

# OptiMUS-0.3: Using large language models to model and solve optimization problems at scale

Ali AhmadiTeshnizi<sup>a</sup>, Wenzhi Gao<sup>b</sup>, Herman Brunborg<sup>b</sup>, Shayan Talaei<sup>a</sup>, Connor Lawless<sup>a</sup>, Madeleine Udell<sup>a,b</sup>

---

**Abstract.** Optimization problems are pervasive in sectors from manufacturing and distribution to healthcare. However, most such problems are still solved heuristically by hand rather than optimally by state-of-the-art solvers because the expertise required to formulate and solve these problems limits the widespread adoption of optimization tools and techniques. We introduce a Large Language Model (LLM)-based system designed to formulate and solve (mixed integer) linear programming problems from their natural language descriptions. Our system is capable of developing mathematical models, writing and debugging solver code, evaluating the generated solutions, and improving efficiency and correctness of its model and code based on these evaluations. OptiMUS-0.3 utilizes a modular structure to process problems, allowing it to handle problems with long descriptions and complex data without long prompts. Experiments demonstrate that OptiMUS-0.3 outperforms existing state-of-the-art methods on easy datasets by more than 22% and on hard datasets (including a new dataset, NLP4LP, released with this paper that features long and complex problems) by more than 24%.

**Key words:** Large Language Models, Optimization Modeling, AI Agents, Mixed Integer Linear Programming, Optimization Tools, Human-Computer Interaction

---

## 1. Introduction

Optimization problems are ubiquitous in domains from operations and economics to engineering and computer science. While major advances in optimization algorithms over the last several decades have led to reliable and efficient optimization methods for a wide variety of structured optimization problems, including linear programming (LP) and mixed-integer linear programming (MILP), optimization modeling — transforming a domain problem into a mathematical optimization problem — still requires expert knowledge. According to a recent survey, 81% of Gurobi’s commercial solver users have advanced degrees, with 49% of them holding a degree in operations research (Gurobi Optimization 2023). This expertise gap prevents many organizations from using optimization, even when it could significantly improve their operations. Examples include inventory management in supermarkets, patient operations in hospitals, transportation policies in small municipalities, energy management in local solar farms, and operations in small businesses or NGOs (Saghafian et al. 2015, Aastrup and Kotzab 2010, Yao et al. 2020, Shakoor et al. 2016). Automating optimization modeling would allow enterprises to improve efficiency using optimization techniques, even when they cannot afford to access optimization experts.

In this paper, we take a first step towards an automated system that can democratize access to advanced optimization modeling expertise. While the ultimate aim of this line of research is to enable real-world stakeholders without any optimization background to formulate and solve their own complex decision-making

problems, we study a necessary precursor to such a system - an automated system to accelerate the workflow of existing optimization practitioners. We specifically focus on helping novice and intermediate optimization users develop, tune, and implement (mixed-integer) linear optimization problems, a flexible and broad class of optimization problems that have been applied successfully in a wide range of real-world use cases such as scheduling and production planning. Our work aims to improve the productivity of optimization experts, akin to the productivity gains realized by software engineers with systems such as GitHub CoPilot, and lay the foundation for more robust systems that can ultimately cater to non-expert users. Given that our focus is on existing optimization practitioners, we assume for the remainder of this paper that the user has already identified a problem for which MILP is an effective tool and can describe the problem in a level of detail such that modeling is possible. Aiding users to find the correct modeling paradigm (e.g., constraint programming, conic optimization) and supporting them in concretizing vague problems remains an important open problem for future research.

Large language models (LLMs) offer a promising way to automate optimization modeling, and increase the reach of the powerful algorithms developed by the operations research community. However, these LLM models are still in their infancy and suffer from major flaws that prevent deployment in important applications. Most problematically, it is difficult to ensure correctness and completeness of the LLM output. We highlight four challenges:

- **Long problem descriptions.** Realistic optimization problems can be exceedingly long and complex: for example, the documentation for the energy system problem described in Holzer et al. (2023) is 60 pages long. Unfortunately, LLMs have a limited context size, and even long-context models perform less well as the input context grows (Liu et al. 2023, Levy et al. 2024). Consequently, LLMs tend to make more mistakes as the length of the problem description increases and perform poorly on complex problems.
- **Large problem data.** The specification of an optimization problem often involves large amounts of data, such as customer attributes or sales of goods. Efficient formulations that ensure fast solve times are essential to handle large problem data and competing with expert-level solutions. Previous approaches to optimization modeling using LLMs, which pass numerical data to the LLM directly and use simple formulations, are thus restricted to the simplest of toy problems.
- **Hallucination.** LLMs are known to *hallucinate*: they may produce answers that sound reasonable, but are incorrect. In the context of optimization, the generated solver code may hallucinate constraints that incorrectly model the problem, or it may hallucinate API calls that do not exist, resulting in code that cannot run. It is especially challenging to verify whether the solution is correct, supposing the code runs without error. For instance, if the solver claims the solution is unbounded, perhaps a constraint has been accidentally omitted from the formulation.
- **Bad models.** The solve time for an optimization problem can depend on the particular modeling formulation chosen, and on how the structure of the problem is communicated to the solver. Optimization

experts spend much of their effort modeling the problem to enhance the efficiency of the solution method. A challenge for LLMs is to produce not just an accurate model, but good code that solves the problem quickly.

This paper studies the feasibility of using LLMs to expand the reach of operations research techniques and analytic methods for decision making by studying the performance of OptiMUS, an automated system we have developed that uses LLMs combined with traditional solvers to model and solve (mixed integer) linear programs (MILPs). Our studies of OptiMUS demonstrate both the promise and pitfalls of the LLM approach to quantitative decision making, and help clarify the effectiveness of LLMs for optimization modeling as of 2024. More broadly, through ablation studies, this paper shows how structured reasoning pipelines provide a clear pathway to address current limitations and enhance the reliability of LLM-based systems in optimization.

