

# Solver-Informed RL: Grounding Large Language Models for Authentic Optimization Modeling

Yitian Chen<sup>1\*</sup>, Jingfan Xia<sup>1,2\*</sup>, Siyu Shao<sup>1,3</sup>, Dongdong Ge<sup>4†</sup>, Yinyu Ye<sup>4,5</sup>

<sup>1</sup> Cardinal Operations, China

<sup>2</sup> Shanghai University of Finance and Economics

<sup>3</sup> The University of Hong Kong

<sup>4</sup> Antai School of Economics and Management, Shanghai Jiao Tong University

<sup>5</sup> Department of Management Science and Engineering, Stanford University

chenyitian@shanshu.ai, jf.xia@163.sufe.edu.cn,

siyu\_shao@connect.hku.hk, dedge@sjtu.edu.cn, yyye@stanford.edu

## Abstract

Optimization modeling is fundamental to decision-making across diverse domains. Despite progress in automating optimization formulation from natural language descriptions, Large Language Models (LLMs) often struggle to generate formally correct and usable models against hallucinations, posing a challenge for reliable automation. Inspired by the success of Reinforcement Learning (RL) in enhancing Large Reasoning Models, we present Solver-Informed Reinforcement Learning (SIRL), a novel framework that significantly improves the authenticity of LLMs for optimization modeling using Reinforcement Learning with Verifiable Reward by leveraging external optimization solvers as verifiers. These verifiers automatically assess the executable code and the instance-level mathematical model represented by the associated LP file, yielding precise and comprehensive feedback signals—including syntax, feasibility, and solution quality, serving as direct rewards for the RL process. This automated verification process, particularly from classic optimization solvers, also underpins our instance-enhanced self-consistency method to synthesize high-quality training data. Extensive experiments on diverse public benchmarks demonstrate that SIRL achieves state-of-the-art performance, substantially outperforming existing methods in generating accurate and executable optimization models. Our code is publicly available at <https://github.com/Cardinal-Operations/SIRL>.

## 1 Introduction

Optimization modeling provides a powerful framework for decision-making across diverse fields, from logistics and finance to engineering and machine learning [1, 2, 3, 4, 5]. Despite the maturity and powerful capabilities of modern optimization solvers, such as Gurobi [6], COPT [7] and CPLEX [8], translating complex real-world problems into precise mathematical models and executable optimization codes remains a significant bottleneck, often requiring substantial domain expertise and manual effort [2, 9].

The advent of Large Language Models (LLMs), (e.g., GPTs [10], Gemini [11, 12], Deepseek [13]), offers a promising avenue to automate or assist in this intricate mathematical modeling and code generation process, potentially democratizing access to optimization solvers. Yet, ensuring the

\*Both authors contribute equally to this research.

†Corresponding authors.

correctness, feasibility, and solver-compatibility of LLM-generated optimization models presents a significant and increasingly active research challenge. In general, existing approaches that leverage LLMs for optimization modeling can fall broadly into two categories. Firstly, prompt-based or agent-based approaches utilize the frozen capabilities of powerful foundation LLMs [10, 12, 11, 13], guide the models to extract relevant problem information, generate the corresponding mathematical models and executable code for optimization solvers [14, 15, 16, 17]. While accessible, these methods do not adapt the underlying model parameters and can be sensitive to prompt design and the chosen foundation LLMs. A distinct line of work addresses training open-source LLMs [18, 19] to enhance their capabilities for optimization modeling using offline learning approaches. This encompasses techniques such as Supervised Fine-Tuning (SFT) on demonstrations data, and alignment techniques (e.g., DPO [20], KTO [21]) on preference data. The success of existing offline methods [22, 23, 24] depends on carefully curated datasets. These curated datasets, comprising problem descriptions paired with either detailed mathematical model and code demonstrations or comparative human preference labels, are typically generated through human annotation [24] or synthesis [22, 23]. Training on these curated datasets enables offline learning approaches to capture the data’s stylistic and structural patterns, such as mathematical formulation and solver code, and achieve good performance. However, since their training objective focuses on mimicking demonstrations or aligning preferences, these methods still struggle to inherently guarantee functional correctness or solution feasibility required for reliable execution by solvers.

Recently, significant progress has been made in enhancing LLMs’ reasoning capabilities, driven by advances in scale and techniques such as Chain-of-Thought (CoT) prompting [25]. Building upon this, reinforcement learning (RL) techniques [26, 27, 28, 29] have further equipped LLMs to tackle complex tasks requiring sophisticated, multi-step reasoning processes, leading to the development of powerful Large Reasoning Models (LRMs) [30, 31, 32], whose capabilities are exemplified by their high performance on challenging mathematical problems (e.g., at the IMO level) [31, 30]. A key technique enabling these advances in LRMs, particularly for tasks like mathematical reasoning and coding, is reinforcement learning with verifiable rewards (RLVR) [31, 33]. RLVR directly optimizes LLMs’ policies using objective feedback from verified outcomes, such as the execution of generated code against test cases or the validation of mathematical solutions [33, 31, 32, 34, 35].

For optimization tasks, solving a real-world problem involves a multistep process encompassing problem analysis and reasoning, mathematical modeling, and code implementation [14, 15, 16, 17, 22, 23, 24], resembling a CoT reasoning process. The LLM output, comprising the mathematical model and the solver code, is verifiable using external optimization solvers. Verification involves steps such as syntax checks, feasibility assessment through model solving, and comparison of the objective value against known optimal value. These verifiable checks provide rich objective reward signals, enabling RLVR [36, 31, 33, 37] to directly optimize the LLM generation towards producing correct, feasible, and high-quality optimization output. This creates a powerful synergy between LLM reasoning and objective verification.

To the best of our knowledge, this is the first application of RLVR to directly enhance LLMs’ proficiency in optimization modeling. The rewards, including feasibility status, objective function value, and mathematical model statistics from the LP file, are obtained by executing the generated code. Specifically, the LP file, which is a standard format for optimization solvers, provides an explicit instance-level representation of the generated mathematical model. This enables a solver-informed evaluation that ensures accurate assessment of the LLM’s performance and validity on both mathematical modeling and outcome correctness.

Our main contributions are fourfold: (1) We introduce a simple yet effective instance-enhanced self-consistency method for synthesizing high-quality training data for optimization tasks. (2) We introduce SIRL, an automated RLVR framework for LLMs in optimization modeling with a novel surrogate function. By enabling a balance between diverse reasoning exploration and the requirements for accuracy and validity in mathematical models and code, this function leads to a significant improvement in the authenticity and reliability of the generated optimization solutions. (3) We demonstrate how classical optimization solvers can serve as effective and powerful tools for both enhancing the data synthesis process and providing rich reward signals for the proposed SIRL framework. (4) Through extensive experiments on diverse public benchmarks, our 7B-parameter model, trained by the proposed SIRL framework, achieves state-of-the-art performance, significantly outperforming other offline learning and agent-based methods in generating correct and reliable models.

## 2 Related work

Our work builds upon and contributes to several research areas, primarily LLMs for optimization, synthetic data generation for LLM training, the paradigm of reinforcement learning with tool-verified feedback mechanisms.

**LLMs for optimization modeling.** The application of LLMs for optimization modeling has emerged as a prominent research direction. Early work mainly relied on prompt engineering techniques, including agent-based prompting [16] and multi-agent reasoning [15, 17], but was limited by careful prompt design and the capability of the foundation LLMs. More recent work focuses on offline learning approaches to adapt LLM parameters using specialized mathematical modeling datasets. For instance, ORLM [22] and OptMATH [23] employed Supervised Fine-Tuning (SFT) with datasets constructed via semi-automated synthetic data generation workflows, LLMOPT [24] introduced an alignment learning framework with multi-instruction tuning and a self-correction mechanism.

**LLM-Based data synthesis.** Fine-tuning LLMs on specialized tasks necessitates high-quality datasets, a resource-intensive requirement often necessitating domain expertise. Data synthesis offers a scalable solution. Examples of general-purpose synthesis techniques include self-instruction [38] and the WizardInstruct series [39, 40, 41]. Specifically, to improve reasoning capabilities, recent work focuses on synthesizing reasoning trajectory data, exemplified by the work Star [42], rstar [43], Phi [44], aiming to enhance complex reasoning ability in tasks such as math and code. Within the domain of optimization modeling, several approaches have been explored: ORLM [22] introduced a semi-automated method for synthesizing operations research datasets based on self-instruct frameworks [38]; ReSocratic [45] employs reverse data synthesis via back-translation; and OptMATH [23] generates controllable complexity data with verification.

Ensuring correctness and formal adherence is a central challenge when synthesizing high-quality data for tasks requiring rigorous output. Techniques addressing this include LLM-as-a-judge [46], which leverages foundation models for evaluation, rejection sampling [47, 23], which provides automated verification, and self-consistency [38], which relies on multi-sample voting for robustness.

**Reinforcement learning with verifiable reward.** Reinforcement Learning from Human Feedback (RLHF) [28] marked a significant step in aligning LLMs, typically using reward models trained on human preferences [48, 49]. However, for tasks where the desired output has objectively verifiable properties, such as mathematical reasoning [50, 49, 51] and code generation [52, 53], relying solely on subjective human preference is suboptimal and can struggle with objective correctness and potential reward hacking [54, 55, 56]. For example, a model might generate code with correct syntax but a flawed function. RLVR has demonstrated significant success in enhancing LLMs’ performance across these domains, as evidenced by strong results on benchmarks like GSM8K [49], MATH [52], AIMO [51], HumanEval [52], CodeForce [53], and has been a key technique in developing highly capable large reasoning models (e.g., OpenAI-O1 [30], DeepSeek-R1 [31], Kimi-k1.5 [32]) and state-of-the-art coding models (e.g., Tulu-3 [33]).

**External tools as verifier.** Leveraging external tools or formal verification mechanisms to validate LLM-generated structured outputs, particularly in domains that require high fidelity and correctness, is becoming an increasingly critical area of research [57, 37, 58]. In mathematical theorem proving, the Lean proof assistant [59] acts as a verifier for LLM-generated proofs [58, 57, 60, 61]. Google’s AlphaGeometry series [62, 63] combines reasoning LLMs with symbolic engines DDAR, achieving breakthroughs in the IMO-level mathematical tasks. For formal domains like mathematics and programming that demand high output fidelity, integrating code interpreters is commonly employed to improve LLMs’ performance and ensure correctness. The basic approach involves generating complete solutions and validating them post-hoc via external execution using interpreters or compilers [52, 64]. Feedback from these validation checks provides a signal for model refinement. More advanced methods either integrate tool interactions directly within the LLMs’ reasoning process (e.g., ToRA [65], MARIO [66]) or focus on enabling the LLMs to learn effective strategies for using tools autonomously (e.g., START [34], ReTool [35]).

