

OPTIVER: UNLEASHING THE POWER OF LLMs FOR OPTIMIZATION MODELING VIA DUAL-SIDE VERIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Building mathematical optimization models is critical in operations research (OR), while it requires substantial human expertise. Recent advancements have utilized large language models (LLMs) to automate this modeling process. However, existing works often struggle to verify the correctness of the generated optimization models, without checking the rationality of the constraints and variables or the validity of solutions to the generated models. This hampers the subsequent verification and correction steps, and thus it severely hurts the modeling accuracy. To address this challenge, we propose a novel LLM-based framework with Dual-side Verification (OptiVer) from both structure and solution perspectives, thereby improving the modeling accuracy. The structure-side verification ensures that the modeling structure of the generated optimization models aligns with the original problem description, accurately capturing the problem’s constraints and requirements. Meanwhile, the solution-side verification interprets and evaluates the validity of the solutions, confirming that the optimization models are logically and mathematically sound. Extensive experiments on several popular benchmarks demonstrate that our approach significantly outperforms the state-of-the-art, achieving over 20% improvement in accuracy.

1 INTRODUCTION

Optimization problems are foundational to operations research (OR), with wide-ranging applications in manufacturing (Jayal et al., 2010), transportation (Yin, 2002), and service industries (Berman et al., 1994). In practice, OR problem statements are typically specified in natural language. Practitioners must therefore (i) translate these descriptions into an appropriate mathematical optimization model (defining objectives, decision variables, and constraints) and (ii) implement the solver code (e.g., SCIP (Achterberg, 2009), Gurobi (Gurobi Optimization, 2021), or Pyomo (Bynum et al., 2021; Hart et al., 2011)) to obtain solutions. This workflow is labor-intensive, demands substantial domain expertise in problem context, mathematical modeling, and code-level implementation or debugging, and is consequently costly and time-consuming (Ahmaditeshnizi et al., 2024).

Given the impressive capabilities of large language models (LLMs) in natural-language understanding and domain knowledge acquisition, a growing number of works have employed LLMs to automate the processes of modeling, programming, and debugging. Existing approaches can be broadly grouped into two categories. The first category, prompt-based methods, relies on pre-trained LLMs (e.g., GPT-4 (OpenAI, 2023) and GPT-4o (OpenAI, 2024)), which are prompted to construct mathematical models incrementally. In practice, these methods are often implemented into a carefully designed framework, such as multi-agent cooperation (Xiao et al., 2024; Ahmaditeshnizi et al., 2024) and Monte Carlo tree search (Astorga et al., 2025). The second category enhances modeling capabilities through fine-tuning, which involves constructing a large, labeled dataset for LLM training (Huang et al., 2025; Wu et al., 2025; Jiang et al., 2025; Chen et al., 2025; Lu et al., 2025).

Beyond these two approaches, recent research has explored self-correction strategies to improve the modeling performance (Jiang et al., 2025; Xiao et al., 2024; Ahmaditeshnizi et al., 2024). These methods trigger correction primarily from error messages produced during code execution. However, such strategies confine self-correction to code-level issues, and the underlying model can remain

054 flawed even when the code runs without failure. To develop more effective model verification
 055 methods, we characterizes incorrect models by the following features. First, incorrect models often
 056 overlook indispensable constraints that the textual problem description does not explicitly state. For
 057 example, when formulating a maximum flow problem, an LLM might omit flow-balance constraints
 058 at intermediate nodes simply because they are not explicitly mentioned, yielding a structurally
 059 incomplete model. Second, incorrect models can yield solutions that are infeasible or violate basic
 060 logical principles, even though the solver executes successfully and reports an improved objective.

061 To address these challenges, we propose a novel
 062 multi-agent framework with Dual-side Verification (OptiVer) from both the modeling structure
 063 and solution perspectives, improving the modeling
 064 accuracy. This approach moves beyond simple
 065 code-execution signals by translating both the
 066 optimization model and its resulting solutions
 067 into natural language and evaluating their
 068 semantic correctness. To the end, OptiVer introduces two novel evaluation metrics for self-
 069 correction: modeling structure consistency and
 070 solution validity. (1) **Consistency in the modeling structure** ensures that the mathematical
 071 formulation (its variables, constraints, and objective) is a complete and faithful translation of
 072 the original problem description. To evaluate
 073 this metric, one LLM agent performs a back-
 074 translation that abstracts the generated model
 075 into a compact, multi-level description of its
 076 components. A second agent then aligns this
 077 abstraction with the structure derived from the
 078 original specification to reveal omissions or
 079 mismatches. (2) **Solution validity** assesses whether the solution of the produced model is logically
 080 and contextually sound for the real-world task. To assess it, one agent interprets the numeric
 081 solution in natural language, explaining its meaning in the context of the real-world problem. Another
 082 agent then critiques that interpretation to expose logical absurdities or mathematical violations that
 083 code-execution checks miss. Finally, we use the verification feedback for model refinement.

086 As illustrated in Figure 1, extensive experiments on five popular benchmarks showcase that our
 087 approach significantly outperforms the state-of-the-art, achieving an average improvement of approxi-
 088 mately 10% in solving accuracy. Notably, OptiVer is designed as a plug-and-play framework, capable
 089 of effectively verifying and refining optimization models generated by any existing pre-trained or
 090 fine-tuned OR LLMs.

092 **2 RELATED WORK**

094 **Automated Optimization modeling** In practice, the OR problems often arise from real-world
 095 situations, which are typically described in natural language. Consequently, automated optimization
 096 modeling has emerged as a critical area aimed at reducing the labor and time costs associated
 097 with the modeling process (Chen et al., 2023; Li et al., 2023). Notable early efforts in this field
 098 include the NL4Opt competition (Ramamonjison et al., 2021). Since then, several benchmarks
 099 have been introduced to evaluate performance, such as ComplexOR (Xiao et al., 2024), NLP4LP
 100 (Ahmaditeshnizi et al., 2024), Mamo (Huang et al., 2024), IndustryOR (Huang et al., 2025) and
 101 Optibench Yang et al. (2025); Wang et al. (2024). Recent research primarily falls into two categories:
 102 prompt-based methods and fine-tuned methods. Prompt-based approaches utilize pre-trained large
 103 language models (LLMs) with carefully crafted prompts to iteratively construct models. For instance,
 104 Chain-of-Experts (Xiao et al., 2024) and OptiMUS (Ahmaditeshnizi et al., 2024) frameworks employ
 105 multi-agent cooperation, while some other methods Astorga et al. (2025) leverage Monte Carlo tree
 106 search techniques to explore potential models. To further enhance the modeling capabilities of LLMs,
 107 researchers also work on fine-tuning these models with extensive OR and modeling knowledge
 (Huang et al., 2025; Wu et al., 2025; Jiang et al., 2025; Chen et al., 2025). For example, LLaMoCo
 (Ma et al., 2024) utilizes an instruction tuning framework to adapt LLMs for solving optimization

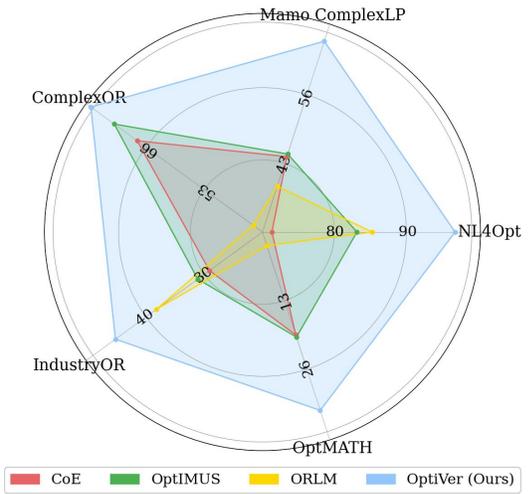


Figure 1: OptiVer outperforms other baselines in solving accuracy (SA) across the benchmarks.

problems in a code-to-code manner. ORLM (Huang et al., 2025) trains open-source LLMs specifically designed for optimization modeling and solver code development. Additionally, advanced techniques such as KTO (Ethayarajh et al., 2024) and data augmentation have been introduced to improve model training (Wu et al., 2025; Jiang et al., 2025).

3 MOTIVATED RESULTS AND CASE ANALYSIS

Challenges We identify two key challenges in existing LLM-based optimization modeling methods. (1) LLMs struggle to identify the modeling structure within the problems, including missing or incorrect constraints and errors in variable definitions. The prevalence of such mistakes is notable in benchmarks, with 36.0% in NL4Opt and 12.6% in ComplexOR as pointed out in OptiMUS (Ahmaditeshnizi et al., 2024). (2) LLMs can only find the errors in the solver codes, but can hardly find the errors in the optimization models. OptiMUS points out that ‘‘Coding errors are easier to identify and fix. In contrast, identifying bugs in the formulation requires deeper reasoning and is harder.’’ In existing methods, the debugging module is typically activated only when solver codes produce execution errors. To illustrate these challenges, we use GPT4o-mini in Figure 2, involving a maximum flow problem (MF).

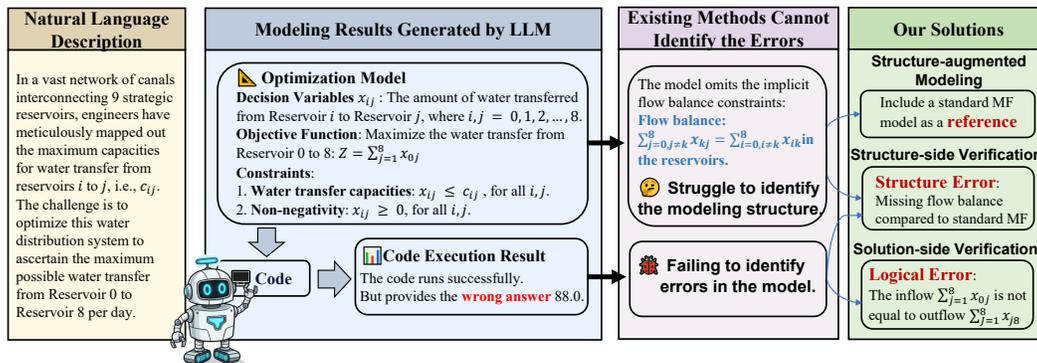


Figure 2: The two challenges we observed in existing optimization modeling methods.

Observations on the Modeling Structures To specify the definition and the usage of modeling structures, we have the following observation. The LLM cannot find the flow balance constraints at first. However, the model can correctly identify the relevant problem classifications. When we prompt the model to formulate the relevant problem classification (Maximum Flow Problem in this case), it successfully identifies the flow balance constraints.

The core principle of structure-augmented modeling is to **leverage similar standard optimization models as a reference** to identify a problem’s implicit constraints. In Operations Research, many problems in similar scenarios share characteristics with optimization models in conventional problem classifications, such as the Vehicle Routing Problem or the Maximum Flow Problem. These classic types have conventional mathematical formulations—including standard variables and assumed constraints—which this paper refers to as **modeling structures**. Even when a new problem does not neatly fit a standard problem classification, referencing the modeling structure of a similar, well-understood problem helps the LLM uncover these implicit relationships, which are often crucial for a correct formulation.

4 METHODOLOGY

Our work investigates how the modeling process can be enhanced through effective verification methods on both the structural and solution sides. An overview of the framework is presented in Figure 4. We define multi-level modeling structures in Section 4.1, followed by a detailed explanation of

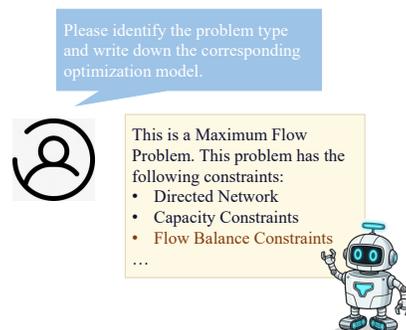


Figure 3: The influence of recalling a problem classification.