Our contributions in this paper are as follows:

- We curate a new dataset, NLP4LP, a comprehensive open-source dataset of 355 optimization problems that ensures adequate coverage of a wide variety of problem types, by combining problems from existing datasets with new problems. Our dataset includes problems that range in difficulty and description length, including real-world problems that are an order of magnitude longer than the problem instances in other MILP modeling datasets. Table 1 compares NLP4LP to existing datasets and Section 3 describes NLP4LP. Our dataset also includes many problems that are less likely to have been used to train existing LLMs, reducing the risk of data leakage compared to previous optimization modeling datasets.
- We develop a modular, LLM-based agent to model and solve optimization problems, which we call OptiMUS (or, occasionally, OptiMUS-0.3 to distinguish it from prior versions). OptiMUS-0.3 employs a connection graph that allows it to process each constraint and objective independently, which allows the system to solve problems with long descriptions and large data files without excessively long prompts.
- We develop several modules designed to improve the performance of OptiMUS and study their impact, including:
  1. *Self-reflective error correction.* We ask OptiMUS to evaluate and correct its output, and assess its own confidence in its output, allowing it to fall back to a more powerful LLM or to user feedback when OptiMUS is unsure its answer is correct.
  2. *Advanced optimization modeling techniques.* We teach OptiMUS advanced optimization modeling techniques by prompting it to identify important structures in the problem, including structure of problems and solver features such as special ordered sets, and to use those structures to model and solve the problem more efficiently.
- We perform ablation studies of the components of our framework to demonstrate which elements improve performance and which can actually degrade performance. We find OptiMUS-0.3 beats the previous state-of-the-art methods on existing datasets by over 22% and on our more challenging test dataset by 24%, and GPT-4o alone by over 38% and 40% on easy and challenging datasets respectively.

- We develop a publicly available web application that allows practitioners to try OptiMUS. This webapp enables human-in-the-loop automated optimization modeling: users can edit LLM outputs and OptiMUS can flag low-confidence LLM outputs to solicit corrections from the user.

An initial version of this work was published in a conference proceeding (AhmadiTeshnizi et al. 2024), that introduced a smaller version of the NLP4LP dataset, and a less robust LLM agent-based system for modeling optimization problems that we will refer to as OptiMUS-0.2. The new enlarged version of the NLP4LP dataset which accompanies this paper has been open-sourced and includes a modular design that enables future researchers to improve on different sub-problems within our framework such as parameter extraction, modeling, or coding. Compared to the earlier version, the modular design of OptiMUS-0.3 substantially improves the robustness of the system for more complicated modeling problems. New modules that enable the system to leverage advanced solver functionality, such as special-order-set constraints, can dramatically speed up the solution time for the associated optimization problems. Additional LLM techniques in OptiMUS-0.3 such as self-reflective error correction improve performance. Finally, this paper is the first to present the OptiMUS webapp.

## 2. Related Work

Formulating an optimization model is often a challenging task even for experts in optimization. Different formulations can lead to significantly different solving times and enable the use of different solvers or solution techniques (Boyd and Vandenberghe 2004). One important skill for an optimization expert is to identify assumptions or relaxations that allow for casting the problem as a well-studied problem type, such as MILP, which enables the use of well-developed optimization solvers. Crafting such an efficient formulation often requires specialized knowledge (Zohrizadeh et al. 2020, Low 2013, Roubíček 2020, Luo et al. 2010, Krarup and Pruzan 1983).

Given the formulation, an optimization expert must choose a solver. Each solver has a distinct interface and capabilities, with associated benefits and downsides (Achterberg 2019, Diamond and Boyd 2016, CPLEX User’s Manual 1987). However, the user manuals for these solvers are often hundreds of pages, making them difficult to understand and use.

*Progress in LLMs.* Recent progress in Natural Language Processing (NLP) has led to the development of LLMs useful for tasks such as answering questions, summarizing text, translating languages, and coding (OpenAI 2023, Touvron et al. 2023, Chowdhery et al. 2022, Wei et al. 2023, Gao et al. 2023, Borgeaud et al. 2022). Connections to other software tools extend the reach and accuracy of LLMs, as demonstrated by plug-ins for code writing and execution (Paranjape et al. 2023, Wei et al. 2023). Yang et al. (2023) use LLMs to directly generate solutions to optimization problems without calling traditional solvers through prompt optimization to improve performance. The approach is limited to small problems since the performance of LLMs degrades as the input context grows, even for explicitly long-context models (Liu et al. 2023).

*Chatbots for Optimization.* In a recent paper, Chen et al. (2023) developed a chatbot to help users detect and fix infeasible optimization problems expressed in `Pyomo` code and serves as an AI assistant rather than as a solver. Li et al. (2023) designed a chatbot to answer natural-language queries about a supply chain optimization model. Alibaba Cloud (2022) also developed a chatbot to facilitate optimization modeling, but there is no public paper or documentation available on it. Lawless et al. (2024) explored the use of LLMs to allow users to customize a simple constraint programming model in the context of meeting scheduling, but does not model more general MILP optimization problems.

*Benchmark-driven Optimization Modeling.* More closely related to our approach, Ramamonjison et al (2023) introduced a dataset of 1101 natural language representations of LP problems. They proposed a two-stage mapping from the natural-language representation to the problem formulation using an intermediate representation. Ramamonjison et al (2022) designed a system to simplify and improve the modeling experience for operations research, but did not offer an end-to-end solution. Xiao et al. (2023) presented a multi-agent cooperative framework called Chain of Experts (CoE) to automatically model and program complex operation research (OR) problems, and evaluated it on NL4Opt and another more complex dataset, ComplexOR, introduced in that paper. Astorga et al. (2024) present a Monte-Carlo Tree Search (MCTS) based approach that decomposes the modeling problem into stages (e.g., modeling variables, then objectives, then constraints, etc.) and uses MCTS to explore the space of plausible models. As mentioned in the paper, this approach is not directly comparable to OptiMUS as it outputs a set of functionally distinct models as opposed to a single formulation and answer. Concurrent with the preparation of this paper, Tang et al. (2024) introduced a new approach to OR modeling with LLMs that uses fine-tuning on a semi-synthetic dataset rather than prompt optimization or agentic models, which they call ORLM (Operations Research Language Model). Yang et al. (2024) also explore a synthetic data generation and fine-tuning approach for a broader class of optimization problems. These fine-tuning approaches are complementary to our work, as any progress towards more capable LLMs for optimization modeling can be combined with our modular prompting framework.

Traditional MILP solvers generally benchmark against the `MIPLIB` benchmark (Gleixner et al. 2021), which offers a diverse collection of MILP problems in standard form. Unfortunately, most of these problems are not associated with a natural-language description, and so cannot be used to study optimization modeling as we do in this paper.

