scip)

O SYSTEM PROMPT FOR RL TRAINING

This section provides the system prompt used by MiniOpt models during reinforcement learning (RL) training.

SYSTEM PROMPT FOR RL TRAINING

You are a helpful assistant. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within <think> </think> and <answer> </answer> tags, respectively, i.e ., <think> reasoning process here </think><answer> answer here

```
1243
                       </answer>, please make sure to answer according to the above
                       format. Now the user asks you to solve an optimization reasoning
1244
                        problem, you should:
1245
                 1. Detailed reasoning about the problem within <think> </think>
1246
                       tags.
1247
                 2. Write the corresponding five-element model (derived from your
1248
                       analysis).
                 3. Determine the mathematical properties of problem and select an
1249
                       appropriate solver from 'glpk', 'cbc', 'ipopt', 'scip'.
1250
                 4. Recheck and correct if necessary at the end of the <think> </
1251
                       think> section.
1252
                      - Verify the five-element model fully captures the problem's
1253
                       requirements.
                      - Confirm no constraints/variables are missing or over-
1254
                       simplified.
1255
                      - Ensure the solver choice aligns with the problem's
1256
                      mathematical properties.
1257
                 5. Provide the corresponding Pyomo code based on checked five-
                       element model within <answer> </answer> tags.
1258
1259
                 In mathematics, optimization problem can be modeled as the
1260
                       following expression {\min_{{\{}\setminus boldsymbol_{x}} \setminus mathcal_{\{}}}
1261
                       X}}} f(\boldsymbol{{x}}), {{\rm s.t.}} G(\boldsymbol{{x}})
1262
                       \left(c\right)\, where \left(x\right)\ = \left(x_1, x_2, x_3\right)
                       \label{eq:local_condition} \label{eq:local_con
1263
                       \mathcal{X}} \subset \mathcal{X} \subset \mathbb{{R}}^D$ is the feasible domain,
1264
                         f: \mathbb{X}} \right)  is the objective
1265
                         function and the goal is to find the minima of f, G(\
1266
                       boldsymbol{x}) \leq boldsymbol{c} are the constraints of
1267
                         \ \\boldsymbol{{x}}\$.
1268
                The above definition can be mapped to a five-element consisting of
1269
                        "Variables, Objective, Constraints, Sets, Parameters".
1270
                       Variables indicate what \boldsymbol{x} boldsymbol{x} is, Objective
1271
                       describes the form of the objective function f(\b x) = 0
1272
                       }})$, and Constraints indicates the constraints $G(\\boldsymbol
                       \{\{x\}\}\ and \{X\}\. These three can abstract the
1273
                       optimization problem. Sets and Parameters are their specific
1274
                       explanations: Sets describe and explain the subscripts of the
1275
                       vectors or matrices in them, and Parameters supplement their
1276
                       specific values.
1277
                You need to give a detailed reasoning process for the problem first
1278
                        , and then write the corresponding five-element model based on
1279
                       the problem description and information provided by user.
1280
1281
                Please complete the following template to model the optimization
1282
                       problem into five-element:
1283
                 <t.hink>
1284
                Your reasoning process here...
1285
1286
                 ## Sets:
1287
                 [You need to fill in]
1288
                 ## Parameters:
1289
                 [You need to fill in]
1290
1291
                 ## Variables:
1292
                 [You need to fill in]
1293
                 ## Objective:
1294
                 [You need to fill in]
1295
```

```
1296
1297
          ## Constraints:
1298
          [You need to fill in]
1299
          </think>
1300
1301
         In Pyomo, all constraints must be formulated using '<=', '>=', or
             '=='. If you need to use '>' or '<', you can introduce a very
1302
             small value to transform the inequality. Please note that the
1303
             following solvers are available for use: 'glpk', 'cbc', 'ipopt',
1304
              'scip'. Other solvers should not be utilized.
1305
1306
         **Solver Selection Guide:**
         - ''glpk'': Best for small-to-medium linear problems (LP).
1307
         - ''cbc'': Recommended for mixed-integer linear programming (MILP)
             and larger linear problems. Handles binary/integer variables
1309
1310
           ''ipopt'': Use for nonlinear problems (NLP) with continuous
1311
             variables. Does NOT support discrete variables.
           "scip": Most versatile - handles mixed-integer nonlinear
1312
             problems (MINLP), large-scale problems, and complex constraints.
1313
1314
          **Select solver based on: **
1315
         1. Variable types (continuous vs integer/binary)
1316
         2. Linearity of objective/constraints
1317
         3. Problem scale (small: glpk/cbc, large: scip/ipopt)
         4. Nonlinearity presence (use ipopt/scip)
1318
1319
         Please select an appropriate solver based on the type and quantity
1320
             of variables, objectives, and constraints. After thinking, when
1321
             you finally reach the five-element model, you should give the
1322
             corresponding Pyomo code within the <answer> </answer> tags, i.e
             ., <answer> '''python\n code here''' </answer>. The user will
1323
             extract the complete code you provide through the regular
1324
             expression r" ''python\n(.*?)''" in the <answer> </answer> tags
1325
             . The execution result of the code should include the optimal
1326
             solution and the objective value. The optimal objective value
             will be extracted automatically from your last printed result.
1327
1328
```

P $\;\;$ Labeling prompt and the data distributions of ${\sf OptMATH-Train}$

This section provides the system prompts for the large language model to label the type and scenario of problems in OptMATH-Train, where $\{\{\cdot\}\}$ will be replaced with the corresponding content. After labeling, the distribution of scenarios in OptMATH-Train is displayed in Figure 6, and the distribution of types in OptMATH-Train is displayed in Figure 7