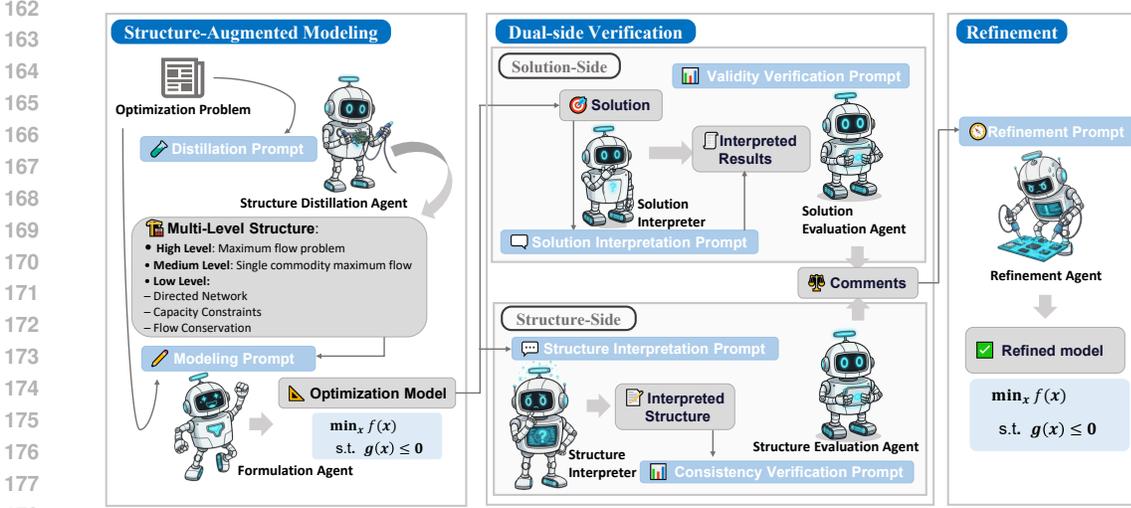


Figure 4: Our OptiVer framework begins by distilling the multi-level structures from the natural language description. These extracted structures are then combined, allowing the formulator to generate an initial model. Then, OptiVer conducts a dual-side verification and refinement process.

structure-side verification in Section 4.2 and solution-side verification using a multi-agent cooperation framework in Section 4.3. We first introduce some notations in this work as follows.

Let \mathcal{D} represent the space of natural problem descriptions, and let \mathcal{M} denote the model space encompassing all possible optimization models. The modeling process can be viewed as a mapping from the problem description $D \in \mathcal{D}$ to an optimization model $M \in \mathcal{M}$. The descriptions follow the distribution $p(D)$, and we explicitly define the modeling space as the Cartesian product of the variable space, constraint space, and objective function space. Let $p(M)$ be the distribution representing the optimization model. The mutual information is a widely used theoretical metric for the similarity of the distributions. The mutual information $I(D, M)$ is a canonical metric used to quantify the dependence between such distributions. In information theory, mutual information is defined as

$$I(M, D) = \sum_{M \in \mathcal{M}, D \in \mathcal{D}} p(M, D) \log \left(\frac{p(M, D)}{p(M)p(D)} \right) \quad (1)$$

A high mutual information indicates that the model accurately captures the constraints and requirements stated in the text. Conversely, information loss (low mutual information) corresponds to common modeling errors, such as missing implicit constraints or incorrect variable definitions.

4.1 STRUCTURE-AUGMENTED MODELING

(1) Motivation of Multi-Level Structure: Coarse-to-fine structure Analysis Before the modeling process, human experts first analyze the problem description to identify a similar conventional problem classification as a reference. Next, they determine the variant of the classification that best aligns with the description. Finally, they assess special requirements in the description. This analysis follows a coarse-to-fine approach, from high-level to low-level structure analysis.

(2) Multi-Level Modeling Structure Our framework begins by distilling the modeling structures from the natural language descriptions. As discussed in Section 3, understanding the problem type is crucial in the modeling process, as it serves as a foundational template for developing optimization models. Inspired by the coarse-to-fine structure analysis process used by human experts, we further refine the concept of modeling structures by introducing the idea of multi-level modeling structures.

- **High-Level Structure:** This represents the fundamental problem type within OR, such as the maximum flow problem, set covering problem, vehicle routing problem, and knapsack problem. Each of these problem types is associated with a basic optimization model.

- **Medium-Level Structure:** This pertains to the classical classification or variants of fundamental problem types. For instance, variations of the maximum flow problem include multi-source MF, multi-commodity MF, minimum-cost MF, and MF in undirected graphs. Each variant is associated with a specific modified optimization model derived from the basic model.
- **Lower-Level Structure:** This level encompasses constraints in the classical optimization model as well as specific requirements that extend beyond classical models. For instance, standard constraints might include capacities and flow balance constraints in the MF, while special requirements could involve flow capacities that fluctuate over time.

As we mentioned in Section 3, even highly complex and unique industrial problems are often variants or combinations of fundamental problem classifications recognized in operations research. The **high-level** and **medium-level structures** in our framework are designed to capture this foundational core, providing a solid starting point for the modeling process. Second, and most critically, the framework’s **low-level structure** is specifically designed to provide the necessary flexibility to handle unique, real-world contexts. This level is not confined to a specific problem classification and is intended to capture the nuanced, problem-specific constraints and requirements that extend beyond classical formulations. This design allows the framework to represent the unique aspects of any given problem, rather than forcing it into a rigid, predefined category. We denote the multi-level modeling structure as S . Below, we provide an example of the modeling structure we have defined.

Example: The Structure Schema Extracted by LLMs

- **High Level:** Maximum flow problem
- **Medium Level:** Single commodity maximum flow
- **Low Level:**
 - **Directed Network:** The flow is directed from one reservoir to another.
 - **Capacity Constraints:** Each edge (connection between reservoirs) has a maximum capacity.
 - **Flow Conservation:** The amount of water entering any intermediate reservoir must equal the amount leaving, except for the source and sink.

(3) Structure Distillation and Structure-Augmented Modeling We use two pre-trained LLMs (implemented by GPT4o-mini in this work) as agents to complete the structure distillation and initial modeling tasks, guided by designed prompts. To distill the multi-level modeling structure S from the natural language description D , we introduce an LLM agent, called the structure distillation agent. The agent takes as input the problem description and outputs the formatted structure context. Then, we call a formulation agent to generate an initial optimization model M guided by prompts combining the problem description and modeling structure, i.e.,

$$S = \text{Distillation_Agent}(D), \quad M = \text{Formulation_Agent}(D, S). \quad (2)$$

4.2 STRUCTURE-SIDE: STRUCTURE INTERPRETATION AND CONSISTENCY VERIFICATION

(1) Motivation Structure-side verification finds the modeling errors by detecting any deviation from the established, correct formulation for a known class of problems, catching errors of omission where the LLM may overlook fundamental constraints and variables required for that problem classification. Inspired by dual learning in machine translation (He et al., 2016), we assert that a correct model must meet the following consistency criterion: when we translate the optimization model back into the space of modeling structure, the resulting context should semantically correspond to the modeling structure directly derived from the problem description.

(2) Structure Interpretation and Consistency Verification We introduce a structure interpretation agent and an evaluation agent to complete the structure verification task. The two agents are also guided with specific prompts. First, a structural interpretation agent performs a "back-translation". It takes the generated mathematical model M and converts it back into its abstract modeling structure \tilde{S} . Next, a structural evaluation agent acts as a critic. It compares the interpretation agent’s output \tilde{S} with the original structure S derived from the problem description to check for semantic consistency.

The evaluation agent’s output is twofold: a binary consistency score c_c and a detailed comment that highlights any discrepancies. This comment provides specific, actionable feedback that is later used to refine the model. The process can be formally summarized as

$$\tilde{S} = \text{StruInterp_Agent}(M), \quad (com, c_c) = \text{StruEval_Agent}(S, \tilde{S}), \quad (3)$$

where $c_c = 1$ indicates consistency, while the comment com details any differences found between the structures, guiding the subsequent refinement step.

Finally, we propose an analysis of the structure-side verification using mutual information. Suppose that \tilde{S} is the interpreted structure from optimization model M . The structure-side verification aims to improve the consistency between structures S from the natural language description and that \tilde{S} interpreted from the optimization model, i.e., the mutual information $I(S, \tilde{S})$.

Proposition 4.1. *We have $I(S, \tilde{S}) \leq I(D, M)$. Thus, the structure-side verification optimizes the lower bound of the mutual information between the problem description and the optimization model.*

The structure-side verification is improving the lower bound.

Proposition 4.2. *Let $\epsilon > 0$ be a fixed threshold representing the minimum information overlap required for a model to be deemed structurally consistent. Define the Bernoulli random variable c_c corresponding to the binary output of our Structure Evaluation Agent, $c_c \triangleq \mathbf{1}_{\{I(S, \tilde{S}) \geq \epsilon\}}$, where $I(S, \tilde{S})$ denotes the mutual information between the distilled structure S and the interpreted structure \tilde{S} . Then, the expected mutual information satisfies, $\mathbb{E}[I(S, \tilde{S})] \geq \epsilon \mathbb{E}[c_c]$.*

4.3 SOLUTION-SIDE: SOLUTION INTERPRETATION AND VALIDITY VERIFICATION

(1) Motivation This method works because it grounds the abstract mathematical model by assessing whether its solution is logically feasible. A model may be syntactically correct and yield a numerical answer, yet that answer could violate the fundamental logic of the original problem (e.g., suggesting more water flows out of a reservoir than flows in). We argue that the semantic content of the solution itself is a far richer source for identifying errors. Solution-side verification enhances performance because it is designed to catch logical errors that are invisible to systems that only check for solver execution errors. The core of this verification is to leverage the **common-sense and logical reasoning capabilities** of such LLMs for improved error detection.

(2) Solution Interpretation and Validity Verification Given an optimization model M , OptiVer executes the solver code and obtains the optimal solution \mathbf{x} . Then, OptiVer performs solution-side verification using two LLM agents, guided by designed prompts. The first agent, a solution interpreter, translates the raw numerical solution \mathbf{x} into a meaningful natural language description \tilde{D} based on the original problem context D . Next, the second agent, a solution evaluation agent, scrutinizes this description to identify any logical or mathematical errors. This agent’s output includes a binary validity score c_v and, crucially, a detailed comment com that provides specific feedback on any flaws found. The score value is 1 if the evaluator recognizes the validity of solution \mathbf{x} , and 0 otherwise. This verification process can be formally summarized as:

$$\tilde{D} = \text{SolInterp_Agent}(\mathbf{x}, D), \quad (com, c_v) = \text{SolEval_Agent}(D, \tilde{D}). \quad (4)$$

Similar to the analysis of the structure-side verification, we have the following analysis of the solution-side verification. Suppose that \tilde{D} is the interpreted solution from the optimization model M . During the solution-side verification, we improve the mutual information $I(D, \tilde{D})$.

Proposition 4.3. *We have $I(D, \tilde{D}) \leq I(D, M)$. Thus, the solution-side verification optimizes the lower bound of the mutual information between the problem description and the optimization model.*

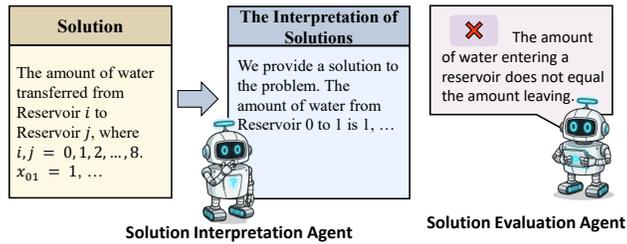


Figure 5: An example of solution verification.

Table 1: Comparison of our method and the baselines across five popular benchmarks. Throughout the experiments, we compare the solving accuracy (SA) of the methods.

	NL4Opt	Mamo ComplexLP	ComplexOR	IndustryOR	OptMATH
Reasoning LLMs					
DeepSeek-R1	82.6	67.2	68.4	32.0	33.1
OpenAI-o1	87.1	66.3	68.4	36.0	32.5
Fine-tuned Method					
ORLM	85.1*	38.8*	42.1*	38.0*	2.6*
Evo-Step	84.4*	61.6*	-	36.3*	-
LLMOPT	80.3*	44.1*	72.7*	29.0*	12.5*
OptMATH	95.9*	54.1*	-	-	34.9*
SIRL	96.3*	62.1*	-	33.0*	29.0*
Prompt-based Method					
Standard	64.6	27.9	31.5	24.0	15.6
CoT	69.3	34.5	36.8	27.0	18.6
CoE	71.3	44.5	68.4	29.0	19.8
OptiMUS	83.0	45.0	73.6	31.0	20.2
OptiVer (Ours)	96.5	66.7	78.9	45.0	34.3

Values marked with * are from the original or reproduced papers. , and - are with missing data because the model has not been publicly released.

Proposition 4.4. *Let $\epsilon' > 0$ be a threshold for solution validity and define $c_v \triangleq \mathbf{1}_{\{I(D, \tilde{D}) \geq \epsilon'\}}$, representing the binary output of the Solution Evaluation Agent. Then $\mathbb{E}[I(D, \tilde{D})] \geq \epsilon' \mathbb{E}[c_v]$.*

4.4 REFINEMENT

Based on the feedback, we refine the optimization model. The refinement agent within the OptiVer framework is a specialized LLM-based component responsible for correcting and enhancing the initial optimization model based on insights from the dual-side verification process. Guided by a refinement prompt, the agent takes the current formulation as input and produces a refined optimization model, represented as $M' = \text{Ref_Agent}(D, S, M, com)$.