Our work proposes a direct analogy: just as Lean verifies mathematical proofs and code compilers verify code for general math and coding problems, classical optimization solvers [6, 7, 8, 67] serve as the natural, powerful, and domain-specific objective oracles, which can both enhance the synthesis of high-quality optimization task data and provide rich reward signals for the proposed RLVR framework.

### 3 Method

Let the training dataset be  $\mathcal{D} = \{(x_i, y_i^*)\}_{i=1}^N$ , where  $x_i$  represents the natural language description of the  $i$ -th optimization problem and  $y_i^*$  is the corresponding ground-truth optimal objective function value. We model the problem solver as an LLM policy  $\pi_\theta$ , parameterized by  $\theta$ . Given an input problem description  $x$ , the policy  $\pi_\theta$  generates a response  $\mathbf{z}$  containing sequences of reasoning process leading to an objective function value  $y$ , derived via a mapping function  $g(x, \mathbf{z})$ .

To guide the learning process, we introduce the SIRL framework, which incorporates a verifiable reward function  $r(x, \mathbf{z}, y^*)$ . This function quantifies the quality of the derived objective value  $y$  for the problem  $x$ , using the ground truth  $y^*$  as a reference. Our goal is to optimize the policy parameter  $\theta$  to maximize the expected reward:

$$\max_{\theta} \mathbb{E}_{(x, y^*) \sim \mathcal{D}, \mathbf{z} \sim \pi_\theta(\cdot | x), y \sim g(x, \mathbf{z})} [r(x, \mathbf{z}, y^*)]. \quad (1)$$

In the following subsection, we outline the key components of our framework: the data synthesis pipeline used to construct the training dataset, the solver-informed reinforcement learning method with its surrogate function design tailored for optimization tasks and the two-stage training curriculum including reward design.

#### 3.1 Data synthesis framework

**Overall framework.** The data synthesis process starts from curated seed data and integrates LLM generation with subsequent steps including self-instruction, augmentation, evaluation, consistency assessment, and filtering. While our work primarily follows the OR-Instruct pipeline in ORLM [22], this approach contrasts with previous work synthesizing complete (question, model, code) sequences [22, 23, 68] by focusing on generating high-quality (question, answer) pairs, where the answer is the optimal value by executing an optimization solver.

As illustrated in Figure 1, we firstly sample questions from the curated seed data, combine them contextually with scenarios sampled from a predefined list, and generate new problems that remain structurally similar to the original. Subsequently, we implement an augmentation phase to increase the challenge related to semantic interpretation, mathematical modeling complexity, and solution complexity, obtaining a larger corpus of extended questions. Then, the LLM-as-a-judge approach [46] validates the generated problems for practical relevance and semantic consistency. Following the validation of the problem, the LLM is employed to generate mathematical models and corresponding executable code. This code is then executed to produce the objective value, feasibility status and the instance-level mathematical models represented by LP files.

To ensure high correctness of the final answer, we assign multiple LLM roles (10 roles) per problem [69, 70] and apply a novel instance-enhanced self-consistency framework when generating answers. Furthermore, an iterative reflection and refinement process [71] is employed to address execution issues, regenerating or refining code upon errors, and regenerating the model and code when infeasible solutions are encountered.

Finally, guided by the principle ‘‘Less is More’’ [72, 73], we filter the (question, answer) pairs. Samples are excluded if a baseline Qwen-32B-Instruct model [19] achieves an 80% pass rate (8/10 attempts across different roles) in generating executable code matching the optimal value, as these instances are considered too elementary. The retained pairs are incorporated into the final training dataset.

In the next part, we detail our novel, simple yet effective instance-enhanced self-consistency method used within our data synthesis framework.

**Instance-enhanced self-consistency.** Relying solely on majority voting of final results for self-consistency in optimization modeling can be limiting, potentially ignoring embedded model information. We enhance this by integrating structural data extracted from the instance’s LP file. The LP files are chosen as they formally encode key model properties (e.g., variable types, objective direction), providing a formalized, implementation-agnostic representation of the instance-level mathematical model. Specifically, after executing the generated code associated with a role  $r$  and obtaining the corresponding LP file, we extract the following features:

- $O_r$ : The final objective function value.

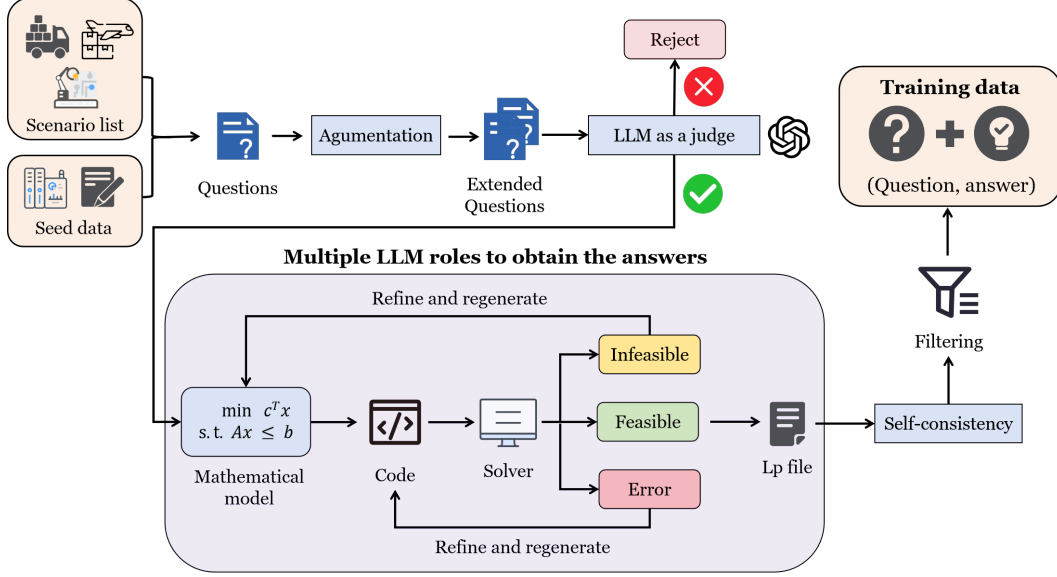


Figure 1: An overview of the data synthesis process.

- $D_r \in \{\max, \min\}$ : The optimization direction (maximization or minimization).
- $N_{bin,r}$ : The count of binary (0-1) variables.
- $N_{int,r}$ : The count of general integer variables (distinct from binary variables).

These statistics provide detailed, model-level insights that supplement the final numerical outcome  $O_r$ . Let  $R$  be the set of roles that generated responses for a given question. We assign a score  $S(r)$  to the response of each role  $r \in R$  using a weighted voting mechanism that measures consensus with other roles. We define a consensus function  $\psi(X_r)$  for a feature  $X \in \{O, D, N_{bin}, N_{int}\}$  as the count of the roles  $r' \in R$  whose corresponding feature value  $X_{r'}$  is identical to  $X_r$ :

$$\psi(X_r) = |\{r' \in R \mid X_{r'} = X_r\}|. \quad (2)$$

The final score  $S(r)$  for the response from the role  $r$  is calculated as a weighted sum reflecting consensus across the extracted features:

$$S(r) = w_1 \cdot \sqrt{\psi(O_r)} + w_2 \cdot \sqrt{\psi(D_r)} + w_3 \cdot \sqrt{\psi(N_{bin,r})} + w_4 \cdot \sqrt{\psi(N_{int,r})}. \quad (3)$$

The weights  $(w_1, w_2, w_3, w_4)$  determine the relative contribution of the individual consensus components to the final result. In our current implementation, all weights are set to 1, giving equal importance to each consensus component.

Figure 2 provides an illustrative example of this enhanced self-consistency method. For a given question, multiple LLM roles generate “Math model + Code” trajectories. Each code is executed by a solver, producing a solution and the corresponding LP file. From the LP file, we extract structural features (objective value consensus  $\psi(O_r)$ , optimization direction  $\psi(D_r)$ , variable counts  $\psi(N_{bin,r}), \psi(N_{int,r})$ ). The final score is calculated according to Equation 3. Finally, the objective value  $O_{r^*}$  from the response achieving the highest score, where  $r^* = \arg \max_{r \in R} S(r)$ , is selected as the definitive answer to the question.

### 3.2 SIRL: Solver-Informed Reinforcement Learning

**RLVR for LLMs in optimization modeling.** In optimization tasks, to obtain the optimal value from a problem description  $x$ , a complex reasoning process is involved with distinct stages: first, analyzing the problem description  $x$  to identify the key information, such as optimization problem type and its core components; second, constructing the mathematical formulation which typically involves the parameter set, objective functions, and constraints; and finally, generating the corresponding executable code, which is then executed by the optimization solver to produce the objective value  $y$

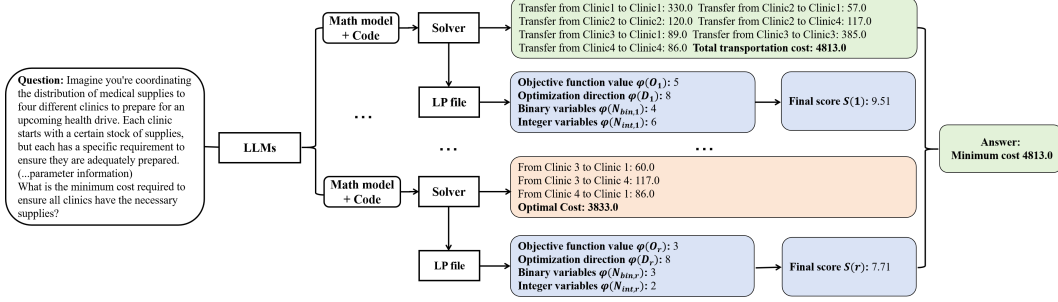


Figure 2: An example of the instance-enhanced self-consistency method.

and other relevant output (e.g., decision variable values and solution status). To address this complex reasoning process that integrates both mathematical modeling and code implementation, we utilize the Chain of Thought (CoT) [25] method, where the LLM policy  $\pi_\theta$  generates a sequence of intermediate thoughts  $\mathbf{z} = (\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^m)$  that serve as a crucial bridge between the initial problem  $x$  and the final result  $y$ .