### 3. Dataset

As part of this paper, we introduce a comprehensive open-source dataset of 355 optimization problems we call NLP4LP. Our goal in creating NLP4LP is to provide the community with examples they can use to design optimization modeling tools using natural language systems and to assess their quality. The problems in NLP4LP are partitioned into 1) an easy dataset that contains only LP problems with short descriptions and

scalar parameters, and 2) a hard dataset that contains both LP and MILP problems with longer descriptions and multi-dimensional parameters (Table 1). Unlike other modeling datasets, each instance in the NLP4LP dataset is accompanied by associated code to run the instance, and ground-truth intermediary representations of the problem (i.e., extracted parameters and targets, a list of clauses of the problem represented in natural language,  $\text{\LaTeX}$ , and code, and a solution and optimal value for given problem data). This additional structure enables future researchers to investigate and improve different modules of the OptiMUS framework such as parameter extraction or code implementation, and even test these improvements with the OptiMUS webapp, without implementing the full software stack from scratch. Compared to the initial version of the NLP4LP dataset presented in AhmadiTeshnizi et al. (2024), this version contains more instances from a wider breadth of domains and more complex problems from published papers, and includes additional information for each instance, such as code, as discussed earlier. The dataset is available on Hugging Face and in the supplementary material.

Given the complexity of the task, and the components of each instance (description, sample feasible data, and solution code), gathering high-quality optimization instances is expensive. Unlike traditional ML datasets, NLP4LP does not provide a training set with the same distribution as the test set. As the cost of collecting each instance is high, we prioritize creating a large test set rather than reducing the size of the test set to provide training data. Consequently, NLP4LP is partitioned into 1) a development (dev) set of 22 problems (12 easy and 10 hard) and 2) a test set of 332 problems (277 easy and 55 hard). The dev set is intended to be used to develop optimization modeling tools, while the test set should be used to evaluate such tools. In our experiments and system design, we did not use any part of the test set for development, ensuring an unbiased evaluation of the models.

Optimization modeling is used in many domains. To ensure that NLP4LP comprehends a wide variety of use cases, we first created a list of important types of optimization problems, such as scheduling, cutting, routing, blending, and packing, and a list of common application domains such as sports, government, retail, agriculture, and energy. We tagged each problem NLP4LP with all relevant labels, and continued to gather instances until the dataset contains at least two instances per label. Instances in our dataset are drawn from several sources: some were created by our research team, inspired by problems in textbooks, lecture notes, and research papers, and some are drawn from existing datasets on LP modeling (Bertsimas and Tsitsiklis 1997, Williams 2013, Nace 2020, Xiao et al. 2023, Ramamonjison et al 2023). Most problems have around four labels, and some have as few as two or as many as eight. For more information on problem types and domains refer to Appendix F. Real-world problems are often much longer and more complicated than textbook problems. For the test set, we have included ten long real-world problems from scientific papers. These problems are an order of magnitude longer than the instances in the other datasets (in terms of number of variables and constraints) and seriously challenge the modeling capabilities of existing LLM systems

**Table 1** A comparison on different aspects of complexity for various datasets. The unit for description length is characters

Dataset	Description Length	Instances (#MILP)	Multi-dimensional Parameters
NL4Opt	$518.0 \pm 110.7$	1101 (0)	×
ComplexOR	$497.1 \pm 247.5$	37 (12)	✓
NLP4LP Easy (Ours)	$507.2 \pm 102.6$	287 (0)	✓
NLP4LP Hard (Ours)	$912.3 \pm 498.2$	68 (18)	✓

(including ours). Table 1 shows that our dataset contains instances that are substantially longer than the problem instances in other MILP modeling datasets.

One major concern in evaluation of LLM-based systems is data leakage: has the LLM been trained on the data in the purported test set? This issue is serious and has been known to contaminate the results of several important studies (Oren et al. 2023). In our study, we take several measures to mitigate the risk of leakage, balancing this concern against the goal of broad coverage in our test set.

- Our dataset is guarded by Captcha, and only authenticated Hugging Face users can access it after agreeing to terms and conditions. These safeguards prevent the dataset from being used to train future LLMs.
- Many problems in the dataset are inspired by textbook examples and recent publications, yet their text is entirely original. Moreover, the sources of these questions do not include solutions (either in  $\text{\LaTeX}$  or code). We verified that the content of these questions does not appear anywhere on the internet using a plagiarism detection tool, which yielded a mean originality score of 84.9% and a median originality score of 100%. See Appendix F for details.
- Some LPs in our dataset are drawn from other datasets created prior to the training cutoff for many models (Ramamonjison et al 2023). However, these instances are rephrased and modified in our dataset. Moreover, reporting results on these problems is standard practice (Meta AI 2024, Google 2024). None of the MILPs, which represent the most challenging problems in our dataset, are sourced from existing datasets.

## 4. Methodology

An optimization problem proposes to maximize or minimize a given function  $f$  over a set of allowable inputs  $\mathcal{X}$ . This paper focuses on modeling mixed-integer linear optimization problems (MILP) where  $f$  is a linear function and the set of allowable inputs  $\mathcal{X}$  can be represented by a set of linear inequalities over integer and real valued variables. A given MILP comprises three key components: (1) a set of  $n$  discrete and real variables (i.e., the inputs that can be changed), (2) a set of linear inequalities and equalities that define different constraints on the input variables, and (3) a linear objective function that defines what the problem aims to minimize or maximize. A MILP can be written mathematically as

$$\underset{\{x\}}{\text{minimize}} \quad \sum_{j=1}^n c_j x_j$$

$$\begin{aligned}
&\text{subject to} \quad \sum_{j=1}^n a_{ij}x_j \leq b_i, \quad i = 1, \dots, m_1 \\
&\quad \quad \quad \sum_{j=1}^n a_{ij}x_j = b_i, \quad i = 1, \dots, m_2 \\
&\quad \quad \quad x \in \mathbb{Z}^{n-r} \times \mathbb{R}^r.
\end{aligned}$$

A feasible point  $x^*$  that minimizes the objective function is called a (optimal) solution and has associated optimal objective value  $c^T x^*$ . For brevity, we refer to the equations defining constraints and the objective as different *clauses* in the optimization model. A *parameter* of an optimization model is any numeric coefficient in the optimization model (i.e., elements of  $A$  or  $c$ ). The goal of the OptiMUS system is to go from a natural language description of an optimization problem (e.g., a paragraph of text describing a business problem), which we call the *problem description*, to a decision (assignment of numerical values to each variable) that solves the problem. We assume that the problem description contains all the information needed to model the optimization problem. In practice, users may give ambiguous or incomplete descriptions of an optimization problem (see Wasserkrug et al. (2024), Lawless et al. (2024) for a discussion). Helping users to resolve these ambiguities automatically is an exciting direction for future work.