5 EXPERIMENTS

Benchmarks We use five real-world operations research benchmarks: NL4Opt (Ramamonjison et al., 2021), Mamo ComplexLP (Huang et al., 2024), ComplexOR (Xiao et al., 2024), IndustryOR (Huang et al., 2025) and OptMATH (Lu et al., 2025). The NL4Opt benchmark, released for the NeurIPS 2022 NL4Opt competition, consists of 289 elementary linear programming problems. Mamo ComplexLP 211 problems. ComplexOR is a comprehensive dataset including linear and mixed-integer programming. In alignment with the studies by (Ahmaditeshnizi et al., 2024) and (Jiang et al., 2025), we focus on 19 specific problems from this dataset. IndustryOR has 100 challenging problems from various industry scenarios. OptMATH has 166 challenging problems.

Implementation and Baselines In our experiments, we utilized the GPT4o-mini to implement the agents in our method and all the prompt-based baselines. For the implementation of OptiVer, please see Appendix F for the prompts of each agent. In our experiments, we compare OptiVer with four available prompt-based methods and five fine-tuned operations research LLMs. The four prompt-based baselines include Standard, Chain-of-Thoughts (CoT) (Wei et al., 2022), Chain-of-Experts (CoE) (Xiao et al., 2024), and OptiMUS (Ahmaditeshnizi et al., 2024). The Standard baseline represents the output of GPT without any optimization of its reasoning processes. We include five fine-tuned open-source operations research language models as baselines, including ORLM (Huang et al., 2025) (based on LLaMA-3-8B model), Evo-Step (Wu et al., 2025) (based on LLaMA-3-8B model), LLMOPT (Jiang et al., 2025) (based on Qwen1.5-14B), OptMATH (Lu et al.,

Table 2: Ablation studies on (1) each component and (2) each level of modeling structures in OptiVer.

Method	NL4Opt	Mamo	ComplexLP	ComplexOR	IndustryOR
Ablation for the Components					
OptiVer w/o struaug	91.8		54.2	63.1	34.0
OptiVer w/o stru-side	91.5		55.2	68.4	29.0
OptiVer w/o sol-side	91.8		53.1	68.4	41.0
Ablation for Each Level of the Structure					
OptiVer w/o high	92.9		59.9	78.9	41.0
OptiVer w/o medium	91.1		57.0	73.6	38.0
OptiVer w/o low	91.1		54.9	73.6	30.0
OptiVer (full)	96.5		66.7	78.9	45.0

2025) (based on Qwen2.5-32B), and SIRL (Chen et al., 2025) (based on Qwen2.5-7B) trained with reinforcement learning. Additionally, we also compare our results with the pre-trained reasoning model DeepSeek-R1 (DeepSeek-AI, 2025) and OpenAI-o1 (OpenAI, 2024).

Metrics Consistent with existing research, we employed solving accuracy (SA) to evaluate performance. Specifically, SA represents the proportion of problems for which the methods successfully identify the optimal solutions. The higher value of SA implies better performance.

5.1 MAIN RESULTS

To demonstrate the effectiveness of our method, we conduct experiments comparing solving accuracy (SA) between our approach and baseline methods across various benchmarks. The results presented in Table 1 indicate that our method significantly outperforms the baselines, achieving an approximate 20% improvement in solving accuracy compared to Standard. For the challenging benchmarks, our method consistently delivers outstanding performance. This demonstrates that OptiVer exhibits strong generalization capabilities across both easy and difficult scenarios. Furthermore, OptiVer achieves performance better than state-of-the-art reasoning LLMs, such as DeepSeek-R1 and OpenAI-o1, despite relying on a much weaker base model, GPT4o-mini. Please see Appendices D and E for case and error analysis.

5.2 ABLATION STUDIES

(1) The Effects of Each Component of OptiVer In this section, we examine the effects of the three components of OptiVer: structure-augmented modeling, structure-side verification, and solution-side verification. To assess their contributions, we implement three variants of OptiVer. The first variant, OptiVer w/o stru-aug, omits the introduction of a modeling structure to enhance the modeling process. For structure-side verification, instead of interpreting the model in structural terms, we instruct an LLM agent to provide a narrative explaining the meaning of the variables, constraints, and objectives. The second variant, OptiVer w/o stru-side, does not implement structure-side verification at all. The third variant, OptiVer w/o sol-side, excludes the solution-side verification process. The results, presented in Table 2, reveal a significant drop in performance in the absence of any of these components, highlighting their essential roles in the modeling process.

(2) The Effects of Each Level of Modeling Structures Next, we investigate the impact of each level within our proposed modeling structure. The variant OptiVer w/o high/medium/low level excludes the use of high, medium, and low-level structures. The experimental results in Table 2 demonstrate that all three levels contribute positively to overall performance, with the medium and low-level structures showing particularly pronounced improvements.

Takeaway Critically, the framework’s “low-level structure” is specifically designed to provide the necessary flexibility to handle unique, real-world contexts. This level is not confined to a specific problem type and is intended to capture the nuanced, problem-specific constraints and requirements that extend beyond classical formulations. This design allows the framework to represent the unique aspects of any given problem, rather than forcing it into a rigid, predefined category.

Table 3: We build OptiVer on different baselines and backbone models.

Method	NL4Opt	Mamo ComplexLP	ComplexOR	IndustryOR
Different Baselines				
ORLM	78.2	41.2	36.8	39.0
ORLM+OptiVer	92.3	59.6	73.6	42.0
OptiMUS	83.0	45.0	73.6	31.0
OptiMUS+OptiVer	96.1	61.0	78.9	45.0
Different Backbones				
GPT-4o	79.4	45.0	57.8	27.0
GPT-4o+OptiVer	97.5	66.3	78.9	48.0
Qwen2.5-14B	70.3	41.2	57.8	31.0
Qwen2.5-14B+OptiVer	85.8	56.3	68.4	39.0

5.3 BUILDING ON DIFFERENT BASELINES AND LLMs

(1) Improving different Baselines: OptiVer was applied to the outputs of three foundational baselines: OptiMUS and the fine-tuned ORLM model. In each case, OptiVer’s verification and refinement process enhanced the initial models generated by these baseline methods. **(2) Improving different Base LLMs:** To illustrate that the framework is not reliant on a specific backbone model, we conducted experiments using various base LLMs with OptiVer. This approach highlights how performance scales with the capabilities of the underlying model, including stronger models (e.g., GPT-4o) and weaker models (e.g., Qwen2.5-14B) The results consistently indicated significant performance gains, as shown in Table 3. Please refer to Appendix C for detailed experiment settings.

5.4 QUANTITY ANALYSIS OF VERIFICATIONS

Critical Components of Verifications The interpretation and evaluation agents are essential components of the verification process, as they determine whether OptiVer can effectively identify errors in the modeling process. We conducted extensive ablation studies to quantitatively assess the accuracy and reliability of these agents. Our experiments were specifically designed to evaluate their ability to distinguish between correct and incorrect models.

Table 4: Verification Precision

Verification Type	Easy	Medium	Hard
Structure Verification	92%	89%	83%
Solution Verification	93%	91%	86%

Table 5: Verification Recall

Verification Type	Easy	Medium	Hard
Structure Verification	86%	79%	68%
Solution Verification	83%	85%	73%

Experiment design We utilized the IndustryOR dataset for evaluation, which consists of three difficulty levels (easy, medium, and hard) that allow us to test the generalization capabilities of OptiVer across varying problem complexities. *The hard problems can be general problems with complex structures that fall out of the conventional problem classifications.* However, this analysis was labor-intensive, as the IndustryOR dataset does not provide detailed, step-by-step ground-truth labels necessary for our analysis. To ensure the correctness of this evaluation, we resorted to a manual checking process, which is time-consuming. We first manually annotated the optimization models for the sampled problems to establish a ground truth. To facilitate our evaluation, we randomly selected ten problems from each difficulty level. For generating incorrect models, we initially labeled the models and extracted the structures. We then created nine negative samples for each labeled model by randomly deleting or rewriting some of the variables and constraints. This resulted in 30 positive and 270 negative modeling samples. For **structure evaluation**, we collected the interpreted structures from both the positive and negative samples. The evaluator compared these interpreted structures with the ground-truth structures and generated a binary score. For **solution evaluation**, we used the positive and negative samples to generate solutions, which we then interpreted and assessed for reliability.

Table 6: Inference efficiency comparison of different methods.

Metric	Method	NL4Opt	MAMO ComplexLP	ComplexOR	IndustryOR
Inference Time	CoE	58.2	72.5	98.6	79.9
	OptiMUS	64.2	80.3	101.3	88.2
	Ours	52.8	67.6	68.2	59.4
Agent Calls	CoE	13.6	15.7	16.6	14.1
	OptiMUS	10.4	13.8	15.3	13.9
	Ours	9.0	9.1	9.4	9.2
Token Usage	CoE	6745.8	7705.0	9469.8	8751.5
	OptiMUS	7039.4	8248.4	10796.1	9062.7
	Ours	5320.7	7385.2	8457.4	6819.8

Table 7: The token usage of each step.

Step	NL4Opt	MAMO ComplexLP	ComplexOR	IndustryOR
Structure Distillation	492.1	441.1	665.4	501.1
Modeling	1175.5	1944.8	2147.1	1564.7
Structure-side Verification	474.6	597.4	594.2	662.7
Solution-side Verification	302.4	422.1	521.3	347.4
Refinement	728.7	915.5	1013.5	759.0
Coding and Debugging	2147.4	3064.4	3515.9	2984.9

Results The evaluation accuracy and recall rates are presented in Tables 4 and 5. For each difficulty level, we evaluated 10 positive samples and 90 negative samples. Both precision and recall rates for the negative samples are high across the difficulty levels, demonstrating the reliability of the scores. We find that the verification process still performs well for hard problems that cannot be classified into a specific problem type, indicating the strong generalization to general problems.

5.5 COMPARISON OF SOLVING EFFICIENCY

Efficiency Definition We examine the solving efficiency of OptiVer in comparison to the prompt-based baselines CoE and OptiMUS by analyzing the average time, token usage and agent calls to solve a problem. We use the **same solver (Gurobi)** for all methods. This ensures fairness in efficiency comparisons. The solving time in Table 6 contains the modeling time using LLMs and the execution time of the solver. The solver execution time is short (under 0.01 seconds) and can be neglected during this process. Thus, the solving time in Table 6 reflects the modeling time by LLMs. The results presented in Table 6 indicate that OptiVer achieves significantly shorter solving times, showcasing its high efficiency. Furthermore, we analyzed the internal cost of our verification steps. The key finding in Table 7 is that the extra verification steps are computationally friendly. The modeling and coding part takes over 60% of the tokens.

The Reason why OptiVer is efficient Compared to other prompt-based baselines, **OptiVer has a simpler workflow**. CoE and OptiMUS are based on the multi-agent cooperation framework. The workflow of these methods is automatically controlled by a management agent. The insufficient decision-making ability of the management agent may lead to suboptimal decision chains. In the experiments, we find that this method may repeatedly call the same agent. For example, for certain complex problems, CoE may call the terminology interpreter again and again. In contrast, OptiVer does not include such a management agent, leading to a simpler workflow.

6 CONCLUSION

In this paper, we propose an LLM-based verification framework designed to enhance the accuracy of automated mathematical modeling tasks. In the structure-side verification, we assess the modeling structures of the current model to ensure structural consistency. Meanwhile, in the solution-side verification, we interpret the solution within the context of the problem descriptions, aiming to identify any logical or mathematical errors in the models. Extensive experiments demonstrate the effectiveness of our method across a wide range of benchmarks.