Specifically, a well-designed system prompt structures the sequence of thoughts  $\mathbf{z}$  into segments reflecting the defined reasoning, modeling, and code generation stages:  $(\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^{m-2})$  contains the analysis and reasoning process (e.g., the identification of optimization problem type, algorithmic choices, or reasoning steps towards the final model structure);  $\mathbf{z}^{m-1}$  contains the concise mathematical modeling; and  $\mathbf{z}^m$  contains the derived executable code. The final value  $y$  is obtained deterministically by extracting and executing code in  $\mathbf{z}^m$ , formally represented as  $y = g(x, \mathbf{z})$ , where  $g$  denotes the deterministic code execution function.

At the token level, each thought  $\mathbf{z}^j$  is realized as a sequence of tokens  $\mathbf{z}^j = (z_1^j, \dots, z_{T_j}^j)$ . The token  $z_t^j$  within this sequence is sampled autoregressively from the model’s policy  $\pi_\theta(\cdot | x, \mathbf{z}^1, \dots, \mathbf{z}^{j-1}, z_1^j, \dots, z_{t-1}^j)$ , conditioned on the initial input  $x$ , all previously completed thoughts  $(\mathbf{z}^1 \dots \mathbf{z}^{j-1})$ , and all tokens generated so far in the current thought.

**Surrogate function design: Partial KL.** To maximize the expected verifiable reward objective defined in Equation (1), we employ REINFORCE++ [74], a robust policy gradient algorithm that incorporates key techniques from Proximal Policy Optimization [75].

In each training iteration, a batch of data  $\{x_i, y_i^*\}_{i=1}^n$  is sampled from the training dataset. Then, for each  $x_i$ , the policy  $\pi_\theta$  is used to sample a set of  $K$  complete response trajectories  $\{\mathbf{z}_{i,k}\}_{k=1}^K$ , where each  $\mathbf{z}_{i,k}$  is a sequence of tokens generated autoregressively and composed of a sequence of thoughts. For simplicity, we denote the collected batch as  $\mathcal{B}$ , where each tuple  $(x, \mathbf{z}, y^*) \in \mathcal{B}$  consists of  $x$ , the input problem description,  $\mathbf{z}$ , the generated response trajectory, and  $y^*$ , the ground-truth objective value.

The algorithm updates the policy parameters  $\theta$  by maximizing the following clipped surrogate objective:

$$\mathcal{J}^{\text{Reinforce++}}(\theta) = \frac{1}{|\mathcal{B}|} \sum_{(x, \mathbf{z}, y^*) \in \mathcal{B}} \left[ \frac{1}{\sum_{j=1}^m T_j} \sum_{j=1}^m \sum_{t=1}^{T_j} \min \left( \rho_t^j A_t^j, \text{clip} \left( \rho_t^j, 1 - \epsilon, 1 + \epsilon \right) A_t^j \right) \right],$$

where  $\epsilon$  is the clipping hyperparameter;  $\rho_t^j = \frac{\pi_\theta(z_t | x, \mathbf{z}^{<j}, \mathbf{z}_{<t}^j)}{\pi_{\theta_{old}}(z_t | x, \mathbf{z}^{<j}, \mathbf{z}_{<t}^j)}$  is the probability ratio of generating token  $z_t$  under the new policy versus the reference policy  $\pi_{\theta_{old}}$ ; and  $A_t^j$  denotes the token-level advantage computed for token  $z_t^j$ .

Building on this algorithmic structure, the per-timestep reward signal  $\hat{r}_t^j$  and its corresponding advantage  $A_t^j$  that is normalized across the mini-batch are defined as follows:

$$\begin{aligned} \hat{r}_t^j &= \mathbb{I}(z_t^j = [\text{EOS}])r(x, \mathbf{z}, y^*) - \beta \text{KL}(j, t), \\ A_t^j &= (\hat{r}_t^j - \mu_{\hat{r}_t^j}) / \sigma_{\hat{r}_t^j}, \end{aligned} \quad (4)$$

where  $\mathbb{I}(s_t^j = [\text{EOS}])$  is an indicator function that assigns the reward  $r(x, \mathbf{z}, y^*)$  only when  $z_t^j$  is the end-of-sequence token. A token-level KL penalty component,  $\text{KL}(j, t)$ , is included to penalize policy deviations from the reference policy  $\pi_{old}$ .

To reconcile the tension between exploratory reasoning trajectories diversity (which may deviate significantly from the reference model distribution) and strict adherence to mathematical formulation/solver syntax requirements in optimization tasks, we propose **Partial KL**. This novel design selectively applies the KL penalty to the mathematical formulation  $\mathbf{z}^{m-1}$  and solver code  $\mathbf{z}^m$  segments. The value for the KL term,  $\text{KL}(j, t)$ , within these segments is computed using the unbiased estimator described in [76]:

$$\text{KL}(j, t) = \begin{cases} \rho_t^j - \log \rho_t^j - 1 & j \in \{m-1, m\}, \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

The rationale for the **Partial KL** design that utilizes selective KL regularization is twofold:

1. **Exploration in reasoning:** For reasoning steps  $(\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^{m-2})$ , the KL penalty is omitted. This fosters exploration, enabling the policy to better understand the problem background and identify diverse reasoning paths and implicit constraints [77, 78].
2. **Stability in modeling and code generation:** For the critical mathematical formulation  $\mathbf{z}^{m-1}$  and solver code  $\mathbf{z}^m$  segments, the KL penalty ensures the generated output remains well-structured and adheres to expected formats, preventing policy collapse while still allowing gradual improvement guided by the reward.

Our SIRL framework, which incorporates **Partial KL**, enables the policy to achieve effective exploration and understanding of the problem while improving code execution accuracy, and yields high-quality outputs for optimization tasks.

### 3.3 Reward design and training scheme

The success of our SIRL framework hinges on its verifiable reward function, implemented as a staged, rule-based system [31, 33]. Integrated with the optimization solver, this system provides objective verification signals used within a two-stage curriculum [79, 80] to progressively train the model and enhance its optimization modeling capabilities.

Given a question  $x$ , generated trajectories  $\mathbf{z}$ , ground-truth answer  $y^*$ , the two-stage reward function  $r(x, \mathbf{z}, y^*)$  is defined as follows:

$$r(x, \mathbf{z}, y^*) = \begin{cases} R_{\text{format}}(\mathbf{z}) + R_{\text{exec}}(\mathbf{z}) + R_{\text{accur}}(x, \mathbf{z}, y^*) & \text{Stage-1,} \\ R_{\text{format}}(\mathbf{z}) + R_{\text{exec}}(\mathbf{z}) + R_{\text{accur}}(x, \mathbf{z}, y^*) + R_{\text{bonus}}(x, \mathbf{z}, y^*) & \text{Stage-2.} \end{cases} \quad (6)$$

In stage-1, we focus on building the model’s fundamental capabilities in formulating and solving standard optimization problems. The reward function comprises three key components: format, execution, and accuracy. Emphasis is placed on the execution component via additional incentives for correctly executed code. This ensures that the generated models are both mathematically sound and executable. Building on fundamental capabilities, Stage 2 aims to tackle more complex problems through the bonus reward  $R_{\text{bonus}}$ , which is based on the generated mathematical model associated with the LP file and designed to incentivize advanced modeling techniques (e.g., Big-M, nonlinear formulations) which are crucial for complex and challenging problems. This bonus is granted only when two conditions are met: (1) the generated solution is correct, and (2) it incorporates advanced modeling strategies. The complete reward function formulation is detailed in the Appendix.

## 4 Experiments

### 4.1 Instance-enhanced self-consistency

We evaluated different self-consistency approaches on the Qwen2.5 models [19] (7B-Instruct and 32B-Instruct) to assess the effect of leveraging instance-level information. The value-based self-consistency method (val\_sc) is a direct adaptation of the standard self-consistency approach where

Table 1: Performance of value-based and instance-enhanced self-consistency on Qwen2.5 Models.

Metric	NL4OPT		MAMOEasy		MAMOCComplex		IndustryOR		OptMATH		Average	
	7B	32B	7B	32B	7B	32B	7B	32B	7B	32B	7B	32B
pass@1	65.7%	67.8%	81.9%	83.7%	17.1%	26.5%	19.0%	23.0%	4.1%	16.6%	37.6%	43.5%
val_sc@5	69.8%	70.2%	85.3%	86.0%	25.1%	35.1%	25.0%	32.0%	5.7%	22.3%	42.2%	49.1%
inst_sc@5	70.6%	71.0%	85.4%	86.3%	29.9%	38.4%	26.0%	33.0%	13.0%	27.5%	45.0%	51.2%
<b>Diff (inst-val)@5</b>	<b>1.1%</b>	<b>1.1%</b>	<b>0.1%</b>	<b>0.3%</b>	<b>19.1%</b>	<b>9.4%</b>	<b>4.0%</b>	<b>3.1%</b>	<b>128.1%</b>	<b>23.3%</b>	<b>30.5%</b>	<b>7.5%</b>
val_sc@10	68.6%	73.5%	85.6%	85.9%	29.9%	38.4%	28.0%	34.0%	9.8%	27.5%	44.4%	51.9%
inst_sc@10	69.0%	72.2%	85.7%	86.0%	32.2%	39.3%	30.0%	36.0%	16.6%	34.2%	46.3%	53.5%
<b>Diff (inst-val)@10</b>	<b>0.6%</b>	<b>-1.8%</b>	<b>0.1%</b>	<b>0.1%</b>	<b>7.7%</b>	<b>2.3%</b>	<b>7.1%</b>	<b>5.9%</b>	<b>69.4%</b>	<b>24.4%</b>	<b>17.0%</b>	<b>6.2%</b>

the final score of different roles depends only on the final objective function value. The instance-enhanced self-consistency method (inst\_sc) also includes structural information within the generated optimization models, augmenting the consensus mechanism. The consensus function of the final objective function value, optimization direction, the count of binary variables and general integer variables are given the same weight in Equation 3.

Table 1 indicates that self-consistency through majority voting outperforms the baseline single-pass generation (pass@1). Both val\_sc and inst\_sc methods demonstrate consistently higher accuracy than the pass@1 baseline. Furthermore, a comparative analysis between the two self-consistency variants suggests that incorporating instance-level information (optimization direction, variable counts) into the voting mechanism provides a more robust measure of consensus, leading to improved selection of correct solutions compared to relying solely on the final objective value.