OptiMUS-0.3 decomposes these two tasks into a sequence of LLM-powered steps that incrementally constructs an optimization model and implements it. An overview of the complete workflow of OptiMUS-0.3 is outlined in Fig. 1. At a high level, OptiMUS-0.3 begins by 1) identifying parameters from the problem description and writing a natural-language description of each clause in the problem. It then 2) formulates both the clauses and decision variables of the optimization model with mathematical precision, using  $\text{\LaTeX}$  code. Finally, 3) OptiMUS-0.3 uses the  $\text{\LaTeX}$  model of the optimization problem and the extracted parameters to generate code snippets for each clause in Python and assembles the snippets into a single runnable code file. By running this file, OptiMUS-0.3 outputs an optimal solution and optimal objective value. Throughout this process, OptiMUS-0.3 maintains a *state* that documents what is known about the optimization model, as well as a *connection graph* that tracks which variables appear in which constraint (detailed in Section 4.1). The current state is supplied in the prompt of every LLM-powered component (detailed in Section 4.2). Each LLM-powered component of OptiMUS-0.3 also incorporates an Error Correction (EC) module (detailed in Section 4.3) that catches common errors and improves the accuracy of the modeling and implementation.

Our current implementation writes code that uses Gurobi (Achterberg 2019), a leading MILP solver, and its associated Python API gurobipy. However, we have designed OptiMUS-0.3 to be easy to extend to other modeling languages and backends. We expect it to perform equally well with other solvers (e.g., CPLEX, SCIP) that have well-documented APIs, as well as with other modeling languages such as cvxpy (Diamond and Boyd 2016).



Many of the advances in optimization solvers in the last several decades rely on the fact that solvers are faster when they can exploit particular structures within the optimization problem. One of the main tasks of an optimization expert is to identify these structures so as to choose an appropriate solver or parameter settings. However, detecting useful structures in a given optimization problem can be very challenging, since it requires both a deep understanding of the mathematical structure of the problem and the capabilities of existing solvers. In view of this challenge, OptiMUS-0.3 includes two specialized optimization modules: 1) a Structure Detection Agent, that detects variables (e.g., Indicator variables) and constraints (e.g., SOS constraints) with special structure that can be exploited by the solver, and 2) a customized optimization coding agent, that can leverage callbacks in the solver and perform basic column and constraint sifting to help improve the runtime of the code. These elements represent a first step towards automating more advanced decomposition algorithms with an automated modeling framework.

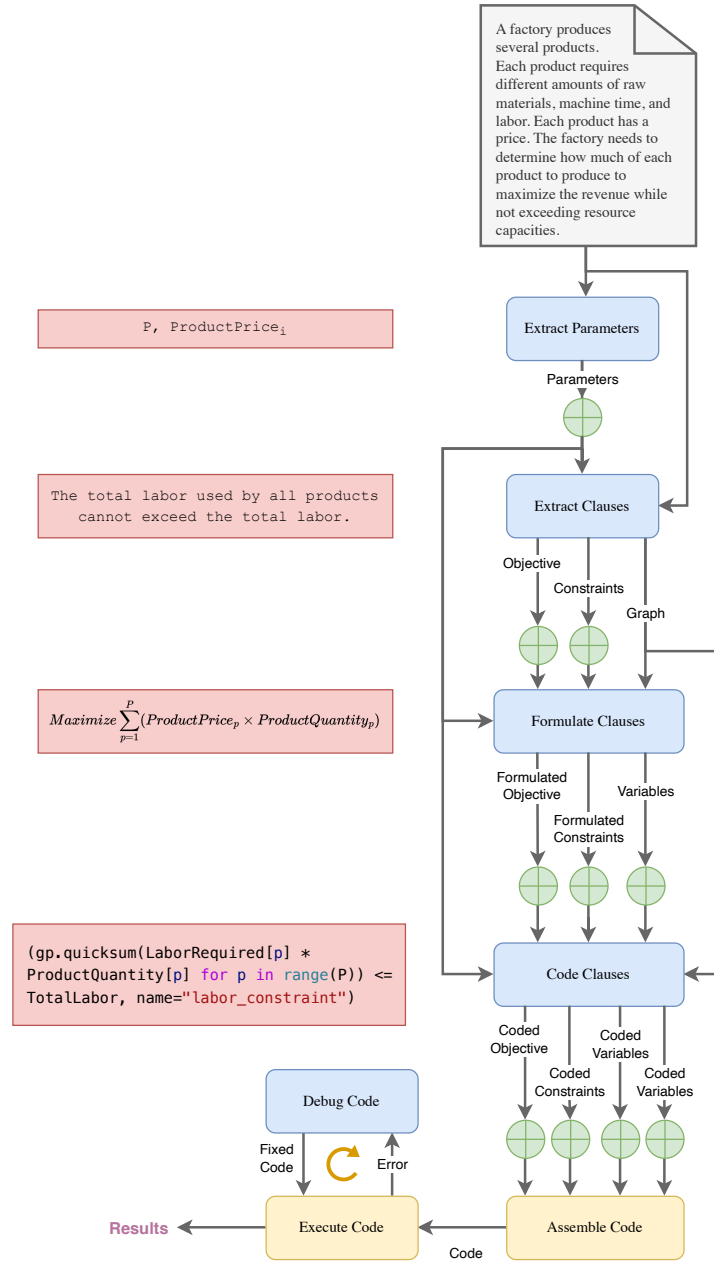
OptiMUS-0.3 can be run as an end-to-end pipeline or through an interactive web application designed to facilitate human-system collaboration. The OptiMUS web application uses the same underlying components as the end-to-end framework but provides additional opportunities for users to audit and edit the outputs of each stage of the OptiMUS pipeline. We defer a deeper discussion of the system to Section 6, and mention specific interaction mechanisms available in the user interface earlier when relevant.