540 ETHICS STATEMENT.

541
542 This work is designed to explore the significance of the verification process in LLM optimization
543 modeling. We do not foresee any direct, immediate, or negative societal impacts of our research.
544

545 REPRODUCIBILITY STATEMENT.

546
547 All the results in this work are reproducible. We have discussed the implementation details in Section
548 5. We also present our prompts for each agent in Appendix F.
549

550 REFERENCES

- 551
552 Tobias Achterberg. Scip: solving constraint integer programs. *Mathematical Programming Computa-*
553 *tion*, 1:1–41, 2009.
554
- 555 Ali Ahmaditeshnizi, Wenzhi Gao, and Madeleine Udell. OptiMUS: Scalable optimization modeling
556 with (MI)LP solvers and large language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine
557 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings*
558 *of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine*
559 *Learning Research*, pp. 577–596. PMLR, 21–27 Jul 2024. URL [https://proceedings.](https://proceedings.mlr.press/v235/ahmaditeshnizi24a.html)
560 [mlr.press/v235/ahmaditeshnizi24a.html](https://proceedings.mlr.press/v235/ahmaditeshnizi24a.html).
- 561 Nicolás Astorga, Tennison Liu, Yuanzhang Xiao, and Mihaela van der Schaar. Autoformulation of
562 mathematical optimization models using llms. In *International Conference on Machine Learning,*
563 *ICML*, Proceedings of Machine Learning Research, 2025.
- 564 Oded Berman, Zvi Ganz, and Janet M. Wagner. A stochastic optimization model for planning
565 capacity expansion in a service industry under uncertain demand. *Naval Research Log-*
566 *istics (NRL)*, 41(4):545–564, 1994. doi: [https://doi.org/10.1002/1520-6750\(199406\)41:](https://doi.org/10.1002/1520-6750(199406)41:4<545::AID-NAV3220410407>3.0.CO;2-Z)
567 [4<545::AID-NAV3220410407>3.0.CO;2-Z](https://doi.org/10.1002/1520-6750(199406)41:4<545::AID-NAV3220410407>3.0.CO;2-Z). URL [https://onlinelibrary.](https://onlinelibrary.wiley.com/doi/abs/10.1002/1520-6750%28199406%2941%3A4%3C545%3A%3AAID-NAV3220410407%3E3.0.CO%3B2-Z)
568 [wiley.com/doi/abs/10.1002/1520-6750%28199406%2941%3A4%3C545%](https://onlinelibrary.wiley.com/doi/abs/10.1002/1520-6750%28199406%2941%3A4%3C545%3A%3AAID-NAV3220410407%3E3.0.CO%3B2-Z)
569 [3A%3AAID-NAV3220410407%3E3.0.CO%3B2-Z](https://onlinelibrary.wiley.com/doi/abs/10.1002/1520-6750%28199406%2941%3A4%3C545%3A%3AAID-NAV3220410407%3E3.0.CO%3B2-Z).
- 570 Michael L. Bynum, Gabriel A. Hackebeil, William E. Hart, Carl D. Laird, Bethany L. Nicholson,
571 John D. Sirola, Jean-Paul Watson, and David L. Woodruff. *Pyomo—optimization modeling in*
572 *python*, volume 67. Springer Science & Business Media, third edition, 2021.
573
- 574 Hao Chen, Gonzalo E. Constante-Flores, and Can Li. Diagnosing infeasible optimization problems
575 using large language models, 2023. URL <https://arxiv.org/abs/2308.12923>.
- 576 Yitian Chen, Jingfan Xia, Siyu Shao, Dongdong Ge, and Yinyu Ye. Solver-informed rl: Grounding
577 large language models for authentic optimization modeling, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2505.11792)
578 [abs/2505.11792](https://arxiv.org/abs/2505.11792).
- 579 DeepSeek-AI. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learn-
580 ing, 2025. URL [https://github.com/deepseek-ai/DeepSeek-R1/blob/main/](https://github.com/deepseek-ai/DeepSeek-R1/blob/main/DeepSeek_R1.pdf)
581 [DeepSeek_R1.pdf](https://github.com/deepseek-ai/DeepSeek-R1/blob/main/DeepSeek_R1.pdf).
582
- 583 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Model
584 alignment as prospect theoretic optimization. In *Forty-first International Conference on Machine*
585 *Learning*, 2024. URL <https://openreview.net/forum?id=iUwHnoENnl>.
- 586 LLC Gurobi Optimization. Gurobi optimizer. URL <http://www.gurobi.com>, 2021.
587
- 588 William E Hart, Jean-Paul Watson, and David L Woodruff. Pyomo: modeling and solving mathemati-
589 cal programs in python. *Mathematical Programming Computation*, 3(3):219–260, 2011.
- 590 Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, and Wei-Ying Ma. Dual
591 learning for machine translation. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Gar-
592 nett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Asso-
593 [ciates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2016/file/5b69b9cb83065d403869739ae7f0995e-Paper.pdf)
[2016/file/5b69b9cb83065d403869739ae7f0995e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/5b69b9cb83065d403869739ae7f0995e-Paper.pdf).

- 594 Chenyu Huang, Zhengyang Tang, Shixi Hu, Ruoqing Jiang, Xin Zheng, Dongdong Ge, Benyou
595 Wang, and Zizhuo Wang. Orlm: A customizable framework in training large models for automated
596 optimization modeling. *Operations Research*, 2025. doi: 10.1287/opre.2024.1233. URL <https://doi.org/10.1287/opre.2024.1233>.
597
- 598 Xuhan Huang, Qingning Shen, Yan Hu, Anningzhe Gao, and Benyou Wang. Mamo: A mathematical
599 modeling benchmark with solvers. *CoRR*, abs/2405.13144, 2024.
600
- 601 A.D. Jayal, F. Badurdeen, O.W. Dillon, and I.S. Jawahir. Sustainable manufacturing: Modeling and
602 optimization challenges at the product, process and system levels. *CIRP Journal of Manufacturing
603 Science and Technology*, 2(3):144–152, 2010. ISSN 1755-5817. doi: [https://doi.org/10.1016/
604 j.cirpj.2010.03.006](https://doi.org/10.1016/j.cirpj.2010.03.006). URL [https://www.sciencedirect.com/science/article/
605 pii/S1755581710000131](https://www.sciencedirect.com/science/article/pii/S1755581710000131). Sustainable Development of Manufacturing Systems.
- 606 Caigao Jiang, Xiang Shu, Hong Qian, Xingyu Lu, Jun Zhou, Aimin Zhou, and Yang Yu. LLMOPT:
607 Learning to define and solve general optimization problems from scratch. In *The Thirteenth
608 International Conference on Learning Representations*, 2025. URL [https://openreview.
609 net/forum?id=9OMvtboTJg](https://openreview.net/forum?id=9OMvtboTJg).
- 610 Beibin Li, Konstantina Mellou, Bo Zhang, Jeevan Pathuri, and Ishai Menache. Large language models
611 for supply chain optimization, 2023. URL <https://arxiv.org/abs/2307.03875>.
612
- 613 Hongliang Lu, Zhonglin Xie, Yaoyu Wu, Can Ren, Yuxuan Chen, and Zaiwen Wen. Optmath:
614 A scalable bidirectional data synthesis framework for optimization modeling. In *International
615 Conference on Machine Learning, ICML*, Proceedings of Machine Learning Research, 2025.
616
- 617 Zeyuan Ma, Hongshu Guo, Jiacheng Chen, Guojun Peng, Zhiguang Cao, Yining Ma, and Yue jiao
618 Gong. LLaMoCo: Instruction tuning of large language models for optimization code generation.
619 *CoRR*, abs/2403.01131, 2024.
- 620 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
621
- 622 OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
623
- 624 OpenAI. Introducing Openai ol. <https://openai.com/ol/>, 2024.
- 625 Rindranirina Ramamonjison, Timothy T. L. Yu, Raymond Li, Haley Li, Giuseppe Carenini, Bissan
626 Ghaddar, Shiqi He, Mahdi Mostajabdaveh, Amin Banitalebi-Dehkordi, Zirui Zhou, and Yong
627 Zhang. NL4Opt competition: Formulating optimization problems based on their natural language
628 descriptions. In *NeurIPS 2022 Competition Track*, pp. 189–203, Virtual, 2021.
- 629 Zhuohan Wang, Ziwei Zhu, Yizhou Han, Yufeng Lin, Zhihang Lin, Ruoyu Sun, and Tian Ding.
630 Optibench: Benchmarking large language models in optimization modeling with equivalence-
631 detection evaluation, 2024. URL <https://openreview.net/forum?id=KD9F5Ap878>.
632
- 633 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V
634 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models.
635 In *Advances in Neural Information Processing Systems 35*, pp. 24824–24837, New Orleans, LA,
636 2022.
- 637 Yang Wu, Yifan Zhang, Yurong Wu, Yuran Wang, Junkai Zhang, and Jian Cheng. Step-opt: Boosting
638 optimization modeling in llms through iterative data synthesis and structured validation, 2025.
639 URL <https://arxiv.org/abs/2506.17637>.
- 640 Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin
641 Fu, Tao Zhong, Jia Zeng, Mingli Song, and Gang Chen. Chain-of-experts: When LLMs meet
642 complex operations research problems. In *The Twelfth International Conference on Learning
643 Representations*, 2024. URL <https://openreview.net/forum?id=HobyL1B9CZ>.
644
- 645 Zhicheng Yang, Yiwei Wang, Yinya Huang, Zhijiang Guo, Wei Shi, Xiongwei Han, Liang Feng, Linqi
646 Song, Xiaodan Liang, and Jing Tang. Optibench meets resocratic: Measure and improve LLMs for
647 optimization modeling. In *The Thirteenth International Conference on Learning Representations*,
2025. URL <https://openreview.net/forum?id=fsDZwS49uY>.

648 Yafeng Yin. Multiobjective bilevel optimization for transportation planning and management prob-
649 lems. *Journal of Advanced Transportation*, 36(1):93–105, 2002. doi: [https://doi.org/10.1002/atr.](https://doi.org/10.1002/atr.5670360106)
650 [5670360106](https://doi.org/10.1002/atr.5670360106). URL [https://onlinelibrary.wiley.com/doi/abs/10.1002/atr.](https://onlinelibrary.wiley.com/doi/abs/10.1002/atr.5670360106)
651 [5670360106](https://onlinelibrary.wiley.com/doi/abs/10.1002/atr.5670360106).
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

USE OF LLMs

Large language models (LLMs) were used to aid writing polish, including refining sentence phrasing, logical flow, and prose clarity, without altering original meanings or technical details. We use LLM to generate the robot logos in Figure 2, 3, 4, and 5. LLMs did not participate in core research tasks (e.g., experiment design, data processing, model training, result analysis, or drafting key technical content).

A PROOF OF PROPOSITIONS

A.1 PROOF OF PROPOSITION 4.1

Proof. We first prove the inequality $I(D, M) \geq I(S, M)$. Using the chain rule, we obtain $I(D, M) = I(S, M) + I(D, M | S) - I(S, M | D)$. We analyze each term in this decomposition within the context of our framework.

We show $I(S, M | D) = 0$. In our framework, the structure S is derived deterministically (or via a high-fidelity extraction) from the description D by the Distillation Agent. Therefore, given D , S is fixed and provides no additional uncertainty reduction for M . Thus, the conditional mutual information is zero. We show $I(D, M | S) \geq 0$. By the non-negativity property of mutual information, this term is non-negative. Substituting these values back into the decomposition yields the fundamental inequality: $I(D, M) \geq I(S, M)$.

Then, we prove the inequality $I(S, M) \geq I(S, \tilde{S})$. Using the chain rule, we obtain $I(S, M) = I(S, \tilde{S}) + I(S, M | \tilde{S}) - I(S, \tilde{S} | M)$. We analyze each term in this decomposition within the context of our framework. We show $I(S, \tilde{S} | M) = 0$. In our framework, the structure \tilde{S} is derived deterministically from the description M . Thus, the conditional mutual information is zero. We show $I(S, M | \tilde{S}) \geq 0$. By the non-negativity property of mutual information, this term is non-negative. Substituting these values back into the decomposition yields the fundamental inequality: $I(S, M) \geq I(S, \tilde{S})$.

□

A.2 PROOF OF PROPOSITION 4.2

Proof. We prove that the mutual information between the problem description and the model is an upper bound for the mutual information between the description and the interpreted solution. Using the chain rule for mutual information, we obtain $I(D, M) = I(D, \tilde{D}) + I(D, M | \tilde{D}) - I(D, \tilde{D} | M)$. We analyze each term in this decomposition within the context of our framework:

$I(D, \tilde{D} | M) \geq 0$ (**Information Leakage**). In our implementation, the Solution Interpretation Agent uses the original description D to ensure the interpreted solution \tilde{D} uses the correct terminology (e.g., mapping variable x_{ij} back to "Reservoir i to j "). Thus, there is some information flow from D to \tilde{D} that is not mediated by the symbolic model M . By the non-negativity property, this term is non-negative.