## 4.2 Main results

Table 2: Performance comparison of models on benchmarks.

Types	Models	Acc (pass@1)					Macro AVG
		NL4OPT	MAMO Easy	MAMO Complex	IndustryOR	OptMATH	
Baseline	GPT-3.5-turbo	78.0%*	79.3%*	33.2%*	21.0%*	15.0%*	45.3%*
	GPT-4	89.0%*	87.3%*	49.3%*	33.0%*	16.6%*	55.0%*
	Deepseek-V3	95.9%*	88.3%*	51.1%*	<b>37.0%*</b>	<b>32.6%*</b>	61.0%*
Agent-based	Chain-of-Experts	64.2%*	-	-	-	-	-
	OptiMUS	78.8%*	77.2%*	43.6%*	31.0%*	20.2%*	49.4%*
Offline-learning	ORLM-LLaMA-3-8B	85.7%*	82.3%*	37.4%*	24.0%*	2.6%*	46.4%
	LLMOpt-Qwen2.5-14B	80.3%*	89.5%*	44.1%*	29.0%*	12.5%*	51.1%
	OptMATH-Qwen2.5-7B	94.7%*	86.5%*	51.2%*	20.0%*	24.4%*	55.4%
Online-RL	SIRL-Qwen2.5-7B	<b>96.3%</b>	<b>90.0%</b>	<b>62.1%</b>	33.0%	29.0%	<b>62.1%</b>

Values marked with \* are from original or reproduced papers with the criterion: relative error  $< 10^{-6}$ .

We evaluated the performance of the proposed SIRL framework on four benchmarks: NL4OPT [81], MAMO [82], IndustryOR [22] and OptMATH [23]. Performance is assessed based on the pass@1 accuracy(acc). Following the rigorous evaluation protocol proposed by OptMATH [23], a solution is considered valid if the relative error is less than  $1e-6$ .

Table 2 presents the main results. The SIRL-Qwen2.5-7B model is trained by the proposed SIRL framework starting from Qwen2.5-7B-Instruct [19] without prior SFT. Compared to existing baselines (foundation LLMs [10, 13], agent-based approaches [15, 16], and offline-learning methods [22, 24, 23]), our model yields consistent and significant performance gains over all (7B,14B)-level models trained by other offline learning methods, and achieves comparable results to the 671B parameter foundation model Deepseek-V3 [13] across all benchmarks. This outcome highlights the efficiency of our proposed SIRL mechanism in enhancing LLMs’ abilities for optimization formulation and code solving, demonstrating its ability of tackling complex optimization modeling challenges. A detailed overview of datasets, evaluation criteria and experimental setup, including prompt structures, training hyperparameter, and decoding strategies can be found in the Appendix.



### 4.3 Ablation Study

In this section, we present two ablation studies to examine the impact of the surrogate function design based on the `Partial KL` strategy and the proposed two-stage reward mechanism.

Table 3: Ablation study on different surrogate function designs.

Type	MAMO Complex		IndustryOR		OptMATH	
	Acc(pass@1)	ER	Acc(pass@1)	ER	Acc(pass@1)	ER
Partial KL	62.1%	97.2%	33.0%	96.0%	29.0%	88.1%
Full KL	58.3%(↓3.8%)	97.6%(↑0.4%)	30%(↓3.0%)	95.0%(↓1.0%)	26.9%(↓2.1%)	89.1%(↑1.0%)
Without KL	57.8%(↓4.3%)	94.8%(↓2.4%)	29%(↓4.0%)	87.0%(↓9.0%)	27.5%(↓1.5%)	72.5%(↓15.6%)

**Ablation study on different surrogate function designs.** We evaluated three distinct surrogate function designs: (i) `Full KL`: the standard approach applying full KL-divergence regularization against the reference policy; (ii) `Without KL`: an approach omitting KL-divergence regularization, which is popular in RLVR training for mathematical problems [77, 78] such as AIME [51]; (iii) `Partial KL`: our novel design that applies the KL penalty selectively to the mathematical formulation and code segments.

Table 3 reports both the pass@1 accuracy and execution rate (ER), which measures the percentage of generated solutions that successfully compile and return a valid result, across three more challenging datasets. The results show that the proposed `Partial KL` approach achieves the best performance across all benchmarks. In contrast, the `Without KL` design exhibits a dramatically lower execution rate than the other two strategies. This lower rate stems from removing KL divergence: while promoting more diverse exploration, it can lead to introducing irrelevant constraints from the problem background, increasing invalid generations. `Partial KL` resolves this issue by applying KL selectively, improving the execution rate while preserving reasoning diversity. The full comparison on all benchmarks and a detailed qualitative analysis, including a case study, are presented in the Appendix.

Table 4: Performance results of the ablation study on reward design.

Reward Type	Acc (pass@1)				
	NL4OPT	MAMO Easy	MAMO Complex	IndustryOR	OptMATH
Two-stage rewards	96.3%	90.0%	62.1%	33.0%	29.0%
Stage-1 reward only	96.7%(↑0.4%)	87.7%(↓2.3%)	55.0%(↓7.1%)	27.0%(↓6.0%)	26.9%(↓2.1%)
Stage-2 reward only	92.2%(↓4.1%)	88.3%(↓1.7%)	58.3%(↓3.8%)	28.0%(↓5.0%)	31.1%(↑2.1%)

**Ablation study on reward design.** We compared the performance of the proposed two-stage reward mechanism (D) against models trained using only the stage-1 reward and using only the stage-2 reward. As shown in Table 4, using only the stage-1 reward yielded comparatively strong results on simple tasks such as NL4OPT. This indicates that this reward enables the model to learn stable foundational skills for optimization tasks. Meanwhile, employing only the stage-2 reward, which includes a bonus component incentivizing the model to learn advanced strategies, achieves the best performance on the most challenging OptMATH dataset. However, this led to diminished performance on simpler tasks such as NL4OPT. Overall, the integrated two-stage reward mechanism successfully balanced the objectives of the individual stages, resolving the trade-offs observed with single-stage rewards, thereby achieving superior performance in most benchmark tasks, with the exception of the OptMATH dataset.

## 5 Conclusion

In this paper, we present SIRL, a novel RLVR framework addressing the challenge of generating authentic mathematical optimization models with LLMs. The core contributions of this work are its unique surrogate function design, `Partial KL`, which selectively applies KL divergence to

mathematical formulation and code segments, and its two-stage reward system, which leverages optimization solvers for automated verification. The comprehensive signals derived from this verification were valuable for both RL training and enhancing our data synthesis. Extensive experiments showed that the SIRL-trained model achieved superior performance in generating accurate, well-formed optimization models compared to existing methods. More broadly, the proposed techniques are applicable to tasks requiring LLMs to balance exploring diverse reasoning with ensuring solution validity, particularly in tool-augmented tasks.

## References

- [1] Ajay Singh. An overview of the optimization modelling applications. *Journal of Hydrology*, 466:167–182, 2012.
- [2] Andreas Antoniou Wu-Sheng Lu. Practical optimization algorithms and engineering applications, 2007.
- [3] Mazura Mokhtar, A Shuib, and D Mohamad. Mathematical programming models for portfolio optimization problem: A review. *International journal of social, management, economics and business engineering*, 8(2):443–450, 2014.
- [4] Gerard Cornuejols and Reha Tütüncü. Optimization methods in finance. *Carnegie Mellon University, Pittsburgh*, 2005.
- [5] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [6] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2024.
- [7] Dongdong Ge, Qi Huangfu, Zizhuo Wang, Jian Wu, and Yinyu Ye. Cardinal optimizer (copt) user guide. *arXiv preprint arXiv:2208.14314*, 2022.
- [8] IBM Ilog. User’s manual for cplex. <http://www.ilog.com/>, 2010.
- [9] David G Luenberger, Yinyu Ye, et al. *Linear and nonlinear programming*, volume 2. Springer, 1984.
- [10] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [11] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [12] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [13] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [14] Rindranirina Ramamonjison, Haley Li, Timothy T Yu, Shiqi He, Vishnu Rengan, Amin Banitalebi-Dehkordi, Zirui Zhou, and Yong Zhang. Augmenting operations research with auto-formulation of optimization models from problem descriptions. *arXiv preprint arXiv:2209.15565*, 2022.
- [15] Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin Fu, Tao Zhong, Jia Zeng, Mingli Song, et al. Chain-of-experts: When llms meet complex operations research problems. In *The Twelfth International Conference on Learning Representations*, 2023.

- [16] Ali AhmadiTeshnizi, Wenzhi Gao, and Madeleine Udell. Optimus: Scalable optimization modeling with (MI)LP solvers and large language models. *arXiv preprint arXiv:2402.10172*, 2024.
- [17] Nicolás Astorga, Tennison Liu, Yuanzhang Xiao, and Mihaela van der Schaar. Autoformulation of mathematical optimization models using llms. *arXiv preprint arXiv:2411.01679*, 2024.
- [18] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [19] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [20] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [21] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- [22] Chenyu Huang, Zhengyang Tang, Shixi Hu, Ruqing Jiang, Xin Zheng, Dongdong Ge, Benyou Wang, and Zizhuo Wang. Orlm: A customizable framework in training large models for automated optimization modeling. *Operations Research*, 2025.
- [23] Hongliang Lu, Zhonglin Xie, Yaoyu Wu, Can Ren, Yuxuan Chen, and Zaiwen Wen. Optmath: A scalable bidirectional data synthesis framework for optimization modeling. *arXiv preprint arXiv:2502.11102*, 2025.
- [24] 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. *arXiv preprint arXiv:2410.13213*, 2024.
- [25] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [26] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [27] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- [28] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [29] Sara Vera Marjanović, Arkil Patel, Vaibhav Adlakha, Milad Aghajohari, Parishad BehnamGhader, Mehar Bhatia, Aditi Khandelwal, Austin Kraft, Benno Krojer, Xing Han Lü, et al. Deepseek-r1 thoughtology: Let’s< think> about llm reasoning. *arXiv preprint arXiv:2504.07128*, 2025.
- [30] OpenAI. Learning to reason with llms. <https://openai.com/index/learning-to-reason-with-llms/>. Accessed: 2024-09-12.
- [31] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [32] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.