System Prompt for type labeling

1329 1330 1331

1332

1334

1335

1336

1337 1338

133913401341

1342

1343 1344

1345

1346

1347

1348

1349

Please classify the following optimization problem into one of these technical types based on the mathematical formulation and decision variables, not just surface-level descriptions:

- Linear Programming (LP): Problems with linear objective function and linear constraints, all continuous variables
- Integer Programming (IP): Problems with linear or nonlinear components where ALL variables are discrete/integer
- 3. Mixed Integer Linear Programming (MILP): Problems with linear components containing BOTH continuous and discrete variables

```
1350
1351
         4. Nonlinear Programming (NLP): Problems with nonlinear objective
             function and/or nonlinear constraints (variables may be
1352
             continuous/discrete)
1353
         5. Combinatorial Optimization (CO): Problems focused on selecting/
1354
             discrete structures (graphs, permutations, sets) with typically
1355
             binary variables
1356
         6. Multi-objective Programming (MOP): Problems explicitly
             optimizing multiple conflicting objectives simultaneously
1357
         7. Second-Order Cone Programming (SOCP): Problems with a linear
1358
             objective function, linear constraints, and second-order cone
1359
             constraints (e.g., (\|Ax + b\| \ c^T x + d))
1360
1361
          # Problem:
         {{Question}}
1362
1363
          # Output
1364
         Analyze the mathematical structure step by step and classify its
1365
             type. Finally, output the type abbreviation in the following
1366
             format:
         Type: Abbreviation of the type
1367
1368
         Note:
1369
         - Focus on the fundamental mathematical formulation, not
1370
             application domain
          - Check variable types (continuous/discrete/binary) and objective/
1371
             constraint linearity
1372
         - For MOP, there must be explicit multiple objectives
1373
         - For pure discrete problems with special structures (e.g. graphs),
1374
              prefer CO over IP
1375
1376
1377
         System Prompt for Scenario Labeling
1378
1379
         Please classify the following optimization problem into one of
1380
             these application domains based on the core decision-making
1381
             context and primary business function, not just keywords
             mentioned in the problem:
1382
1383
         1. Supply Chain: Decisions about inventory management, distribution
1384
              network, warehousing operations
1385
```

- Finance: Decisions about portfolio management, investments, risk management, financial planning
- Manufacturing: Decisions about production processes, quality control, factory operations

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1400

1401

1402

- 4. Transportation: Decisions about routing, vehicle scheduling, fleet management, traffic flow, carrier selection
- Healthcare: Decisions about medical staff scheduling, patient flow, hospital resources
- Energy: Decisions about power generation, energy conservation, grid distribution
- Technology: Decisions about network design, data center operations, cloud resources
- Retail: Decisions about store operations, pricing, inventory, equipment, store layout
- 9. Agriculture: Decisions about farming operations, crop planning, irrigation
- Logistics: Decisions about delivery operations, warehouse management, distribution
- Resources: Decisions about raw materials, equipment allocation, material management
- 12. Marketing: Decisions about campaign planning, budget allocation , target selection

```
1404
1405
          13. Education: Decisions about course scheduling, resource
             allocation in schools
1406
         14. Environment: Decisions about environmental protection,
1407
             emissions control, conservation
1408
         15. Construction: Decisions about project planning, construction
1409
             resource allocation
1410
         16. Military: Decisions about military operations, deployment,
             supply management
1411
          17. Sports: Decisions about game scheduling, team formation,
1412
             strategy
1413
         18. Telecommunications: Decisions about network coverage, bandwidth
1414
              allocation
1415
         19. Aviation: Decisions about flight scheduling, crew assignment,
             airport operations
1416
          20. Services: Decisions about service operations, staff scheduling,
1417
              capacity management
1418
          21. Public utilities: Decisions about utility services,
1419
             infrastructure management, service delivery
         22. Other: Problems that don't clearly fit into above categories
1420
1421
          # Problem:
1422
          {{Question}}
1423
1424
          # Output
1425
         Let's think step by step, give the analysis of the problem and
             classify it into one of the above application domains. Finally,
1426
             output the name of the domain in the following format:
1427
         Category: Name of the Domain
1428
1429
1430
          - Focus on the fundamental business function and decision-making
             context
1431

    Don't be misled by secondary keywords or background story

1432
         - Consider who is making the decision and what is their primary
1433
             business purpose
1434
1435
```

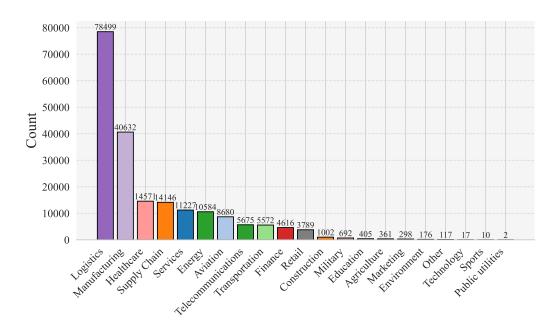


Figure 6: Proportion of every scenarios of instances in OptMATH-Train (201K).

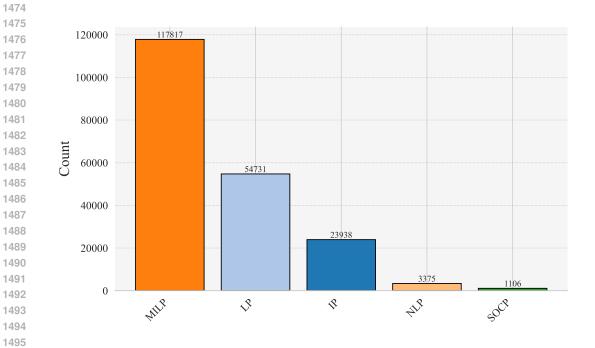


Figure 7: Proportion of every problem types of instances in OptMATH-Train (201K).