$I(D, M | \tilde{D}) \gg 0$ (**Information Loss**). This term represents the information M contains about D that is lost when summarizing the solution into \tilde{D} . The model M represents the entire set of logical constraints and relationships defined in D . The interpreted solution \tilde{D} describes only one specific solution derived from M . Since the general model logic M contains vastly more information about the problem definition D than a single solution instance \tilde{D} , this Information Loss term is strictly positive and significantly dominates the information leakage term ($I(D, M | \tilde{D}) \gg I(D, \tilde{D} | M)$).

Substituting these values back into the decomposition yields the fundamental inequality $I(D, M) \geq I(D, \tilde{D})$. □

756 **A.3 PROOF OF PROPOSITION 4.3**

757
 758 *Proof.* For any non-negative random variable $Z \triangleq I(S, \tilde{S})$ and any threshold $\varepsilon > 0$, the pointwise
 759 inequality $Z \geq \varepsilon \mathbf{1}_{\{Z \geq \varepsilon\}}$ holds. Substituting our terms yields:

760
 761
$$I(S, \tilde{S}) \geq \varepsilon \mathbf{1}_{\{I(S, \tilde{S}) \geq \varepsilon\}} = \varepsilon c_c \quad \text{a.s.}$$

762
 763
 764 Taking expectations on both sides gives $\mathbb{E}[I(S, \tilde{S})] \geq \varepsilon \mathbb{E}[c_c]$. □

765
 766 **A.4 PROOF OF PROPOSITION 4.4**

767
 768 *Proof.* Let $Z' \triangleq I(D, \tilde{D})$. Applying the inequality $Z' \geq \varepsilon' \mathbf{1}_{\{Z' \geq \varepsilon'\}}$ and taking expectations yields
 769 the result. □

770
 771
 772
 773 **B MORE EXPERIMENT RESULTS**

774
 775 **B.1 THE PROBLEMS WE TRY TO ADDRESS IS CRITICAL IN OPTIMIZATION MODELING**

776
 777 The OptiVer framework is designed
 778 and validated for broad applicability
 779 across a wide range of problem do-
 780 mains and model types. Our exper-
 781 imental results provide compelling
 782 evidence for this generalizability.

778
 779
 780
 781
 782
 783

	NL4Opt	Mamo ComplexLP
Missing Constraints	37.3	20.0
Failure Model Debugging	51.8	40.0

784 We demonstrate that the motivations and challenges we address are common and critical in the
 785 optimization modeling field.

786
 787 **The missing constraints** Section 4.5 of the OptiMUS paper (Ahmaditeshnizi et al., 2024) has
 788 summarized and classified common errors, including missing or wrong constraints, incorrect model,
 789 and coding errors. Missing or wrong constraints mean the model fails to extract all the constraints
 790 from the model or generates wrong constraints. An incorrect model means errors, such as defining
 791 binary variables for visiting cities instead of links in TSP. The prevalence of such mistakes is notable
 792 in benchmarks, with 36.0% in NL4Opt and 12.6% in ComplexOR.

793
 794 **Incorrect model debugging** This is also a common challenge. The OptiMUS paper (Ahmaditesh-
 795 nizi et al., 2024) points out that "Coding errors are easier to identify and fix. In contrast, identifying
 796 bugs in the formulation requires deeper reasoning and is harder." In existing methods, the debugging
 797 module is called only when the solver codes raise execution errors.

798
 799 **C EXPERIMENT SETTING IN SECTION 5.3**

800
 801 The plug-and-play capability of OptiVer is supported by both its architectural design and empirical
 802 results across diverse setups. We have built OptiVer based on two modeling baselines (OptiMUS and
 803 ORLM) and two pretrained backbone LLMs (GPT-4o and Qwen2.5-14B). OptiMUS are prompt-
 804 based methods with general LLMs as backbones (we use GPT4o-mini here), which can process
 805 any text inputs. We first extract multi-level structures for the problems using a structure distillation
 806 agent. These extracted structures are then appended to the problem descriptions and sent as input
 807 to the OptiMUS. Once the baselines generate an initial formulation, we proceed with OptiVer’s
 808 verification step. However, ORLM is a fine-tuned model designed to handle only specific input
 809 formats. Therefore, the ORLM model is used solely to provide an initial optimization model, while
 we perform the verification and refinement processes using the GPT4o-mini model.

D CASE STUDY

For the example to explain why OptiVer can mitigate the errors, we provide the following optimization problem with output of CoT and OptiVer.

Problem (simplified version) In a vast network of canals interconnecting 9 strategic reservoirs, engineers have meticulously mapped out the maximum capacities for water transfer from reservoirs i to j , i.e., c_{ij} . The challenge is to optimize this water distribution system to ascertain the maximum possible water transfer from Reservoir 0 to Reservoir 8 per day.

Optimization Formulation

Optimization Model Given by CoT:

- **Decision Variables:** x_{ij} : The amount of water transferred from Reservoir i to Reservoir j , where $i, j = 0, 1, 2, \dots, 8$.
- **Objective Function:** Maximize the water transfer from Reservoir 0 to 8:

$$Z = \sum_{j=1}^8 x_{0j}$$

- **Constraints:**
 1. Water transfer capacities: $x_{ij} \leq c_{ij}$, for all i, j .
 2. Non-negativity: $x_{ij} \geq 0$, for all i, j .

This model is incorrect due to missing flow balance constraints. The verification process is outlined as follows:

- **Structure-Augmented Modeling:** The model references a maximum flow problem. It correctly formulates the flow balance constraint when recalling the standard model.
- **Structure-Side Verification:** The model interprets the current optimization model and compares it with the structure of the original problems.
- **Solution-Side Verification:** If the model lacks flow balance constraints, the obtained solution is represented as $x_{ij} = c_{ij}$. The evaluation agent in OptiVer analyzes the solutions and determines that the inflow does not equal the outflow within the system. Consequently, the evaluation agent identifies this discrepancy as an error.

Optimization Formulation

Optimization Model Given by OptiVer:

- **Modeling structures:**
 - **High Level:** Maximum flow problem
 - **Medium Level:** Single commodity maximum flow
 - **Low Level:**
 1. Directed Network: The flow is directed from one reservoir to another.
 2. Capacity Constraints: Each edge has a maximum capacity.
 3. Flow Conservation: The amount of water entering any intermediate reservoir must equal the amount leaving.
- **Decision Variables:** x_{ij} : The amount of water transferred from Reservoir i to Reservoir j .
- **Objective Function:** Maximize the water transfer:

$$Z = \sum_{j=1}^8 x_{0j}$$

- **Constraints:**
 1. $x_{ij} \leq c_{ij}$, for all i, j .
 2. $x_{ij} \geq 0$, for all i, j .
 3. Flow Conservation: $\sum_{j \neq k}^8 x_{kj} = \sum_{i \neq k}^8 x_{ik}$ for k in the reservoirs

Analysis The modeling structures are proposed to address the challenges of missing constraints. The core of structure-augmented modeling is to identify a similar standard optimization model, and identify the implicit constraints using the standard optimization model as a reference.

E ERROR ANALYSIS ON DIFFERENT PROBLEM TYPES

The results presented in Table 9 clearly demonstrate OptiVer’s strong generalization capabilities, as it consistently and significantly outperforms the CoT baseline across five distinct and challenging problem categories. This robust performance is particularly evident in problem types where standard prompting methods struggle. For instance, on the Capacitated TSP, where CoT achieves a mere 5.13% accuracy, OptiVer boosts performance to 48.72%. Similarly, for Diet, Transportation, and Maximum Flow problems, OptiVer elevates accuracy from the 16-27% range to a much more effective 55-82% range. This shows that OptiVer’s verification process can successfully navigate complex problem structures that are difficult for LLMs to model correctly. Furthermore, even in cases where the CoT baseline is already strong in some problems, such as the Facility Location-Allocation Problem (80.65%), OptiVer still provides a significant improvement, pushing the accuracy to 93.55%. The consistent and substantial performance lift across this diverse set of problems underscores that OptiVer’s adaptive verification framework is a broadly applicable and effective strategy, rather than a technique tailored to a specific problem type.

Table 9: The performance on each problem category on the MAMO ComplexLP dataset

Problem Category	CoT	OptiVer
Diet Problem	27.27%	81.82%
Transportation Problem	23.53%	70.59%
Capacitated TSP	5.13%	48.72%
Maximum Flow Problem	16.28%	55.81%
Facility Location-Allocation Problem	80.65%	93.55%

918 F THE PROMPT DESIGN

919 F.1 STRUCTURE DISTILLATION

```

922 interpretation_prompt=[
923 """
924 You are a mathematical formulator working with a team of optimization
925     ↪ experts. The objective is to tackle a complex optimization problem.
926 """',
927 """
928 Please interpret and explain the following problem description.
929
930 {problem}
931
932 - What is the specific problem type of this OR and CO problem? What
933     ↪ specific kind of OR problem?
934 """',
935 """
936 This is the base formulation of the problem
937
938 {base_formulation}
939
940 - What is the subdivision of different kinds of this problem?
941 - Is this base formulation correct?
942 """',
943 """
944 - Is there any implicit constraints in the problem, including but not
945     ↪ limited to the logical selection relation, if/else and if/then
946     ↪ relation?
947 """',
948 """
949 Please summarize and write in JSON Format. For 'subdivision', please find
950     ↪ the ones matching this problem description
951
952 ```json
953 {{
954   "problem_type": ...,
955   "specific_type": ...,
956   "subdivisions": {{
957     subdivision 1: description,
958     subdivision 2: description,
959     ...
960   }},
961   "implicit_constraints": {{
962     implicit constraint 1: description,
963     implicit constraint 2: description,
964     ...
965   }},
966 }}
967 ```
968
969 - Note that I'm going to use python json.loads() function to parse the
970     ↪ json file, so please make sure the format is correct (don't add ', '
971     ↪ before enclosing '}' or ']' characters.
972 - Generate the complete json file and don't omit anything.
973 - Use '```json' and '```' to enclose the json file.
974 """
975 ]

```

F.2 STRUCTURE-AUGMENTED MODELING

```

972
973
974 formulation_prompt = [
975     """
976     You are an expert mathematical formulator and an optimization professor
977         ↪ at a top university. Your task is to model the problem in the
978         ↪ standard LP or MILP form.
979     """,
980     """
981     Here is the description of the problem to be formulated.
982
983     {problem}
984
985     - Please summarize the parameters and their tensor sizes.
986     - Please explain the definition of the parameters.
987     - Please keep the answer brief and concise.
988     """,
989     """
990     please write in JSON Format. Make sure the bracket is closed, especially
991         ↪ when processing the matrices. Do not transpose the matrices and
992         ↪ keep the shape of the matrices.
993
994     {{
995         "parameters": [
996             {
997                 "symbol": "mathematical symbol of the parameters",
998                 "definition": "definition of the parameters",
999                 "value": the value of the parameters,
1000                 "shape": [],
1001             },
1002             {
1003                 "symbol": "mathematical symbol of the parameters",
1004                 "definition": "definition of the parameters",
1005                 "value": the value of the parameters,
1006                 "shape": [],
1007             },
1008             ...
1009         ],
1010     }}
1011
1012     - Use CamelCase and full words for new variable symbols, and do not
1013         ↪ include indices in the symbol (e.g. ItemsSold instead of itemsSold
1014         ↪ or items_sold or ItemsSold_i)
1015     - Note that I'm going to use python json.loads() function to parse the
1016         ↪ json file, so please make sure the format is correct (don't add ',,'
1017         ↪ before enclosing '}}' or ']' characters.
1018     - Use '```json' and '```' to enclose the json file.
1019     """,
1020     """
1021     Here are some of the cases when we need auxiliary variables. Do we need
1022         ↪ to include auxiliary binary variables in the formulation?
1023
1024     - Logical Conditions: When a decision depends on a binary condition (e.g.,
1025         ↪ whether to open a facility or not, use a kind of transportation or
1026         ↪ not ,and so on), auxiliary binary variables can represent these
1027         ↪ conditions.
1028     - Modeling step costs: Using binary variables involves creating a
1029         ↪ mathematical formulation where costs change based on specific
1030         ↪ thresholds or levels of activity.
1031     - Disjunctive Constraints: When a problem involves "either-or" situations,
1032         ↪ binary variables can be used to model these disjunctions
1033         ↪ effectively (Combined with the big M method).
1034     - Capacity Constraints: In problems involving limited resources, binary
1035         ↪ variables can indicate whether a resource is being utilized or not,
1036         ↪ allowing for better modeling of capacity.