- [33] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- [34] Chengpeng Li, Mingfeng Xue, Zhenru Zhang, Jiaxi Yang, Beichen Zhang, Xiang Wang, Bowen Yu, Binyuan Hui, Junyang Lin, and Dayiheng Liu. Start: Self-taught reasoner with tools. *arXiv preprint arXiv:2503.04625*, 2025.
- [35] Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjuan Zhong. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*, 2025.
- [36] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [37] Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan J Prenger, and Animashree Anandkumar. Leandajo: Theorem proving with retrieval-augmented language models. *Advances in Neural Information Processing Systems*, 36:21573–21612, 2023.
- [38] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khoshdel, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- [39] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- [40] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- [41] Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023.
- [42] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- [43] Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. rstar-math: Small llms can master math reasoning with self-evolved deep thinking. *arXiv preprint arXiv:2501.04519*, 2025.
- [44] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.
- [45] 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. *arXiv e-prints*, pages arXiv–2407, 2024.
- [46] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- [47] Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and Jingren Zhou. Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825*, 2023.
- [48] Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.

- [49] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [50] 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. *arXiv preprint arXiv:2103.03874*, 2021.
- [51] AI-MO Team. AI-MO Validation and AIME Test Set Dataset. Hugging Face Datasets, 2024.
- [52] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [53] Shanghaoran Quan, Jiayi Yang, Bowen Yu, Bo Zheng, Dayiheng Liu, An Yang, Xuancheng Ren, Bofei Gao, Yibo Miao, Yunlong Feng, et al. Codeelo: Benchmarking competition-level code generation of llms with human-comparable elo ratings. *arXiv preprint arXiv:2501.01257*, 2025.
- [54] Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471, 2022.
- [55] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [56] Lilian Weng. Reward hacking in reinforcement learning. lilianweng. github. io, 2024.
- [57] Huajian Xin, ZZ Ren, Junxiao Song, Zhihong Shao, Wanbiao Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, et al. Deepseek-prover-v1. 5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. *arXiv preprint arXiv:2408.08152*, 2024.
- [58] Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.
- [59] Leonardo de Moura and Sebastian Ullrich. The lean 4 theorem prover and programming language. In *Automated Deduction—CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July 12–15, 2021, Proceedings 28*, pages 625–635. Springer, 2021.
- [60] Huaiyuan Ying, Zijian Wu, Yihan Geng, Jiayu Wang, Dahua Lin, and Kai Chen. Lean workbook: A large-scale lean problem set formalized from natural language math problems. *arXiv preprint arXiv:2406.03847*, 2024.
- [61] Bohan Li, Jiaxuan Wang, Cong Fang, Fan Yang, Jing Xiong, Yingkan Xu, Shiyang Li, Xiang Li, Songfang Huang, and Jingren Zhou. DeepSeek-Prover-V2: Advancing Formal Mathematical Reasoning via Reinforcement Learning for Subgoal Decomposition. *arXiv preprint arXiv:2504.21801*, 2025.
- [62] Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.
- [63] Yuri Chervonyi, Trieu H Trinh, Miroslav Olšák, Xiaomeng Yang, Hoang Nguyen, Marcelo Menegali, Junehyuk Jung, Vikas Verma, Quoc V Le, and Thang Luong. Gold-medalist performance in solving olympiad geometry with alphageometry2. *arXiv preprint arXiv:2502.03544*, 2025.
- [64] Chengpeng Li, Guanting Dong, Mingfeng Xue, Ru Peng, Xiang Wang, and Dayiheng Liu. Dotamath: Decomposition of thought with code assistance and self-correction for mathematical reasoning. *arXiv preprint arXiv:2407.04078*, 2024.
- [65] Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. Tora: A tool-integrated reasoning agent for mathematical problem solving. *arXiv preprint arXiv:2309.17452*, 2023.

- [66] Minpeng Liao, Wei Luo, Chengxi Li, Jing Wu, and Kai Fan. Mario: Math reasoning with code interpreter output—a reproducible pipeline. *arXiv preprint arXiv:2401.08190*, 2024.
- [67] MOSEK ApS. Mosek optimization suite. 2019.
- [68] 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. *arXiv preprint arXiv:2407.09887*, 2024.
- [69] Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*, 2024.
- [70] Vedant Shah, Dingli Yu, Kaifeng Lyu, Simon Park, Jiatong Yu, Yinghui He, Nan Rosemary Ke, Michael Mozer, Yoshua Bengio, Sanjeev Arora, et al. Ai-assisted generation of difficult math questions. *arXiv preprint arXiv:2407.21009*, 2024.
- [71] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- [72] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2023.
- [73] Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*, 2025.
- [74] Jian Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.
- [75] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [76] John Schulman. Approximating kl divergence. *John Schulman’s Homepage*, 2020.
- [77] Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- [78] Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang Zhang, Jiacheng Xu, Wei Shen, Siyuan Li, Liang Zeng, Tianwen Wei, Cheng Cheng, Bo An, Yang Liu, and Yahui Zhou. Skywork open reasoner series. <https://capricious-hydrogen-41c.notion.site/Skywork-Open-Reasoner-Series-1d0bc9ae823a80459b46c149e4f51680>, 2025. Notion Blog.
- [79] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- [80] Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *international conference on machine learning*, pages 1311–1320. Pmlr, 2017.
- [81] Rindranirina Ramamonjison, Timothy Yu, Raymond Li, Haley Li, Giuseppe Carenini, Bissan Ghaddar, Shiqi He, Mahdi Mostajabdaveh, Amin Banitalebi-Dehkordi, Zirui Zhou, et al. NL4Opt competition: Formulating optimization problems based on their natural language descriptions. In *NeurIPS 2022 Competition Track*, pages 189–203. PMLR, 2023.
- [82] Xuhan Huang, Qingning Shen, Yan Hu, Anningzhe Gao, and Benyou Wang. MAMO: A mathematical modeling benchmark with solvers. *arXiv preprint*, 2024.
- [83] Mokhtar S Bazaraa, John J Jarvis, and Hanif D Sherali. *Linear programming and network flows*. John Wiley & Sons, 2011.

- [84] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:2409.19256*, 2024.
- [85] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019.

## A Benchmark dataset

We evaluated the performance of our trained model on multiple datasets that include NL4OPT [81], MAMO [82], IndustryOR [22], OptMATH [23]. The number of problems within each dataset is detailed in Table 5. As noted in [23, 24], there are minor errors within these testing datasets. To address this, we rigorously reviewed and corrected the test sets of these benchmarks, updating the questions and corresponding answers to ensure the integrity of our evaluation, with a specific focus on the NL4OPT and IndustryOR dataset.

Table 5: Summary statistics for the optimization problem datasets

Dataset	# of Data
NL4OPT	245
Mamo Easy	652
Mamo Complex	211
IndustryOR	100
OptMATH-Bench	193

- **NL4OPT**: Originating from the NL4OPT competition, this benchmark assesses the conversion of natural language problem statements into executable optimization models. It predominantly features linear programming (LP) tasks drawn from diverse domains, although the complexity of the underlying mathematical structures is relatively consistent. A final set comprising 245 high-quality instances was obtained by filtering low-quality instances [23].
- **MAMO** consists of two subsets: **MAMO EasyLP** and **MAMO ComplexLP**. **MAMO EasyLP** contains 652 instances, primarily high-school level mixed-integer linear programming (MILP) problems designed for fundamental learning; **MAMO ComplexLP** includes 211 more challenging undergraduate-level problems, integrating both LP and MILP concepts suitable for assessing advanced modeling skills.
- **IndustryOR**: Introduced as the first benchmark specifically targeting industrial applications, IndustryOR comprises 100 real-world optimization scenarios sourced from various sectors. It distinguishes itself by incorporating a wider range of problem types, including linear programming (LP), integer programming (IP), mixed-integer programming (MIP), nonlinear programming (NLP), and other specialized formulations, categorized across three difficulty levels.
- **OptMATH**: The benchmark was constructed to address the limitations of existing datasets, particularly their frequent focus on linear problems and the potentially insufficient challenge level. OptMATH features a curated selection of complex mathematical optimization problems designed to effectively differentiate advanced modeling capabilities. Its scope is intentionally broad, encompassing LP, MILP, IP, NLP, Second-Order Cone Programming (SOCP), and other relevant optimization paradigms. Two variants of the dataset exist, with 166 and 193 instances, respectively. We present the results from the larger 193-instance variant in the main results part to ensure a thorough evaluation.

We evaluate performance across these benchmarks using a strict accuracy metric: a generated solution is considered correct if the related difference between its predicted objective value  $y_{pred}$  and the ground truth value  $y_{label}$  satisfies:

$$\frac{|y_{pred} - y_{label}|}{|y_{label}|} < 10^{-6}.$$

Similar to [23], the errors about the type of the variables are ignored. We thoroughly review problems exhibiting such errors. In these cases, the decision variables typically represent time or similar quantities, which can reasonably be treated as integer or continuous variables.



## B Seed data

The seed data consist of 686 real-world industry cases collected from operations research textbooks and a list of 100 industry scenarios [22]. Here is an example:

### Seed data example

**Problem:** \*\*Tourist Restaurant Seating Allocation Problem\*\*

**\*\*Problem Background:\*\***

The tourist restaurant seating allocation problem involves arranging the distribution of 100 seats in the restaurant at two different periods to improve customer satisfaction, optimize service efficiency and maximize restaurant revenue. This is a common management challenge in the catering industry, especially during peak dining periods.

**\*\*Optimization Goals:\*\***

- Maximize customer satisfaction: Improve customer dining experience through reasonable seat allocation.
- Maximize restaurant revenue: Maximize restaurant revenue by optimizing seat allocation.
- Optimize service efficiency: Ensure that the restaurant can operate efficiently at different times and reduce waiting time.

**\*\*Numerical Problem Constraints:\*\***

1. **\*\*Total Seats Limit:\*\***

- At any given period, the total number of seats allocated must not exceed the restaurant's total capacity of 100 seats.

2. **\*\*Fluctuating Demand Across Periods:\*\***

- In Period 1, the seat allocation must be between the minimum demand of 20 seats and the maximum demand of 60 seats.
- In Period 2, the seat allocation must be between the minimum demand of 30 seats and the maximum demand of 70 seats.

**\*\*Objective Function:\*\***

The objective is to maximize customer satisfaction and restaurant revenue, where the combined weight for customer satisfaction and revenue per allocated seat is 1.0 for each.

**\*\*Formulation:\*\***

Let ' $x[i, j]$ ' denote the allocation of seat ' $j$ ' in period ' $i$ ' (where ' $i = 1, 2$ ' for two periods and ' $j = 1, 2, \dots, 100$ ' for total seats).