The remainder of this section details the key components of the OptiMUS-0.3 framework. Section 4.1 describes the underlying state and connection graph of the system. Section 4.2 describes a sample LLM component of the system. The remaining sections describe specialized modules that improve the performance of the system, including error correction and debugging (Section 4.3), specialized modeling tools (Section 4.4), and advanced optimization coding (Section 4.5). OptiMUS-0.3, and all associated components, have been released as open-source code and are publicly available to use through a web application, described in Section 6.

#### 4.1. State

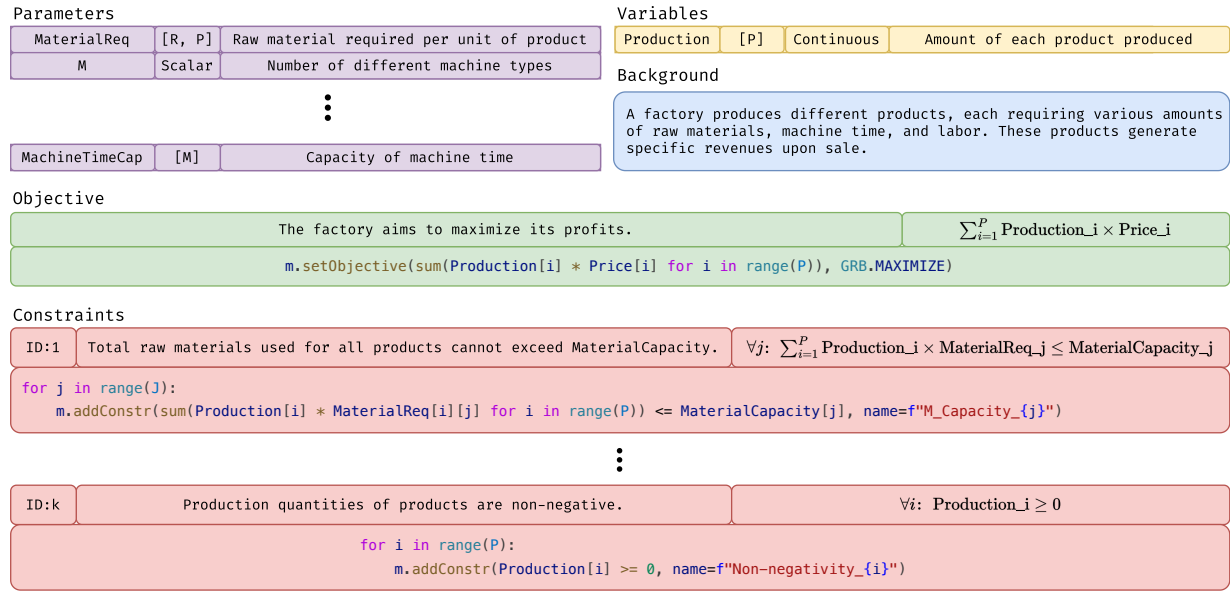
OptiMUS-0.3 manages and modifies the solution using states saved in JSON format. The state consists of the following components:

- **Parameters:** OptiMUS-0.3 can choose a symbol for each parameter, infer its shape, and define the parameter if it is not explicitly included in the problem statement. Numerical data from the problem statement is replaced by a symbol and stored separately for later use, ensuring that the parameter list remains concise and easy to include in future prompts.
- **Clauses:** Each clause (objective or constraint) consists of a natural language description, a  $\text{\LaTeX}$  formulation, and code implementing the clause.
- **Variables:** Like parameters, each variable has a symbol, a shape, and a definition. Unlike parameters, each variable also has a type, which can be *Continuous*, *Integer*, or *Binary*.



**Figure 1** Optimus uses a sequential process with error correction (green circles) to model and solve optimization problems. It extracts the parameters and clauses, formulates them, generates the code for each, and finally synthesizes complete code and runs it. If the code raises an error, Optimus uses an iterative debugging loop to fix the code. Red boxes provide sample outputs of each step of the pipeline. Blue boxes indicate LLM-based components, and yellow boxes indicate deterministic components.

- **Background:** The background is a short string that explains the real-world context for the problem. Including this string in every prompt helps improve common sense reasoning.



**Figure 2** A completed state for a factory production optimization problem

- **Connection Graph:** The connection graph is a bipartite graph  $G = (V, E)$  that links clauses to their associated parameters and variables. Formally, for every variable  $x_j$  (clause  $c_i$ ) in our formulation we create a node  $v_{x_j}$  ( $v_{c_i}$ ). In total, the connection graph consists of  $n + m + 1$  nodes  $v \in V$  corresponding to  $n$  variables,  $m$  constraints and the objective. We create an edge  $e \in E$  between a variable node  $v_{x_j}$  and a clause node  $v_{c_i}$  if the variable participates in the clause (i.e.,  $(v_{x_j}, v_{c_i}) \in E$  if  $a_{ij} \neq 0$ ). We build the connection graph iteratively by adding nodes and edges as more constraints and variables are formulated. This structure is utilized during the coding phase to ensure that only the relevant parameters and constraints are passed to the LLM.

The state is initialized as an empty object, and is completed step by step throughout the process. An example of the completed state for a factory optimization problem appears in Fig. 2.

## 4.2. Sample LLM Component

OptiMUS-0.3 uses LLMs as a flexible tool to perform various tasks for the overall system, including extracting parameters, modeling clauses, and coding the components (see Figure 1). In the interest of brevity, this section outlines the design of a typical LLM component with the OptiMUS framework. Each LLM component is governed by a natural language directive called a *prompt*. Within OptiMUS-0.3, each prompt includes the following key pieces of information:

- **Task Description:** We describe the specific task, e.g., “Your job is to extract natural language constraints for this paragraph defining an optimization problem”.
- **Problem Context:** We include relevant information about the problem and the current solution in the prompt. During the formulation step, this involves the specific clause being addressed, along with

all parameters and variables defined so far. During the coding step, this involves the clause and its formulation, supplemented by related parameters and variables which are dynamically loaded using the connection graph.

- **Examples:** We include a fixed set of sample outputs for the task in the prompt to help the LLM perform In-Context Learning (Dong et al. 2024). For modules that output Python code, we also provide examples detailing the functionality of the gurobipy API.

For a comprehensive list of prompts used in OptiMUS-0.3 see our code<sup>1</sup>. Within the OptiMUS webapp, users can inspect and correct the output of every LLM-powered component of the system.