```

```

1026 - Selection Problems: In scenarios where a fixed set of items or
1027   ↪ variables can be selected (e.g., choosing a subset of projects to
1028   ↪ fund), binary variables indicate the selection status.
1029 - Scheduling Order: When determining the sequence in which tasks are
1030   ↪ performed, binary variables can indicate the order of tasks (e.g.,
1031   ↪ Task A before Task B). This is often used in job-shop scheduling or
1032   ↪ project scheduling.
1033 - Penalty Costs: In scheduling with penalties for delays (like tardiness
1034   ↪ or unmet deadlines), binary variables can help track whether a task
1035   ↪ incurs a penalty, allowing for cost minimization.
1036 - Job Switching: In scenarios where workers or machines can switch
1037   ↪ between tasks, binary variables can indicate if a switch occurs,
1038   ↪ helping to manage transition times and costs.
1039 """
1040 """
1041 This problem is a {problem_type} problem with structures
1042 {structure}
1043 To analyze the description carefully, here is the base formulation of
1044   ↪ this problem (which can be correct or needs to be modified)
1045 {base_formulation}
1046 Now take a deep breath and formulate this problem according to the
1047   ↪ description and base formulation.
1048
1049 - Consider whether we need to introduce auxiliary binary variables, note
1050   ↪ that do not include redundant variables.
1051 - For variables, use integer type for discrete items (such as production,
1052   ↪ unit, people) and continuous ones for continuous items (water,
1053   ↪ land, time, grams, and so on).
1054 - Your formulation should be in LaTeX mathematical format (do not include
1055   ↪ the $ symbols).
1056 - Important: You can not define new parameters. You can only define new
1057   ↪ variables. Use CamelCase and full words for new variable symbols,
1058   ↪ and do not indices in the symbol (e.g. ItemsSold instead of
1059   ↪ itemsSold or items_sold or ItemsSold_i). You can include indices in
1060   ↪ the constraint and objective formulations.
1061 - Make sure that you do not use the numeric number in the formulation
1062   ↪ except when necessary, instead, you use the parameter name (you can
1063   ↪ include indices in the constraint and objective formulations).
1064 - Always use non-strict inequalities (e.g. \\leq instead of <), even if
1065   ↪ the constraint is strict.
1066
1067 Take a deep breath and solve the problem step by step.
1068 """
1069 ]
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

```

1080 F.3 STRUCTURE INTERPRETATION AND STRUCTURE CONSISTENCY VERIFICATION
 1081

```

1082 modification_prompt = [
1083     """
1084     You are an expert mathematical formulator and an optimization professor
1085     ↪ at a top university. Your task is to model and fix the problem in
1086     ↪ the standard LP or MILP form.
1087     """,
1088     """
1089     This is a {problem_type} problem with parameters
1090     {parameters}
1091     The formulation is as follows
1092     {formulation_interpretation}
1093     Does this problem consistent with the characteristics of the following
1094     ↪ structure description? If yes, please say "Yes" directly.
1095     If not, please give your comments to modify the formulation.
1096     {original_problem_interpretation}
1097     """,
1098     """
1099     Please reformulate the problem to make the formulation consistent with
1100     ↪ the structure description.
1101     - Consider whether we need to introduce extra binary variables or
1102     ↪ linearization for a piece-wise linear function.
1103     - Your formulation should be in LaTeX mathematical format (do not include
1104     ↪ the $ symbols).
1105     - Important: You can not define new parameters. You can only define new
1106     ↪ variables. Use CamelCase and full words for new variable symbols,
1107     ↪ and do not include indices in the symbol (e.g. ItemsSold instead of
1108     ↪ itemsSold or items_sold or ItemsSold_i). You can include indices
1109     ↪ in the constraint and objective formulations.
1110     - Make sure that you do not use a numeric number in the formulation
1111     ↪ except where necessary; instead, you use the parameter name (you
1112     ↪ can include indices in the constraint and objective formulations).
1113     - Always use non-strict inequalities (e.g. \\leq instead of <), even if
1114     ↪ the constraint is strict.
1115     Take a deep breath and solve the problem step by step.
1116     """
1117 ]
  
```

1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

F.4 SOLUTION INTERPRETATION AND SOLUTION VALIDITY VERIFICATION

```

1134
1135
1136 solution_prompt = [
1137     """
1138     You are an expert mathematical formulator and an optimization professor
1139         ↪ at a top university. Your task is to model and fix the problem
1140         ↪ using the solution information in the standard LP or MILP form.
1141     """,
1142     """
1143     This is a {problem_type} problem with solutions
1144
1145     {solutions}
1146
1147     The formulation is as follows
1148
1149     {formulation_interpretation}
1150
1151     Please interpret the meaning of the solution.
1152     """,
1153     """
1154     Here is the problem description.
1155
1156     {original_problem_interpretation}
1157
1158     Is this solution the optimal solution? The optimal solution should be
1159         ↪ mathematical sound and logical coherence:
1160     - We cannot find a better solution.
1161     - The solution should meet the constraints of the problem description.
1162
1163     If yes, please say "Yes" directly.
1164     If not, please give your comments to modify the formulation.
1165     """,
1166     """
1167     Please reformulate the problem to make the formulation consistent with
1168         ↪ the structure description.
1169
1170     - Consider whether we need to introduce extra binary variables or
1171         ↪ linearization for a piece-wise linear function.
1172     - Your formulation should be in LaTeX mathematical format (do not include
1173         ↪ the $ symbols).
1174     - Important: You can not define new parameters. You can only define new
1175         ↪ variables. Use CamelCase and full words for new variable symbols,
1176         ↪ and do not include indices in the symbol (e.g. ItemsSold instead of
1177         ↪ itemsSold or items_sold or ItemsSold_i). You can include indices
1178         ↪ in the constraint and objective formulations.
1179     - Make sure that you do not use a numeric number in the formulation
1180         ↪ except where necessary; instead, you use the parameter name (you
1181         ↪ can include indices in the constraint and objective formulations).
1182     - Always use non-strict inequalities (e.g. \\leq instead of <), even if
1183         ↪ the constraint is strict.
1184
1185     Take a deep breath and solve the problem step by step.
1186     """,
1187     ]

```

G MORE EXAMPLES

G.1 NL4OPT

Example in NL4Opt

Natural Language Description:

A fishery wants to transport their catch. They can either use local sled dogs or trucks. Local sled dogs can take 100 fish per trip while trucks can take 300 fish per trip. The cost per trip for sled dogs is \$50 while the cost per trip for a truck is \$100. The budget is at most \$1000 and the number of sled dog trips must be less than the number of truck trips. Formulate an LP to maximize the number of fish that can be transported.

Modeling Structure Given by OptiVer:

```

1 {
2   "problem_type": "Linear Programming",
3   "specific_type": "Resource Allocation Problem",
4   "subdivisions": {
5     "transportation optimization": "Optimizing the transport
6       ↪ of goods or resources using various modes of
7       ↪ transport.",
8     "resource allocation": "Distributing limited resources to
9       ↪ achieve the best possible outcome.",
10    "cost minimization and profit maximization": "Balancing
11      ↪ costs against outputs to maximize profit or benefits
12      ↪ ."
13  },
14  "implicit_constraints": {
15    "non-negative trips": "The number of trips (x and y)
16      ↪ cannot be negative.",
17    "integer trips": "The number of trips should generally be
18      ↪ an integer value."
19  }
20 }

```

Optimization Model Given by OptiVer:

```

1 {
2   "constraints": [
3     {
4       "name": "Budget Constraint",
5       "description": "The total cost of trips must not
6         ↪ exceed the budget.",
7       "formulation": "CostPerSledDogTrip \\.cdot
8         ↪ NumberOfSledDogTrips + CostPerTruckTrip \\.cdot
9         ↪ NumberOfTruckTrips \\.leq TotalBudget"
10    },
11    {
12      "name": "Trip Comparison Constraint",
13      "description": "The number of sled dog trips must
14        ↪ not exceed the number of truck trips.",
15      "formulation": "NumberOfSledDogTrips \\.leq
16        ↪ NumberOfTruckTrips"
17    }
18  ],
19 }

```

```

1242
1243
1244 13     {
1245 14         "name": "Non-negativity Constraints",
1246 15         "description": "The number of trips must be non-
1247 16             ↪ negative.",
1248 17         "formulation": "NumberOfSledDogTrips \\geq 0"
1249 18     },
1250 19     {
1251 20         "name": "Non-negativity Constraints",
1252 21         "description": "The number of trips must be non-
1253 22             ↪ negative.",
1254 23         "formulation": "NumberOfTruckTrips \\geq 0"
1255 24     },
1256 25     {
1257 26         "name": "Integer Constraints for Sled Dog Trips",
1258 27         "description": "The number of sled dog trips must be
1259 28             ↪ an integer.",
1260 29         "formulation": "NumberOfSledDogTrips \\in \\mathbb{Z}
1261 30             ↪ )"
1262 31     },
1263 32     {
1264 33         "name": "Integer Constraints for Truck Trips",
1265 34         "description": "The number of truck trips must be an
1266 35             ↪ integer.",
1267 36         "formulation": "NumberOfTruckTrips \\in \\mathbb{Z}"
1268 37     }
1269 38 ],
1270 39 "variables": [
1271 40     {
1272 41         "definition": "Integer variable representing the
1273 42             ↪ number of sled dog trips.",
1274 43         "symbol": "NumberOfSledDogTrips",
1275 44         "type": "integer",
1276 45         "shape": []
1277 46     },
1278 47     {
1279 48         "definition": "Integer variable representing the
1280 49             ↪ number of truck trips.",
1281 50         "symbol": "NumberOfTruckTrips",
1282 51         "type": "integer",
1283 52         "shape": []
1284 53     }
1285 54 ],
1286 55 "objective": [
1287 56     {
1288 57         "description": "Maximize the total number of fish
1289 58             ↪ transported.",
1290 59         "formulation": "Z = FishPerSledDogTrip \\cdot
1291 60             ↪ NumberOfSledDogTrips + FishPerTruckTrip \\cdot
1292 61             ↪ NumberOfTruckTrips",
1293 62         "objective_sense": "Maximize"
1294 63     }
1295 64 ]

```

Solver Code Given by OptiVer:

```

1296
1297
1298 import json
1299 import numpy as np
1300 import math
1301 import gurobipy as gp
1302 from gurobipy import GRB
1303
1304 # Create a new model
1305 model = gp.Model('model')
1306
1307 # Define parameters
1308 parameters = {
1309     "CostPerSledDogTrip": 50,
1310     "CostPerTruckTrip": 100,
1311     "FishPerSledDogTrip": 100,
1312     "FishPerTruckTrip": 300,
1313     "TotalBudget": 1000
1314 }
1315
1316 # Define variables
1317 NumberOfSledDogTrips = model.addVar(vtype=GRB.INTEGER, name="
1318     ↪ NumberOfSledDogTrips")
1319 NumberOfTruckTrips = model.addVar(vtype=GRB.INTEGER, name="
1320     ↪ NumberOfTruckTrips")
1321
1322 # Add variables to the model
1323 model.update()
1324
1325 # Constraints
1326 # Budget Constraint
1327 model.addConstr(
1328     parameters["CostPerSledDogTrip"] * NumberOfSledDogTrips +
1329     parameters["CostPerTruckTrip"] * NumberOfTruckTrips <=
1330     parameters["TotalBudget"],
1331     name="BudgetConstraint"
1332 )
1333
1334 # Trip Comparison Constraint
1335 model.addConstr(
1336     NumberOfSledDogTrips <= NumberOfTruckTrips,
1337     name="TripComparisonConstraint"
1338 )
1339
1340 # Non-negativity Constraints
1341 model.addConstr(NumberOfSledDogTrips >= 0, name="
1342     ↪ NonNegativitySledDogTrips")
1343 model.addConstr(NumberOfTruckTrips >= 0, name="
1344     ↪ NonNegativityTruckTrips")
1345
1346 # Set the objective
1347 model.setObjective(
1348     parameters["FishPerSledDogTrip"] * NumberOfSledDogTrips +
1349     parameters["FishPerTruckTrip"] * NumberOfTruckTrips,
1350     GRB.MAXIMIZE
1351 )
1352
1353 # Optimize the model
1354 model.optimize()
1355 status = model.status
1356
1357 obj_val = None
1358 # Check whether the model is infeasible, has infinite solutions, or
1359     ↪ has an optimal solution

```

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

```
if status == gp.GRB.INFEASIBLE:  
    obj_val = "infeasible"  
elif status == gp.GRB.UNBOUNDED:  
    obj_val = "unbounded"  
elif status == gp.GRB.OPTIMAL:  
    obj_val = model.objVal  
  
print("Objective Value:", obj_val)
```

G.2 MAMO COMPLEXLP

Example in Mamo ComplexLP

Natural Language Description:

In a scenario involving a salesperson who needs to visit five different cities to conduct business, each city is uniquely numbered from 1 to 5. The salesperson's objective is to minimize the total travel expenses, which could be influenced by factors such as distance, fuel costs, or transportation fees. The salesperson can start their journey from any of these cities but must ensure they visit each city exactly once before returning to their starting point.