**\*\*Objective Function:\*\***

Maximize:

$$\sum_{i=1}^2 \sum_{j=1}^{100} (1.0 + 1.0) x[i, j]$$

**\*\*Constraints:\*\***

1. Total seat allocation in each period should not exceed the total available seats:

$$\sum_{j=1}^{100} x[i, j] \leq 100, \quad \forall i \in \{1, 2\}$$

2. Seat allocation per period must satisfy minimum and maximum demand:

$$20 \leq \sum_{j=1}^{100} x[1, j] \leq 60 \quad \forall \quad 30 \leq \sum_{j=1}^{100} x[2, j] \leq 70$$

**\*\*Expected Output:\*\***

The optimal seating allocation for each period (' $x[i, j]$ ' values) that maximizes customer satisfaction and revenue while adhering to all constraints.

**Answer:** 260.0

## C Prompt templates

In this section, we summarize the prompts used in the framework.

### C.1 Ten roles for self-consistency

In the self-consistency mechanism, we generate multiple candidate solutions for each problem by prompting the LLM with ten distinct roles. Each role represents a specific expert persona, designed to elicit varied approaches to optimization modeling by emphasizing different combinations of expertise in operations research, Python development, and Gurobi solver application. The ten roles are defined as:

#### Ten roles

- **Role 1:** "A highly skilled Python engineer and optimization specialist with deep expertise in operations research and the Gurobi solver."
- **Role 2:** "An optimization expert and Python engineer specializing in operations research and the Gurobi solver."
- **Role 3:** "A Python engineer and optimization specialist with a strong background in operations research and the Gurobi solver."
- **Role 4:** "A skilled Python engineer and optimization specialist proficient in operations research and the Gurobi solver."
- **Role 5:** "A results-driven Python engineer and optimization expert with a strong foundation in operations research and the Gurobi solver."
- **Role 6:** "A seasoned operations research scientist and Python developer, leveraging advanced optimization techniques and the Gurobi solver to tackle complex business challenges."
- **Role 7:** "An innovative optimization modeler and Python programmer, specializing in the development and implementation of high-performance solutions using operations research methodologies and the Gurobi optimization suite."
- **Role 8:** "A pragmatic problem-solver with expertise in operations research, proficient in Python and the Gurobi API, focused on translating real-world scenarios into efficient and scalable optimization models."
- **Role 9:** "A meticulous optimization analyst and Python coder, deeply familiar with the theoretical underpinnings of operations research and the practical application of the Gurobi solver for achieving optimal outcomes."
- **Role 10:** "A strategic optimization architect and Python implementation specialist, with a proven track record of designing and deploying robust operations research solutions powered by the Gurobi optimization engine."

### C.2 LLM-as-a-judge for generated problem evaluation

We employ the LLM-as-a-judge methodology [46] to validate the generated problems for practical relevance and semantic consistency. The prompt utilized for this validation process is detailed below.

### Problem evaluation prompt

You are an expert in operations research. You'll receive an operations research problem. You will analyze it and determine whether the problem is a valid one considering the following aspects:

1. Determine if the problem's language and structure are consistent with typical operations research problem formulations.
2. Assess whether the problem scenario has real-world applicability or practical significance.
3. Identify any semantic inconsistencies, contradictions, or ambiguities within the problem statement.

Below is the operations research problem:

{{Question}}

Please provide your step-by-step analysis and your final judgment for each of these points.

## Analysis Process:

[Your detailed step-by-step reasoning for each point above]

## Final Judgment:

[Yes or No]

### C.3 Refine and regenerate for the error case

Upon encountering code execution errors, we leverage the LLM for refinement and code regeneration [71]. The prompt used for this error correction mechanism is as follows.

### Error regenerate prompt

You are an experienced operations research algorithm engineer. You are presented with an operations research problem and a previous attempt to model and code a solution. That attempt resulted in an error.

Problem Description:

{{Question}}

Previous Code Solution Attempt:

{{Previous code}}

After running the provided code from the previous attempt, the following error occurred:

{{Error output after executing the code}}

Your task:

Based on the information above, please perform the following:

1. Analyze Root Cause & Identify Pitfalls

Thoroughly analyze the root cause of the error.

Summarize potential pitfalls or common mistakes related to this type of code error.

2. Provide Corrected Gurobi Code:

Write the complete and corrected Python code using the 'gurobipy' library to accurately solve the problem.

Please structure your response strictly as follows:

## Cause of the Error and Potential Pitfalls:

[Your detailed analysis of the error's cause and a summary of potential pitfalls.]

## Corrected Gurobi Code:

[Your complete and corrected Gurobi Python code.]

### C.4 Refine and regenerate for the infeasible case

If an infeasible solution is obtained after executing the code using the solver, we leverage the LLM to refine the entire result by regenerating both the mathematical model and its corresponding code [71]. The prompt used for this infeasibility resolution mechanism is as follows.

### Infeasible regenerate prompt

You are an experienced operations research algorithm engineer. You are presented with an operations research problem and a previous attempt to model and code a solution. That attempt resulted in an infeasible solution.

Problem Description:

{{Question}}

Previous Model and Code Solution Attempt:

{{Previous model}}

{{Previous code}}

After running the provided code from the previous attempt, the answer could not provide a feasible solution.

Your task:

Based on the information above, please perform the following:

1. Analyze Root Cause & Identify Pitfalls

Thoroughly analyze the root cause of the infeasibility.

Summarize potential pitfalls or common mistakes related to this type of infeasibility.

2. Provide an Improved Mathematical Model: Develop a mathematical model for correctly modeling this OR problem. This should address the flaws in the previous attempt.

3. Provide Corrected Gurobi Code:

Write the complete and corrected Python code associated with the mathematical model using the 'gurobipy' library to accurately solve the problem.

Please structure your response strictly as follows:

## Cause of the Infeasibility and Potential Pitfalls:

[Your detailed analysis of the infeasibility's cause and a summary of potential pitfalls.]

## Corrected Mathematical Model:

[Your improved mathematical model.]

## Corresponding Gurobi Code:

[Your complete and corrected Gurobi Python code associated with the mathematical model.]

## C.5 System prompt and the LLM responses

We follow the approach of DeepSeek-R1-Zero [31] by providing a specific output format within the system prompt. This format guides the LLMs through a sequential process: it first asks the LLM to analyze the problem, then formulates a mathematical model, and finally outputs the corresponding solver code. To ensure clear delineation and parsability of these steps, we prompt the LLM to explicitly separate each component using distinct tags: <think></think>, <model></model>, <python></python>. An example output response guided by the system prompt is provided below. The corresponding problem was sourced from the NL4OPT dataset.

### System Prompt

**SYSTEM:** You are a helpful Assistant with expertise in operations research and the Gurobi solver. When the User provides an OR question, you will analyze it, build a detailed mathematical model, and provide the Gurobi code to solve it.

Your response should follow these steps:

1. <think> Carefully analyze the problem to identify decision variables, objective, and constraints.</think>

2. <model> Develop a complete mathematical model, explicitly defining: \* Sets \* Parameters \* Decision Variables (and their types) \* Objective Function \* Constraints </model>

3. <python> Provide the corresponding Gurobi Python code to implement the model.</python>

The output must be in Markdown format, with each step enclosed in the specified tags.

**USER:** Answer the following mathematical modeling question:

{{Question}}

Please think step by step.

### Problems

**Problem:** An industrial tire company delivers large tires for equipment to remote engineering sites either by cargo planes or ultrawide trucks. Each cargo plane can transport 10 tires per trip and costs \$1000. Each ultrawide truck can transport 6 tires per trip and costs \$700. The company needs to transport at least 200 tires and has available \$22000. Because most remote sites don't have proper airports, the number of plane trips cannot exceed the number of ultrawide truck trips. How many trips of each should be done to minimize the total number of trips?

## Response

<think> To solve this problem, we need to identify the decision variables, the objective function, and the constraints.

1. **Decision Variables:**

- Let  $x$  be the number of plane trips.

- Let  $y$  be the number of truck trips.

2. **Objective Function:**

We need to minimize the total number of trips, which is  $x + y$ .

3. **Constraints:**

- The total number of tires transported must be at least 200:  $10x + 6y \geq 200$ .

- The total cost must not exceed \$22000:  $1000x + 700y \leq 22000$ .

- The number of plane trips cannot exceed the number of truck trips:  $x \leq y$ .

- Non-negativity constraints:  $x \geq 0, y \geq 0$ .

Since these are linear constraints, we can use linear programming to solve this problem.

</think>

<model>

### Sets

- There are no sets in this problem as it is a simple linear programming problem.

### Parameters

- The number of tires each plane can transport: 10.

- The cost per plane trip: \$1000.

- The number of tires each truck can transport: 6.

- The cost per truck trip: \$700.

### Decision Variables

-  $x$ : Number of plane trips.

-  $y$ : Number of truck trips.

### Objective Function

Minimize the total number of trips:  $x + y$ .

### Constraints

1.  $10x + 6y \geq 200$

2.  $1000x + 700y \leq 22000$

3.  $x \leq y$

4.  $x \geq 0, y \geq 0$

</model>

<python>

```
import gurobipy as gp
```

```
from gurobipy import GRB
```

```
# Create a new model
```

```
m = gp.Model("tire_transport")
```

```
# Decision variables
```

```
x = m.addVar(name="x", vtype=GRB.INTEGER) # Number of plane trips
```

```
y = m.addVar(name="y", vtype=GRB.INTEGER) # Number of truck trips
```

```
# Objective function
```

```
m.setObjective(x + y, GRB.MINIMIZE)
```

```
# Constraints
```

```
m.addConstr(10 * x + 6 * y >= 200, "tire_constraint")
```

```
m.addConstr(1000 * x + 700 * y <= 22000, "cost_constraint")
```

```
m.addConstr(x - y <= 0, "trip_constraint")
```

```
# Optimize the model
```

```
m.optimize()
```

```
# Print the solution
```

```
if m.status == GRB.OPTIMAL:
```

```
    print(f"Optimal solution found:")
```

```
    print(f"Plane trips: {x.x}")
```

```
    print(f"Truck trips: {y.x}")
```

```
    print(f"Total trips: {x.x + y.x}")
```

```
    print(f"Total cost: ${m.objVal}")
```

```
else:
```

```
    print("No optimal solution found.")
```

```
</python>
```

## D Reward function design

The overall two-stage reward function  $r(x, \mathbf{z}, y^*)$  is defined as follows:

$$r(x, \mathbf{z}, y^*) = \begin{cases} R_{\text{format}}(\mathbf{z}) + R_{\text{exec}}(\mathbf{z}) + R_{\text{accur}}(x, \mathbf{z}, y^*) & \text{Stage-1,} \\ R_{\text{format}}(\mathbf{z}) + R_{\text{exec}}(\mathbf{z}) + R_{\text{accur}}(x, \mathbf{z}, y^*) + R_{\text{bonus}}(x, \mathbf{z}, y^*) & \text{Stage-2.} \end{cases}$$

The reward function comprises the following components:

**Format reward**  $R_{\text{format}}(\mathbf{z})$ : the format reward guides the LLM policy to produce response  $\mathbf{z}$  with a specific, parsable structure trajectory defined by the system prompt. This structure segregates the solution trajectory via tags like `<think></think>` for reasoning steps, `<model></model>` for the optimization model, and `<python></python>` for executable code.  $R_{\text{format}}(\mathbf{z})$  is a binary reward (1 or 0) awarded only if  $\mathbf{z}$  strictly includes all required tags in their correct order.