### 4.3. Error Correction

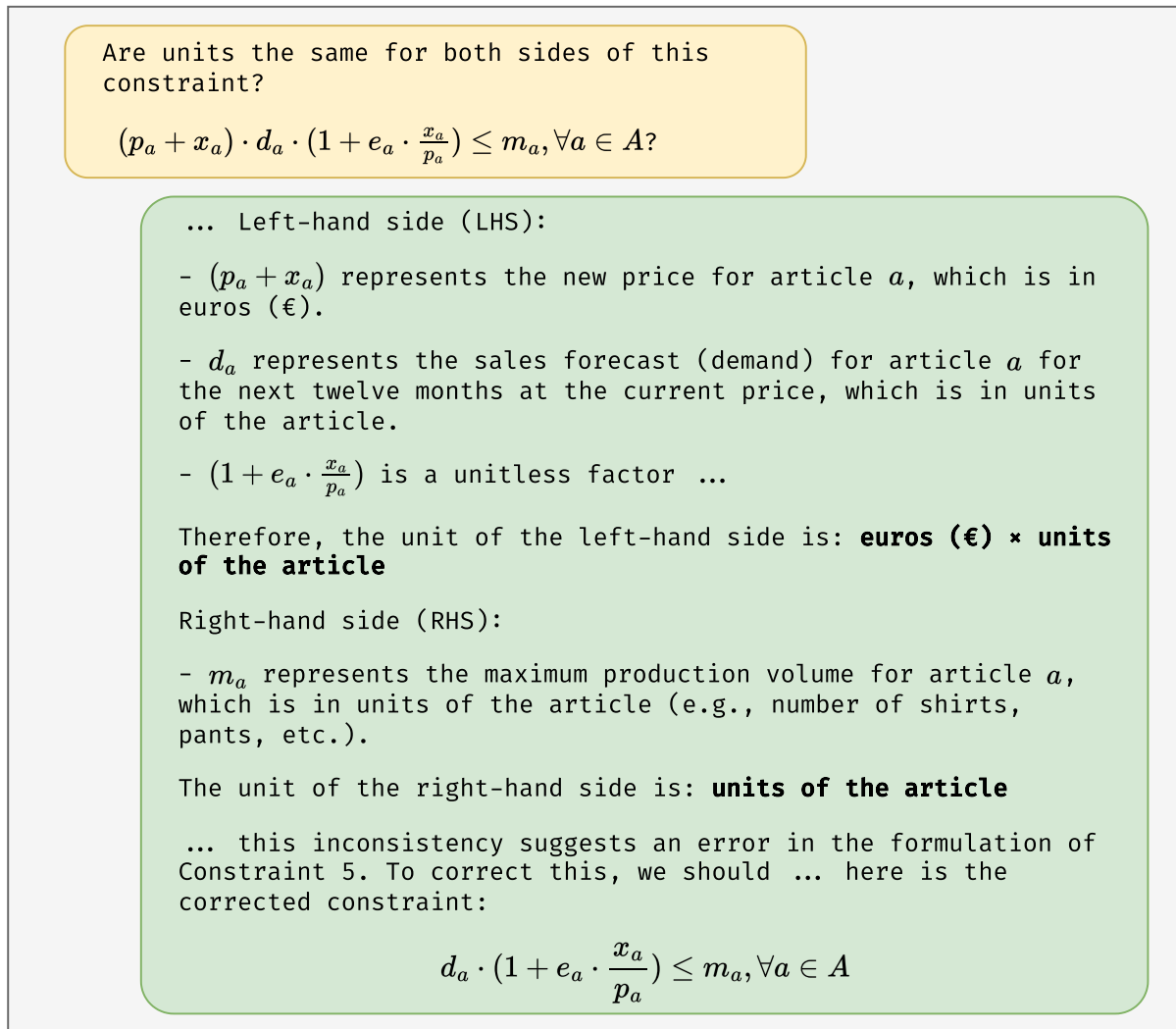
To build a trustworthy and reliable system using LLMs, it is important to mitigate the impact of LLM hallucinations. In the context of optimization modeling, an LLM might generate incorrect parameters, redundant mathematical constraints, or erroneous code. To address hallucinations, OptiMUS-0.3 uses two main error correction techniques: *reflective prompts* and *confidence-based user feedback*.

**4.3.1. Reflective prompts** LLMs can often identify and fix their mistakes by using reflective prompts (Shinn et al. 2023), a strategy by which an LLM is asked to reflect upon its output and decide if it needs to be corrected. We analyzed the most common types of mistakes made by LLMs on the development dataset at each step of the process and designed reflective prompts to address each one (see Appendix B for a full list of reflective prompts used in OptiMUS-0.3). Using these prompts significantly reduces the modeling error rate. For instance, Figure 3 shows an example where the reflexive prompt *Are units the same for both sides of this constraint?* can help correct a modeling error. The resulting improvement to OptiMUS is substantial (see Table 3), and substantiates the claim that domain-specific knowledge can be used to develop LLM applications that outperform general-purpose LLMs.

**4.3.2. Confidence-based user/LLM feedback** In the context of real-world optimization modeling, it is exceedingly rare for a natural-language problem description to correspond to an unambiguous MILP formulation. Rather, the process of optimization modeling generally involves many cycles of formulation and user feedback to arrive at a correct model for the problem. OptiMUS-0.3 simulates this process in order to detect errors and improve its solve rate. OptiMUS-0.3 can use a confidence-based reflective method to decide when to request help with its formulation, either from a user (if available) or another (more powerful) LLM.

First, following the methodology of existing work on confidence calibration in LLMs (Virk et al. 2024), we prompt the system to assess its confidence in its output on a scale of 1 to 5. If the confidence score is less than 5, OptiMUS-0.3 asks the user (or a stronger LLM) to help. The user can decide to either keep the output

<sup>1</sup> included in supplementary materials



**Figure 3** OptiMUS-0.3 can fix its constraint modeling errors when prompted “Are units the same for both sides of  $C$ ?”

(e.g., constraint, code), remove it, or modify it. In the OptiMUS webapp (Section 6), this feature is used to highlight outputs that merit particular attention. In our experiments (Section 5), this feature is turned off by default to facilitate streamlined and reproducible testing without a user in the loop. See Fig. 4 for an example of how a weaker LLM (Llama-3-70B) can call on a stronger LLM (GPT-4o) when it is not confident in its solution.