The travel costs between the cities are as follows:

- From City 1, the travel costs are 58 units to City 2, 15 units to City 3, 75 units to City 4, and 91 units to City 5.
- From City 2, it costs 58 units to City 1, 54 units to City 3, 85 units to City 4, and 11 units to City 5.
- Traveling from City 3, the expenses are 15 units to City 1, 54 units to City 2, 28 units to City 4, and 61 units to City 5.
- From City 4, the costs are 75 units to City 1, 85 units to City 2, 28 units to City 3, and 47 units to City 5.
- Lastly, from City 5, it costs 91 units to City 1, 11 units to City 2, 61 units to City 3, and 47 units to City 4.

Given this setup, what is the minimum total travel cost for the salesperson to visit each city exactly once and then return to the starting city?

Modeling Structure Given by OptiVer:

```

1 {
2   "problem_type": "Combinatorial Optimization Problem",
3   "specific_type": "Traveling Salesman Problem (TSP)",
4   "subdivisions": {
5     "Hamiltonian cycle": "Path that visits each vertex exactly
6       ↪ once and returns to the starting vertex",
7     "Weighted graph": "Graph with weights on edges
8       ↪ representing travel costs between cities",
9     "Directed graph": "Graph where edges have a direction,
10      ↪ indicating the cost of travel from one city to
11      ↪ another"
12   },
13   "implicit_constraints": {
14     "subtour elimination": "Explicit constraints to prevent
15       ↪ subtours in the solution",
16     "start_end city constraint": "Salesperson must start and
17       ↪ end at the same city"
18   }
19 }

```

Optimization Model Given by OptiVer:

```

1 {
2   "constraints": [
3     {
4       "name": "Each City Visited Once",

```

```

1458
1459 5      "description": "Each city must be visited exactly
1460      ↪ once by the salesperson.",
1461 6      "formulation": "\\sum_{j \\in Cities} x_{ij} = 1\\
1462      ↪ quad \\forall i \\in Cities"
1463 7      },
1464 8      {
1465 9      "name": "Return to Start City",
1466 10     "description": "The salesperson must return to the
1467     ↪ starting city after visiting all cities.",
1468 11     "formulation": "\\sum_{i \\in Cities} x_{ji} = 1\\
1469     ↪ quad \\forall j \\in Cities"
1470 12     },
1471 13     {
1472 14     "name": "Subtour Elimination",
1473 15     "description": "Constraints to prevent subtours in
1474     ↪ the solution.",
1475 16     "formulation": "u_i - u_j + (|Cities|) \\cdot x_{ij}
1476     ↪ \\leq |Cities| - 1 \\quad \\forall i, j \\in
1477     ↪ Cities, \\; i \\neq j"
1478 17     },
1479 18     {
1480 19     "name": "Position Constraints",
1481 20     "description": "Position variables must be within
1482     ↪ valid range.",
1483 21     "formulation": "2 \\leq u_i \\leq |Cities| \\quad \\
1484     ↪ forall i \\in Cities"
1485 22     }
1486 23 ],
1487 24 "variables": [
1488 25     {
1489 26     "definition": "Binary variable indicating whether
1490     ↪ the salesperson travels from city i to city j
1491     ↪ .",
1492 27     "symbol": "x_{ij}",
1493 28     "type": "binary",
1494 29     "shape": [
1495 30     5,
1496 31     5
1497 32     ]
1498 33     },
1499 34     {
1500 35     "definition": "Auxiliary continuous variable
1501     ↪ representing the position of city i in the
1502     ↪ tour.",
1503 36     "symbol": "u_i",
1504 37     "type": "continuous",
1505 38     "shape": [
1506 39     5
1507 40     ]
1508 41     }
1509 42 ],
1510 43 "objective": [
1511 44     {
1512 45     "description": "Minimize the total travel cost for
1513     ↪ the salesperson.",

```

```

1512
1513 46         "formulation": "Z = \\sum_{i \\in Cities} \\sum_{j
1514           ↪ \\in Cities} d_{ij} \\cdot x_{ij}",
1515 47         "objective_sense": "Minimize"
1516 48     }
1517 49 ]
1518 50 }

```

Solver Code Given by OptiVer:

```

1523 import json
1524 import numpy as np
1525 import math
1526 import gurobipy as gp
1527 from gurobipy import GRB
1528
1529 # Create a new model
1530 model = gp.Model('model')
1531
1532 # Define parameters
1533 parameters = {
1534     "Cities": [1, 2, 3, 4, 5],
1535     "TravelCosts": [
1536         [0, 58, 15, 75, 91],
1537         [58, 0, 54, 85, 11],
1538         [15, 54, 0, 28, 61],
1539         [75, 85, 28, 0, 47],
1540         [91, 11, 61, 47, 0]
1541     ]
1542 }
1543
1544 Cities = parameters["Cities"]
1545 TravelCosts = parameters["TravelCosts"]
1546 num_cities = len(Cities) # Number of cities
1547
1548 # Define variables
1549 x = model.addVars(num_cities, num_cities, vtype=GRB.BINARY, name="x
1550 ↪ ") # Binary variables for routes
1551 u = model.addVars(num_cities, vtype=GRB.CONTINUOUS, name="u") #
1552 ↪ Continuous variables for subtour elimination
1553
1554 # Constraints: Each City Visited Once
1555 for i in range(num_cities):
1556     model.addConstr(gp.quicksum(x[i, j] for j in range(num_cities)
1557 ↪ if i != j) == 1, f"VisitOnce_Cities[{i}]")
1558
1559 # Constraints: Return to Start City
1560 for j in range(num_cities):
1561     model.addConstr(gp.quicksum(x[i, j] for i in range(num_cities)
1562 ↪ if i != j) == 1, f"ReturnToStart_Cities[{j}]")
1563
1564 # Subtour Elimination Constraints
1565 for i in range(1, num_cities):
1566     for j in range(1, num_cities):
1567         if i != j:
1568             model.addConstr(u[i] - u[j] + num_cities * x[i, j] <=
1569 ↪ num_cities - 1, f"SubtourElimination_Cities[{i}][{j}
1570 ↪ ]")
1571
1572 # Position Constraints

```

```
1566
1567
1568     for i in range(num_cities):
1569         model.addConstr(u[i] >= 2, f"LowerBound_u[{i}]")
1570         model.addConstr(u[i] <= num_cities, f"UpperBound_u[{i}]")
1571
1572     # Objective: Minimize total travel cost
1573     model.setObjective(gp.quicksum(TravelCosts[i][j] * x[i, j] for i in
1574         ↪ range(num_cities) for j in range(num_cities)), GRB.MINIMIZE)
1575
1576     # Optimize the model
1577     model.optimize()
1578     status = model.status
1579
1580     obj_val = None
1581     # Check whether the model is infeasible, has infinite solutions, or
1582     ↪ has an optimal solution
1583     if status == gp.GRB.INFEASIBLE:
1584         obj_val = "infeasible"
1585     elif status == gp.GRB.UNBOUNDED:
1586         obj_val = "unbounded"
1587     elif status == gp.GRB.OPTIMAL:
1588         obj_val = model.objVal
1589
1590     print("Objective Value:", obj_val)
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
```

G.3 COMPLEXOR

Example in ComplexOR

Natural Language Description:

The capacitated warehouse location problem involves determining the optimal locations for a set number of warehouses to service customers at minimum cost, taking into account warehouse capacities, operating costs, and customer demand.

The capacitated warehouse location problem is the problem of locating `NumberOfLocations` warehouses which have to service `NumberOfCustomers` customers, at minimum cost. Each customer has an associated demand `CustomerDemand`. There are constraints on the total demand that can be met from a warehouse, as specified by `WarehouseCapacity`. Costs are incurred when allocating service to customers from warehouses `ServiceAllocationCost`, and warehouses have a fixed operating cost `WarehouseFixedCost`. Additionally, there is a lower limit `MinimumDemandFromWarehouse` on the amount of demand that a warehouse must meet if it is opened, as well as constraints on the minimum `MinimumOpenWarehouses` and maximum `MaximumOpenWarehouses` number of warehouses that can be operational.

The total number of potential warehouse locations is 10. The total number of customers to be serviced is 20. The demand of each customer is [117, 86, 69, 53, 110, 74, 136, 140, 126, 79, 54, 86, 114, 76, 136, 73, 144, 51, 53, 120]. The cost of allocating service from each warehouse to each customer is [[80, 94, 44, 51, 190, 44, 129, 178, 129, 91, 172, 119, 177, 150, 90, 51, 53, 97, 184, 87], [139, 33, 104, 135, 50, 176, 97, 121, 47, 29, 186, 163, 149, 108, 156, 169, 100, 160, 153, 85], [153, 36, 18, 170, 18, 181, 178, 68, 171, 106, 159, 110, 21, 106, 91, 29, 144, 140, 155, 116], [103, 59, 78, 125, 14, 11, 152, 95, 76, 173, 36, 148, 75, 132, 59, 153, 113, 74, 185, 71], [193, 186, 130, 145, 114, 150, 33, 154, 20, 75, 103, 30, 137, 131, 167, 32, 53, 150, 176, 166], [159, 130, 156, 65, 36, 59, 199, 124, 104, 72, 180, 73, 43, 152, 143, 90, 161, 65, 172, 141], [173, 121, 110, 127, 22, 159, 195, 137, 47, 10, 87, 11, 154, 66, 126, 60, 152, 54, 20, 25], [181, 34, 186, 152, 109, 195, 133, 198, 30, 65, 69, 19, 109, 143, 108, 196, 59, 133, 10, 123], [82, 113, 147, 21, 88, 24, 38, 16, 70, 122, 148, 192, 116, 108, 18, 20, 143, 18, 116, 142], [176, 170, 87, 91, 195, 183, 124, 89, 72, 97, 89, 23, 45, 196, 97, 27, 83, 81, 171, 148]]. The total capacity for each warehouse is [3010, 2910, 4530, 4720, 4920, 3750, 4930, 2970, 3310, 2460]. The lower limit on the demand that must be met from a warehouse if it is to be operational is [64, 55, 27, 71, 93, 90, 89, 87, 43, 50]. The minimum number of warehouses that need to be operational is 3. The maximum number of warehouses that can be operational is 8. The fixed operating cost of each warehouse is [8517, 5068, 9433, 6127, 6033, 5966, 7762, 9406, 6602, 7040].

Modeling Structure Given by OptiVer:

```

1 {
2   "problem_type": "Mixed Integer Linear Programming",
3   "specific_type": "Capacitated Warehouse Location Problem",
4   "subdivisions": {
5     "1": "Location Optimization",
6     "2": "Capacity Planning",
7     "3": "Cost Minimization"
8   },
9   "implicit_constraints": {
10    "1": "Each customer's demand must be fully met.",
11    "2": "A warehouse that is opened must meet a specified
12         ↪ minimum demand."
13  }

```

Optimization Model Given by OptiVer:

```

1674
1675
1676 1 {
1677 2   "constraints": [
1678 3     {
1679 4       "name": "Demand Meeting",
1680 5       "description": "Each customer's demand must be fully
1681 6         ↪ met.",
1682 7       "formulation": "\\sum_{j=1}^{\text{NumberOfLocations}} y_{
1683 8         ↪ ij} = \text{CustomerDemand}[i] \\quad \\forall i"
1684 9     },
1685 10    {
1686 11     "name": "Capacity Limit",
1687 12     "description": "The total demand served from each
1688 13       ↪ warehouse cannot exceed its capacity.",
1689 14     "formulation": "\\sum_{i=1}^{\text{NumberOfCustomers}} y_{
1690 15       ↪ ij} \\leq \text{WarehouseCapacity}[j] \\cdot x_j \\
1691 16       ↪ quad \\forall j"
1692 17    },
1693 18    {
1694 19     "name": "Minimum Demand",
1695 20     "description": "A warehouse that is opened must meet
1696 21       ↪ a specified minimum demand.",
1697 22     "formulation": "y_{ij} \\geq
1698 23       ↪ \text{MinimumDemandFromWarehouse}[j] \\cdot x_j \\
1699 24       ↪ quad \\forall j, \\forall i"
1700 25    },
1701 26    {
1702 27     "name": "Minimum Open Warehouses",
1703 28     "description": "At least a minimum number of
1704 29       ↪ warehouses must be operational.",
1705 30     "formulation": "\\sum_{j=1}^{\text{NumberOfLocations}} x_j
1706 31       ↪ \\geq \text{MinimumOpenWarehouses}"
1707 32    },
1708 33    {
1709 34     "name": "Maximum Open Warehouses",
1710 35     "description": "At most a maximum number of
1711 36       ↪ warehouses can be operational.",
1712 37     "formulation": "\\sum_{j=1}^{\text{NumberOfLocations}} x_j
1713 38       ↪ \\leq \text{MaximumOpenWarehouses}"
1714 39    }
1715 40  ],
1716 41  "variables": [
1717 42    {
1718 43     "definition": "Binary variable indicating whether
1719 44       ↪ warehouse j is opened (1) or not (0)",
1720 45     "symbol": "x_j",
1721 46     "type": "binary",
1722 47     "shape": [
1723 48       10
1724 49     ]
1725 50    },
1726 51    {
1727 52     "definition": "Continuous variable representing the
1728 53       ↪ amount of demand served from warehouse j to
1729 54       ↪ customer i",
1730 55     "symbol": "y_{ij}",
1731 56     "type": "continuous",