Let  $\mathcal{T} = \{\text{<think>...</think>, <model>...</model>, <python>...</python>}\}$  be the set of required tag pairs. The reward is:

$$R_{\text{format}}(\mathbf{z}) = \begin{cases} 0.5 & \text{if } \mathbf{z} \text{ contains all tags in } \mathcal{T} \text{ according to system prompt} \\ 0 & \text{otherwise.} \end{cases}$$

This reward is also foundational for enabling the extraction and evaluation of the generated model and code.

**Execution reward**  $R_{\text{exec}}(\mathbf{z})$ : assigns a reward of 1 if the optimization code within response  $\mathbf{z}$  is executable, and 0 otherwise.

$$R_{\text{exec}}(\mathbf{z}) = \begin{cases} 1 & \text{if the code is executable,} \\ 0 & \text{otherwise.} \end{cases}$$

**Accuracy reward**  $R_{\text{accur}}(x, \mathbf{z}, y^*)$ : the accuracy reward evaluates the correctness of the final answer  $y = g(x, \mathbf{z})$  obtained by executing the code in  $\mathbf{z}$ . The answer is considered correct if matches the ground truth  $y^*$  within a tolerance  $|y - y^*| \leq 0.01$ . In the first stage, the reward is defined as

$$R_{\text{accur}}(x, \mathbf{z}, y^*) = \begin{cases} 2 & \text{if the answer is right,} \\ 0 & \text{otherwise.} \end{cases}$$

**Bonus accuracy reward**  $R_{\text{bonus}}(x, \mathbf{z}, y^*)$ : real-world optimization problems frequently involve non-linear relationships or discrete variables, to encourage our model to tackle more complex optimization problems requiring techniques beyond standard Linear Programming (LP), we introduce a bonus reward. By analyzing the LP file generated by the solver code, we can verify whether these advanced techniques (Big-M methods [83], binary variables, or nonlinear formulations) is used. The binary bonus  $R_{\text{bonus}}(\mathbf{z})$  is awarded for output  $\mathbf{z}$  if and only if, both the correct answer derived from  $\mathbf{z}$  is correct and the generated model utilizes advanced modeling techniques detectable through instance analysis.

$$R_{\text{bonus}}(\mathbf{z}) = \begin{cases} 1 & \text{if advanced modeling techniques are used,} \\ 0 & \text{otherwise.} \end{cases}$$

## E Details of experiments

**Training setup.** All experiments in the paper were conducted on a single compute node with eight 80GB NVIDIA H100 GPUs.

Starting from the synthetic dataset, we applied a filtering strategy guided by the principle ‘‘Less is More’’ [72, 73]. Specifically, we excluded (question, answer) pairs if the baseline Qwen-32B-Instruct model [19] achieved a 80% success rate (8/10 attempts across different prompting roles) in generating executable code matching the ground-truth optimal value, as such samples were deemed too trivial. This process yielded approximately 70,000 samples. From this set, we then randomly sampled 10,000 instances to form our training data.

We used Qwen2.5-7B-Instruct [19] as the base model and adapted the Verl framework [84] for reinforcement learning training, modifying its implementation to incorporate our novel surrogate function design with the Partial KL strategy and two-stage reward mechanism.

Table 6: Training Parameters

Type	Parameter	Value
Algorithm	Advantage Estimator	reinforce_plus_plus
Data	Batch size	128
	Learning rate	1e-6
	Max prompt length	2048
	Max response length	8192
	Truncation	left
Actor/Rollout	KL loss type	low_var_kl
	KL loss coefficient	0.005
	Rollout number	8
	PPO mini batch size	8
	PPO micro batch Size per GPU	4
	Clip ratio low	0.20
	Clip ratio high	0.28

The key hyperparameters for SIRL training are detailed in Table 6:

**Decoding strategy.** We employ the top-P (nucleus) decoding strategy [85] for the training and inference phases. The exact sampling hyperparameters used to generate our main results are specified in Table 7:

Table 7: Sampling parameters used for text generation.

Parameter	Value
n	1
Temperature	0.5
Top p	0.9
Max tokens	8192
Repetition penalty	1.02

## F In-depth analysis of the Partial KL strategy

**Full ablation study on different surrogate function designs.** Here, we present a detailed analysis of the Partial KL divergence. Table 8 shows the results of an ablation study on all benchmarks, which, due to page limitations, was not included in the main paper. The results are consistent with those reported in the original manuscript, and the surrogate function design employing the Without KL strategy demonstrates a significantly reduced execution rate compared to the other two designs.

Table 8: Ablation study on Partial KL

Type	NL4OPT		MAMOEasy		MAMOCComplex		IndustryOR		OptMATH	
	Acc(pass@1)	ER	Acc(pass@1)	ER	Acc(pass@1)	ER	Acc(pass@1)	ER	Acc(pass@1)	ER
Partial KL	96.3%	100.0%	90.0%	100.0%	62.1%	97.2%	33.0%	96.0%	29.0%	88.1%
Full KL	95.1% <b>(↓1.2%)</b>	99.2%	89.9% <b>(↓0.1%)</b>	99.7%	58.3% <b>(↓3.8%)</b>	97.6%	30% <b>(↓3.0%)</b>	95.0%	26.9% <b>(↓2.1%)</b>	89.1%
Without KL	92.7% <b>(↓3.6%)</b>	97.1%	88.7% <b>(↓1.3%)</b>	100.0%	57.8% <b>(↓4.3%)</b>	94.8%	29% <b>(↓4.0%)</b>	87.0%	27.5% <b>(↓1.5%)</b>	72.5%

**Case study on different surrogate function designs.** In this part, two case studies are presented to demonstrate how different surrogate function configurations affect LLM-generated responses. The first case study demonstrates the limitations of Without KL strategies, showcasing execution errors in an elementary mathematical modeling question; The second examines how the Partial KL strategy, by selectively removing the KL divergence term, enhances reasoning capabilities to better understand questions, especially those involving logic constraints.

**Case Study 1:** The first example is chosen from the IndustryOR dataset. In this case, we compare the generated code snippets to present the execution error arising from the model trained with the Without KL strategy. We observe that omitting the KL divergence in the code component leads to



execution errors. Specifically, as marked in **red**, the response from the Without KL model attempts to enforce that decision variables are greater than zero, but the implemented code is non-functional since ' $>$ ' not supported between instances of 'Var' and 'int'.

**Case study 2:** The second case study involves a problem from the IndustryOR dataset, characterized by more complex logical constraints. In this case, an incorrect reasoning step (highlighted in **red**) leads to a flawed problem formulation. This category of problems demands a deep understanding of the underlying structure and inherent logical dependencies. We observe that omitting the KL divergence penalty encourages the model to enhance its reasoning capabilities. This fosters more effective exploration of the solution space.

#### Case Study 1: Question instance from the IndustryOR dataset

A product can be processed on any one of the four devices: A, B, C, or D. The preparation completion costs when each device is enabled, the unit production cost for the product, and the maximum processing capacity of each device are shown in Table 5-7. If 2000 units of the product need to be produced, how can the total cost be minimized? Try to establish a mathematical model.

Table 5-7 Device, Prep Completion Cost (Yuan), Unit Production Cost (Yuan/Unit), Maximum Processing Capacity (Units)

Device	Prep Completion Cost (Yuan)	Unit Production Cost (Yuan/Unit)	Maximum Processing Capacity (Units)
A	1000	20	900
B	920	24	1000
C	800	16	1200
D	700	28	1600

#### Case Study 1: Generated response under the Full KL strategy

```
<python>
import gurobipy as gp
from gurobipy import GRB
# Create a new model
m = gp.Model("product_production")
# Decision variables
x_A = m.addVar(vtype=GRB.CONTINUOUS, name="x_A")
x_B = m.addVar(vtype=GRB.CONTINUOUS, name="x_B")
x_C = m.addVar(vtype=GRB.CONTINUOUS, name="x_C")
x_D = m.addVar(vtype=GRB.CONTINUOUS, name="x_D")
y_A = m.addVar(vtype=GRB.BINARY, name="y_A")
y_B = m.addVar(vtype=GRB.BINARY, name="y_B")
y_C = m.addVar(vtype=GRB.BINARY, name="y_C")
y_D = m.addVar(vtype=GRB.BINARY, name="y_D")
# Objective function
m.setObjective(1000 * y_A + 920 * y_B + 800 * y_C + 700 * y_D + 20 * x_A + 24 * x_B + 16 * x_C +
28 * x_D, GRB.MINIMIZE)
# Constraints
m.addConstr(x_A + x_B + x_C + x_D == 2000, "total_production")
m.addConstr(x_A <= 900 * y_A, "capacity_A")
m.addConstr(x_B <= 1000 * y_B, "capacity_B")
m.addConstr(x_C <= 1200 * y_C, "capacity_C")
m.addConstr(x_D <= 1600 * y_D, "capacity_D")
# Optimize the model
m.optimize()
# Print the solution
if m.status == GRB.OPTIMAL:
    print(f"Total cost: {m.objVal}")
    print(f"Units processed on A: {x_A.x}")
    print(f"Units processed on B: {x_B.x}")
    print(f"Units processed on C: {x_C.x}")
    print(f"Units processed on D: {x_D.x}")
    print(f"Device A enabled: {y_A.x}")
    print(f"Device B enabled: {y_B.x}")
    print(f"Device C enabled: {y_C.x}")
    print(f"Device D enabled: {y_D.x}")
else:
    print("No optimal solution found.")
</python>
```