**4.3.3. Debug Code** After the complete code is synthesized and executed by OptiMUS, any errors encountered during execution are passed to the LLM for debugging. At this stage, the LLM has the ability to modify any part of the code as needed. This step is particularly useful for resolving errors that arise from inconsistencies across individual code snippets, such as slight variations in variable names or formatting.

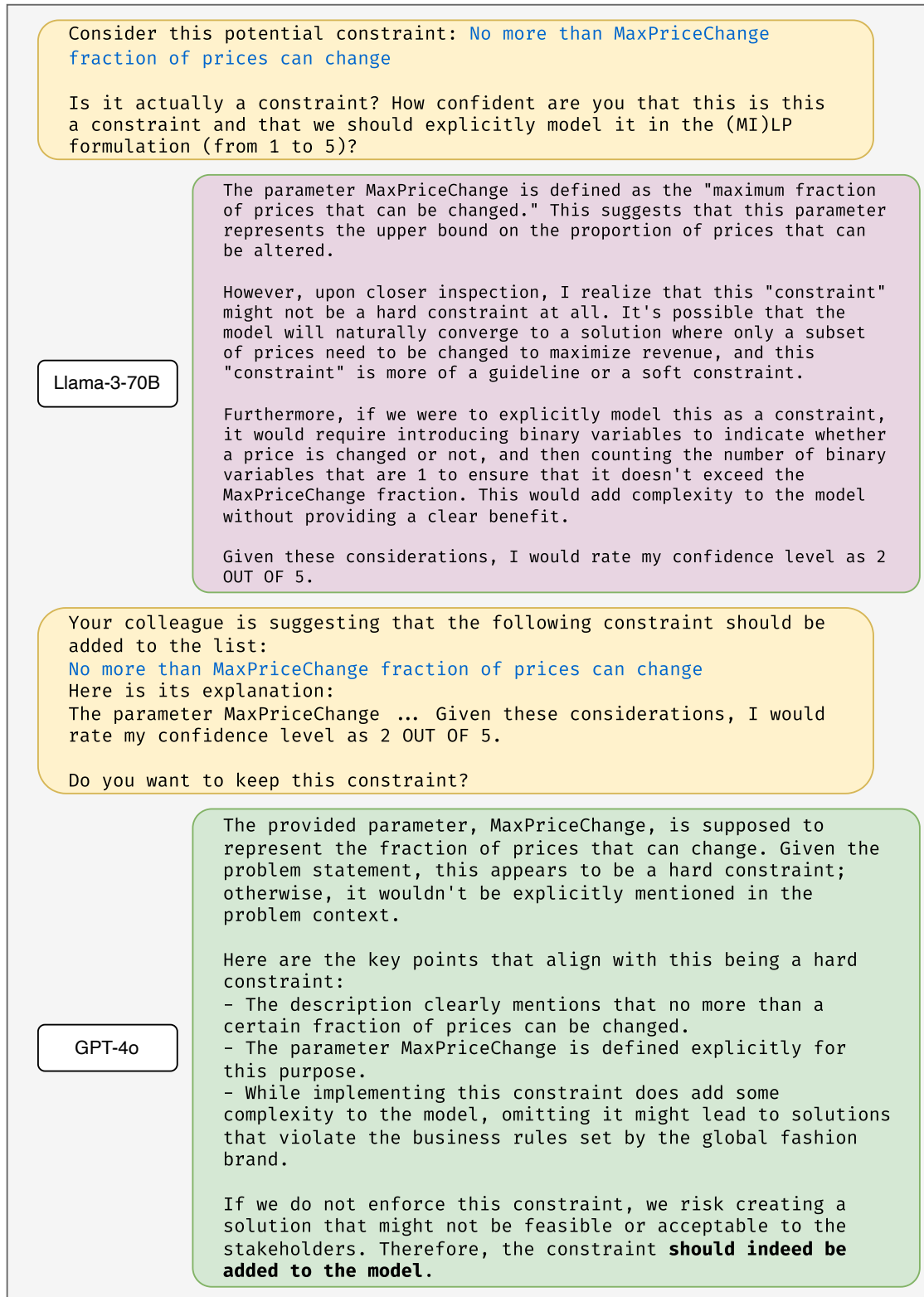


Figure 4 GPT-4o can provide feedback when Llama-3-70B lacks the confidence to identify a constraint.

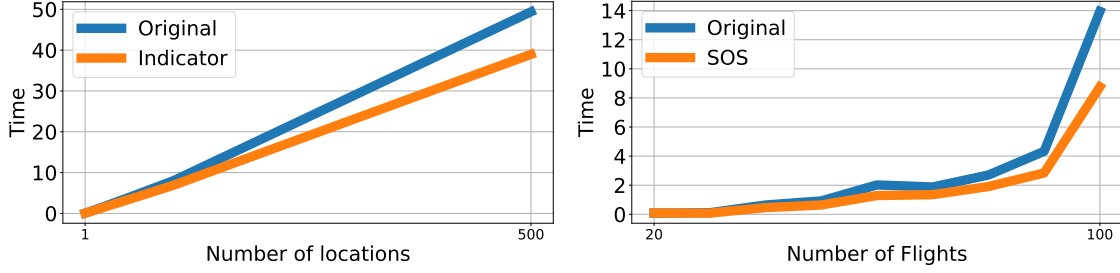
The debugging process is repeated iteratively until the code executes successfully, with a maximum of five attempts.

#### 4.4. Structure Detection Agent

OptiMUS maintains a pool of optimization structures commonly used in optimization software. These structures can be specified explicitly in many modern optimization solver interfaces (Achterberg 2019, CPLEX User’s Manual 1987, Gamrath et al. 2016), which can enhance problem-solving performance and simplify the code. Typical examples of structure include Special Ordered Set (SOS) (Beale and Forrest 1976), indicator variables, semi-continuous variables, and piecewise-linear constraints. At least one of these structures appears in approximately 10% of the NLP4LP dataset (see Appendix C). We expect these structures to be even more common, and thus important to exploit, in industrial applications.