```

```

1728
1729
1730     "shape": [
1731         20,
1732         10
1733     ]
1734 ],
1735 "objective": [
1736     {
1737         "description": "Minimize the total cost of servicing
1738             ↪ customers, including service allocation and
1739             ↪ operating costs of warehouses.",
1740         "formulation": "Z = \\sum_{j=1}^{NumberOfLocations}
1741             ↪ WarehouseFixedCost[j] \\cdot x_j + \\sum_{i=1}^{
1742             ↪ NumberOfCustomers} \\sum_{j=1}^{
1743             ↪ NumberOfLocations} ServiceAllocationCost[j][i]
1744             ↪ \\cdot y_{ij}",
1745         "objective_sense": "Minimize"
1746     }
1747 ]
1748

```

Solver Code Given by OptiVer:

```

1752
1753 import json
1754 import numpy as np
1755 import math
1756 import gurobipy as gp
1757 from gurobipy import quicksum # Import quicksum
1758
1759 # Create a Gurobi model
1760 model = gp.Model('model')
1761
1762 # Parameters
1763 NumberOfLocations = 10
1764 NumberOfCustomers = 20
1765
1766 CustomerDemand = [117, 86, 69, 53, 110, 74, 136, 140, 126, 79,
1767                   54, 86, 114, 76, 136, 73, 144, 51, 53, 120]
1768
1769 ServiceAllocationCost = [
1770     [80, 94, 44, 51, 190, 44, 129, 178, 129, 91, 172, 119, 177, 150,
1771      ↪ 90, 51, 53, 97, 184, 87],
1772     [139, 33, 104, 135, 50, 176, 97, 121, 47, 29, 186, 163, 149, 108,
1773      ↪ 156, 169, 100, 160, 153, 85],
1774     [153, 36, 18, 170, 18, 181, 178, 68, 171, 106, 159, 110, 21, 106,
1775      ↪ 91, 29, 144, 140, 155, 116],
1776     [103, 59, 78, 125, 14, 11, 152, 95, 76, 173, 36, 148, 75, 132,
1777      ↪ 59, 153, 113, 74, 185, 71],
1778     [193, 186, 130, 145, 114, 150, 33, 154, 20, 75, 103, 30, 137,
1779      ↪ 131, 167, 32, 53, 150, 176, 166],
1780     [159, 130, 156, 65, 36, 59, 199, 124, 104, 72, 180, 73, 43, 152,
1781      ↪ 143, 90, 161, 65, 172, 141],
1782     [173, 121, 110, 127, 22, 159, 195, 137, 47, 10, 87, 11, 154, 66,
1783      ↪ 126, 60, 152, 54, 20, 25],
1784     [181, 34, 186, 152, 109, 195, 133, 198, 30, 65, 69, 19, 109, 143,
1785      ↪ 108, 196, 59, 133, 10, 123],

```

```

1782
1783     [82, 113, 147, 21, 88, 24, 38, 16, 70, 122, 148, 192, 116, 108,
1784         ↪ 18, 20, 143, 18, 116, 142],
1785     [176, 170, 87, 91, 195, 183, 124, 89, 72, 97, 89, 23, 45, 196,
1786         ↪ 97, 27, 83, 81, 171, 148]
1787 ]
1788 WarehouseCapacity = [3010, 2910, 4530, 4720, 4920, 3750, 4930, 2970,
1789     ↪ 3310, 2460]
1790
1791 WarehouseFixedCost = [8517, 5068, 9433, 6127, 6033, 5966, 7762,
1792     ↪ 9406, 6602, 7040]
1793
1794 MinimumDemandFromWarehouse = [64, 55, 27, 71, 93, 90, 89, 87, 43,
1795     ↪ 50]
1796
1797 MinimumOpenWarehouses = 3
1798 MaximumOpenWarehouses = 8
1799
1800 # Variables
1801 x = model.addVars(NumberOfLocations, vtype=gp.GRB.BINARY, name="x")
1802 y = model.addVars(NumberOfCustomers, NumberOfLocations, vtype=gp.
1803     ↪ GRB.CONTINUOUS, name="y")
1804
1805 # Objective function
1806 model.setObjective(
1807     quicksum(WarehouseFixedCost[j] * x[j] for j in range(
1808         ↪ NumberOfLocations)) +
1809     quicksum(quicksum(ServiceAllocationCost[j][i] * y[i, j] for j in
1810         ↪ range(NumberOfLocations)) for i in range(
1811         ↪ NumberOfCustomers)),
1812     gp.GRB.MINIMIZE
1813 )
1814
1815 # Constraints
1816 for i in range(NumberOfCustomers):
1817     model.addConstr(
1818         quicksum(y[i, j] for j in range(NumberOfLocations)) ==
1819         ↪ CustomerDemand[i],
1820         name=f"demand_meeting_{i}"
1821     )
1822
1823 for j in range(NumberOfLocations):
1824     model.addConstr(
1825         quicksum(y[i, j] for i in range(NumberOfCustomers)) <=
1826         ↪ WarehouseCapacity[j] * x[j],
1827         name=f"capacity_limit_{j}"
1828     )
1829
1830 for j in range(NumberOfLocations):
1831     model.addConstr(
1832         quicksum(y[i, j] for i in range(NumberOfCustomers)) >=
1833         ↪ MinimumDemandFromWarehouse[j] * x[j],
1834         name=f"minimum_demand_{j}"
1835     )
1836
1837 model.addConstr(
1838     quicksum(x[j] for j in range(NumberOfLocations)) >=
1839     ↪ MinimumOpenWarehouses,
1840     name="minimum_open_warehouses"
1841 )
1842
1843 model.addConstr(

```

```
1836
1837     quicksum(x[j] for j in range(NumberOfLocations)) <=
1838         ↪ MaximumOpenWarehouses,
1839     name="maximum_open_warehouses"
1840 )
1841
1842 # Optimize the model
1843 model.optimize()
1844
1845 # Check the optimization status
1846 status = model.status
1847
1848 obj_val = None
1849 if status == gp.GRB.INFEASIBLE:
1850     obj_val = "infeasible"
1851 elif status == gp.GRB.UNBOUNDED:
1852     obj_val = "unbounded"
1853 elif status == gp.GRB.OPTIMAL:
1854     obj_val = model.objVal
1855
1856 print("Objective Value:", obj_val)
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
```

G.4 COMPLEXOR

Example in ComplexOR

Natural Language Description:

The Knapsack Problem involves selecting the most valuable combination of items to fit in a knapsack without exceeding its weight limit. The Knapsack Problem is a classic optimization problem in operations research and computer science. The problem is to determine the most valuable combination of items to include in a knapsack, given a set of TotalItems with different values and weights represented by ItemValues and ItemWeights respectively, and a maximum weight capacity of the knapsack MaxKnapsackWeight. The goal is to maximize the total value of the items in the knapsack, represented by ItemValues, without exceeding its weight capacity MaxKnapsackWeight. The available kinds of items is 6. The value of each kind of item is [17, 4, 10, 21, 12, 18]. The weight of each item is [23, 6, 14, 30, 15, 25]. The maximum weight capacity of the knapsack is 60. The total weight of the selected items must not exceed MaxKnapsackWeight. Multiple items in one kind can be selected. No more than TotalItems kinds of items can be considered for selection. Maximize the total value of the items in the knapsack.

Modeling Structure Given by OptiVer:

```

1 {
2   "problem_type": "Combinatorial Optimization",
3   "specific_type": "Bounded Knapsack Problem",
4   "subdivisions": {
5     "Subdivision 1": "0/1 Knapsack Problem",
6     "Subdivision 2": "Bounded Knapsack Problem",
7     "Subdivision 3": "Multi-dimensional Knapsack Problem"
8   },
9   "implicit_constraints": {
10    "implicit constraint 1": "x_i must be an integer",
11    "implicit constraint 2": "MaxKnapsackWeight must be non-
    ↪ negative"
12  }
13 }

```

Optimization Model Given by OptiVer:

```

1 {
2   "constraints": [
3     {
4       "name": "Weight Limit Constraint",
5       "description": "The total weight of selected items
    ↪ must not exceed the maximum weight capacity of
    ↪ the knapsack.",
6       "formulation": "\\sum_{i=1}^{\text{TotalItems}}
    ↪ \text{ItemWeights}_i \cdot \text{ItemQuantities}_i \leq
    ↪ \text{MaxKnapsackWeight}"
7     },
8     {
9       "name": "Non-negativity and Integer Constraints",
10      "description": "The quantity of each item selected
    ↪ must be non-negative and integer.",

```

```

1944
1945 11         "formulation": "ItemQuantities_i \\geq 0\\quad \\
1946             ↪ text{and integer for }i = 1,2,\\dots,
1947             ↪ TotalItems"
1948 12     }
1949 13 ],
1950 14     "variables": [
1951 15         {
1952 16             "definition": "Number of items of type i selected (i
1953             ↪ = 1to TotalItems)",
1954 17             "symbol": "ItemQuantities_i",
1955 18             "type": "integer",
1956 19             "shape": []
1957 20         }
1958 21     ],
1959 22     "objective": [
1960 23         {
1961 24             "description": "Maximize the total value of the
1962             ↪ selected items in the knapsack.",
1963 25             "formulation": "Z = \\sum_{i=1}^{TotalItems}
1964             ↪ ItemValues_i \\cdot ItemQuantities_i",
1965 26             "objective_sense": "Maximize"
1966 27         }
1967 28     ]
1968 29 }

```

Solver Code Given by OptiVer:

```

1971
1972 import json
1973 import numpy as np
1974 import math
1975 import gurobipy as gp
1976
1977 # Create a new model
1978 model = gp.Model('model')
1979
1980 # Parameters
1981 TotalItems = 6
1982 ItemValues = [17, 4, 10, 21, 12, 18]
1983 ItemWeights = [23, 6, 14, 30, 15, 25]
1984 MaxKnapsackWeight = 60
1985
1986 # Variables: ItemQuantities_i (integer variables)
1987 ItemQuantities = model.addVars(TotalItems, vtype=gp.GRB.INTEGER,
1988     ↪ name="ItemQuantities")
1989
1990 # Objective: Maximize Z = sum(ItemValues_i * ItemQuantities_i)
1991 model.setObjective(gp.quicksum(ItemValues[i] * ItemQuantities[i]
1992     ↪ for i in range(TotalItems)), gp.GRB.MAXIMIZE)
1993
1994 # Constraints
1995 # Weight Limit Constraint: sum(ItemWeights_i * ItemQuantities_i) <=
1996     ↪ MaxKnapsackWeight
1997 model.addConstr(gp.quicksum(ItemWeights[i] * ItemQuantities[i] for
1998     ↪ i in range(TotalItems)) <= MaxKnapsackWeight, "WeightLimit")
1999
2000 # Non-negativity and Integer Constraints are inherently defined by
2001     ↪ the variable type

```

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

```
# ItemQuantities_i >= 0 and ItemQuantities_i in Z
# Gurobi automatically treats integer variables as non-negative, so
  ↳ no additional constraint is needed for non-negativity.

# Optimize the model
model.optimize()
status = model.status

obj_val = None
# Check whether the model is infeasible, has infinite solutions, or
  ↳ has an optimal solution
if status == gp.GRB.INFEASIBLE:
    obj_val = "infeasible"
elif status == gp.GRB.UNBOUNDED:
    obj_val = "unbounded"
elif status == gp.GRB.OPTIMAL:
    obj_val = model.objVal

print("Objective Value:", obj_val)
```