### Case Study 1: Generated response under the Without KL strategy

```
<python>
from gurobipy import *
# Define the model
model = Model("Product_Optimization")
# Parameters
prep_completion_cost = {"A": 1000, "B": 920, "C": 800, "D": 700}
unit_production_cost = {"A": 20, "B": 24, "C": 16, "D": 28}
max_processing_capacity = {"A": 900, "B": 1000, "C": 1200, "D": 1600}
# Decision variables
x = {}
for i in prep_completion_cost.keys():
    x[i] = model.addVar(name=i, lb=0, ub=max_processing_capacity[i], vtype=GRB.CONTINUOUS)
# Objective function
model.setObjective(quicksum(prep_completion_cost[i] * (x[i] > 0) + unit_production_cost[i] * x[i] for i
in prep_completion_cost.keys()), GRB.MINIMIZE)
# Constraints
model.addConstr(quicksum(x[i] for i in prep_completion_cost.keys()) == 2000, "Total_Production")
# Optimize the model
model.optimize()
# Print the solution
if model.status == GRB.Status.OPTIMAL:
    print("Optimal Solution:")
    for i in prep_completion_cost.keys():
        print(f"Device {i}: {x[i].x}")
    print(f"Total Cost: {model.objVal}")
else:
    print("No optimal solution found.")
</python>
```

### Case Study 1: Generated response under the Partial KL strategy

```
<python>
from gurobipy import *
# Define the parameters
prep_completion_cost = {'A': 1000, 'B': 920, 'C': 800, 'D': 700}
unit_production_cost = {'A': 20, 'B': 24, 'C': 16, 'D': 28}
max_processing_capacity = {'A': 900, 'B': 1000, 'C': 1200, 'D': 1600}
# Create a model
model = Model("ProductProduction")
# Define decision variables
x = {}
y = {}
for device in max_processing_capacity.keys():
    x[device] = model.addVar(vtype=GRB.BINARY, name=f"x_{device}")
    y[device] = model.addVar(name=f"y_{device}")
# Set objective
model.setObjective(quicksum(prep_completion_cost[device] * x[device] +
    unit_production_cost[device] * y[device] for device in max_processing_capacity.keys()),
    GRB.MINIMIZE)
# Add constraints
model.addConstr(quicksum(y[device] for device in max_processing_capacity.keys()) == 2000)
for device in max_processing_capacity.keys():
    model.addConstr(y[device] <= max_processing_capacity[device] * x[device])
# Optimize the model
model.optimize()
# Print the optimal solution
if model.status == GRB.OPTIMAL:
    print("Optimal solution found:")
    for device in max_processing_capacity.keys():
        print(f"Device {device}: x = {x[device].x}, y = {y[device].x}")
    print(f"Minimum cost: {model.objVal}")
else:
    print("No optimal solution found.")
</python>
```

### Case Study 2: Question instance from the IndustryOR dataset

A manufacturing company needs to transport 1800 units of product from the warehouse to three different sales points. The company has four transportation options to choose from: truck, van, motorcycle, and electric vehicle. Since the van and electric vehicle both consume a lot of energy, the company wants to choose only one of these two options. Each trip with a truck generates 100 units of pollution, a van generates 50 units of pollution, a motorcycle generates 10 units of pollution, and an electric vehicle generates 0 units of pollution. The total pollution generated from all trips cannot exceed 2000 units. At least 10 trips must use a truck. Trucks, vans, motorcycles, and electric vehicles can transport 100 units, 80 units, 40 units, and 60 units of product per trip, respectively. The company needs to ensure that the total amount of transported product is at least 1800 units.

## Case Study 2: Generated response under the Full KL strategy

<think>

To solve this problem, we need to identify the decision variables, the objective function, and the constraints.

1. **Decision Variables**:

- Let  $(T)$  be the number of trips using trucks.
- Let  $(V)$  be the number of trips using vans.
- Let  $(M)$  be the number of trips using motorcycles.
- Let  $(E)$  be the number of trips using electric vehicles.

2. **Objective Function**:

- The objective is to minimize the total pollution, which can be expressed as the equation  $(100T + 50V + 10M + 0E)$ .

3. **Constraints**:

- The total product transported must be at least 1800 units:  $(100T + 80V + 40M + 60E \geq 1800)$ .
- The total pollution generated cannot exceed 2000 units:  $(100T + 50V + 10M \leq 2000)$ .
- At least 10 trips must use trucks:  $(T \geq 10)$ .
- Non-negativity constraints:  $(T, V, M, E \geq 0)$ .

4. **Binary Variable for Van or Electric Vehicle**:

- Since the company wants to choose only one of the van or electric vehicle, we can use a binary variable  $(B)$  where  $(B = 1)$  if a van is used and  $(B = 0)$  if an electric vehicle is used. This can be handled by the constraints  $(V \leq B)$  and  $(E \leq 1 - B)$ .

Let's now develop the mathematical model and provide the Gurobi Python code.

</think>

<model>

### Sets

- No sets are needed for this problem.

### Parameters

- Pollution per trip for trucks: 100
- Pollution per trip for vans: 50
- Pollution per trip for motorcycles: 10
- Pollution per trip for electric vehicles: 0

### Decision Variables

- $(T)$ : Number of trips using trucks
- $(V)$ : Number of trips using vans
- $(M)$ : Number of trips using motorcycles
- $(E)$ : Number of trips using electric vehicles
- $(B)$ : Binary variable (1 if van is used, 0 if electric vehicle is used)

### Objective Function

- Minimize total pollution:  $(\text{Minimize } 100T + 50V + 10M)$

### Constraints

1. Total product transported:  $(100T + 80V + 40M + 60E \geq 1800)$
2. Total pollution:  $(100T + 50V + 10M \leq 2000)$
3. At least 10 trips using trucks:  $(T \geq 10)$
4. Binary variable constraint for van or electric vehicle:  $(V \leq B)$  and  $(E \leq 1 - B)$
5. Non-negativity constraints:  $(T, V, M, E, B \geq 0)$

</model>

## Case Study 2: Generated response under the Without KL strategy

<think>

To solve this problem, we need to identify the decision variables, the objective function, and the constraints.

1. **Decision Variables:**

- Let  $(T)$  be the number of trips made by trucks.
- Let  $(V)$  be the number of trips made by vans.
- Let  $(M)$  be the number of trips made by motorcycles.
- Let  $(E)$  be the number of trips made by electric vehicles.

2. **Objective Function:**

The objective is to minimize the total pollution generated, which can be expressed as:

$$\text{Minimize } 100T + 50V + 10M$$

3. **Constraints:**

- The total amount of product transported must be at least 1800 units:  $100T + 80V + 40M + 60E \geq 1800$
- The total pollution generated cannot exceed 2000 units:  $100T + 50V + 10M \leq 2000$
- At least 10 trips must use a truck:  $T \geq 10$
- Non-negativity constraints:  $T, V, M, E \geq 0$  (integer values)

Since the problem involves integer values for the decision variables, we will use integer linear programming (ILP) to solve it.

</think>

<model>

### Sets

- There are no sets in this problem as it is a simple linear programming problem.

### Parameters

- The product to be transported: 1800 units.
- Pollution generated by each type of vehicle: truck (100), van (50), motorcycle (10), electric vehicle (0).
- Product transported by each type of vehicle: truck (100), van (80), motorcycle (40), electric vehicle (60).
- Maximum pollution: 2000 units.
- Minimum truck trips: 10.

### Decision Variables

- $(T)$ : Number of trips made by trucks.
- $(V)$ : Number of trips made by vans.
- $(M)$ : Number of trips made by motorcycles.
- $(E)$ : Number of trips made by electric vehicles.

### Objective Function

$$\text{Minimize } 100T + 50V + 10M$$

### Constraints

1. Product transportation constraint:  $100T + 80V + 40M + 60E \geq 1800$
2. Pollution constraint:  $100T + 50V + 10M \leq 2000$
3. Minimum truck trips:  $T \geq 10$
4. Non-negativity constraints:  $T, V, M, E \geq 0$  (integer values)

</model>

## Case Study 2: Generated response under the Partial KL strategy

<think>

To solve this problem, we need to define the decision variables, objective function, and constraints based on the given transportation options and pollution limits. Here's a step-by-step approach:

1. **Decision Variables:**

- Let  $x_1$  be the number of trips with trucks.
- Let  $x_2$  be the number of trips with vans.
- Let  $x_3$  be the number of trips with motorcycles.
- Let  $x_4$  be the number of trips with electric vehicles.
- Let  $y$  be a binary variable where  $y = 1$  if a van is chosen, and  $y = 0$  if an electric vehicle is chosen.

2. **Objective Function:**

- The objective is to minimize the total pollution generated. The pollution generated by trucks is 100 units per trip, by vans is 50 units per trip, by motorcycles is 10 units per trip, and by electric vehicles is 0 units per trip.

- Therefore, the objective function is:

$$\text{Minimize } 100x_1 + 50x_2 + 10x_3 + 0x_4$$

3. **Constraints:**

- The total product transported must be at least 1800 units.
- The total pollution generated must be at most 2000 units.
- At least 10 trips must use a truck.
- Only one of the van or electric vehicle can be chosen.

</think>

<model>

### Sets

- $I = \{1, 2, 3, 4\}$  where  $I$  is the set of transportation options.

### Parameters

- $P_i$  is the pollution generated by transportation option  $i$ .
- $C_i$  is the capacity of transportation option  $i$ .

### Decision Variables

- $x_1$  is the number of trips with trucks.
- $x_2$  is the number of trips with vans.
- $x_3$  is the number of trips with motorcycles.
- $x_4$  is the number of trips with electric vehicles.
- $y$  is a binary variable where  $y = 1$  if a van is chosen, and  $y = 0$  if an electric vehicle is chosen.

### Objective Function

$$\text{Minimize } 100x_1 + 50x_2 + 10x_3 + 0x_4$$

### Constraints

$$100x_1 + 80x_2 + 40x_3 + 60x_4 \geq 1800 \quad \text{(Total product constraint)}$$

$$100x_1 + 50x_2 + 10x_3 + 0x_4 \leq 2000 \quad \text{(Total pollution constraint)}$$

$$x_1 \geq 10 \quad \text{(At least 10 trips with trucks)}$$

$$y \in \{0, 1\} \quad \text{(Binary variable for van or electric vehicle)}$$

$$x_2 \leq 100y \quad \text{(Van constraint)}$$

$$x_4 \leq 100(1 - y) \quad \text{(Electric vehicle constraint)}$$

</model>