To detect structure in a constraint (or variable), OptiMUS iterates through these structures and formats them into a structure detection prompt. Within each prompt, the LLM is provided with the structure description, explained by an example illustrating how the structure should be exploited. The LLM is asked to decide whether the structure is relevant to the existing formulation. Upon identifying the appropriate structure, the formulation is adjusted to highlight the problem structure. For instance, a set of constraints that indicate only one of a group of decision variables can be non-zero may be reformulated as type-1 SOS constraint. This structure is then conveyed to gurobi via the Python interfaces for special structures such as SOS constraints or indicator variables. This information can either be exploited within the branch and bound solver (e.g., using SOS constraints for branching rules) or allows the solver to automatically reformulate the structure into linear constraints (i.e., a big-M formulation for indicator constraints) (Gurobi Optimization 2023). While it may seem counter-intuitive to identify structure in a MILP only to have that structure re-formulated into linear constraints, in practice we found that automated methods within Gurobi can generate tighter formulations than OptiMUS alone (e.g., a better selection for big-M values). Figure 5 shows two instances where identifying SOS constraints and indicator variables leads to faster performance than a naive implementation. Full details of these examples are included in Appendix C. The structure detection agent is run during the formulate clauses stage of the algorithm.

OptiMUS also maintains a *problem structure* pool that enumerates combinatorial optimization problems with special structure that can be exploited in fast custom solvers (e.g., SAT, network flow, routing problems). For instance, the traveling salesman problem (TSP) can be solved much more efficiently using a specialized solver such as Concorde as opposed to general-purpose MILP solvers (Cook et al. 2011). OptiMUS iterates through this pool of special problem structures given an initial natural language description of the problem. When the optimization problem falls into one of the problem types in the structure pool, OptiMUS identifies the particular problem structure and suggests using a customized solver via a notification in the OptiMUS webapp.



**Figure 5** Impact of Structure Detection on OptiMUS Performance. (Left) Speedup of solving a facility location optimization problem with indicator variables with the structure identified to the solver (indicator) versus with OptiMUS modeling of indicator variable. (Right) Speedup of solving a flight assignment problem with SOS-constraints with the structure identified to the solver (SOS) versus with OptiMUS modeling of the SOS constraints. Full details of problem instances are included in Appendix C.

#### 4.5. Advanced Optimization Coding Agent

For extremely large-scale optimization problems, optimization solvers are often embedded into a high-level optimization framework and called as a subroutine: examples include decomposition algorithms such as column generation, Benders decomposition, and cutting plane methods. To leverage variable and constraint structure in this large-scale context, the subproblem solver must invoke an advanced solver interface (for example, call-backs, model attribute query and analysis). While we leave the problem of developing a fully automated system that can perform advanced decomposition algorithms to future work, we have added several features to OptiMUS that advance this vision. OptiMUS implements a dedicated coding agent that exploits advanced solver functionality such as callbacks, and can iteratively call a solver within a simple variable or constraint sifting scheme (also known as delayed constraint generation). Neither of these modules requires explicitly reformulating the problem (e.g., generating a pricing problem within column generation), but allow for better computational performance than a naive implementation.

Similar to the structure detection agent, the optimization coding agent maintains a series of template prompts for this functionality. Within each prompt, the LLM is provided with the purpose of the template. We illustrate this functionality with an application to sifting, a simple version of constraint or column generation, in Section 4.5.1. The advanced coding agent runs during the Code Clauses stage of OptiMUS (see Fig. 1).

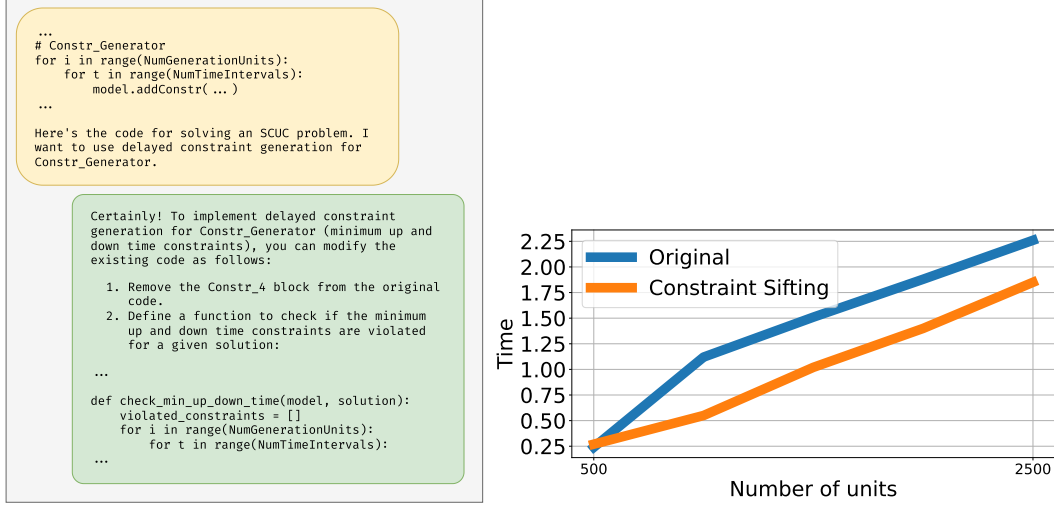
**4.5.1. Sifting** Sifting is a large-scale linear optimization framework initially proposed in (Bixby et al. 1992). Consider the primal linear optimization problem

$$\min_x c^\top x \quad \text{subject to} \quad Ax = b, \quad x \geq 0,$$

where  $A \in \mathbb{R}^{m \times n}$  and  $n \gg m$ . On the primal side, sifting combines the idea of primal simplex and column generation by solving a restricted master problem,

$$\min_{x_S} c_S^\top x_S \quad \text{subject to} \quad A_S x_S = b, \quad x_S \geq 0,$$





**Figure 6** Left: constraint sifting prompt. Right: performance plot for solving to a 5% gap.

for a subset of columns  $S \subseteq \{1, \dots, n\}$ . The restricted problem can be solved by any linear programming subroutine. Its dual variable  $y_S$  will be used to price out dual infeasible columns

$$I_S := \{j : c_j - A_j^\top y < 0\}.$$

Variable sifting works by iteratively updating  $S \leftarrow S \cup I_S$  until  $S = \emptyset$ . The idea can be generalized to the dual side when  $m \gg n$ , where a restricted problem contains a subset of rows. We provide an example of constraint sifting in OptiMUS in Fig. 6 on a security-constrained unit commitment problem. For details of the problem see Appendix D. Most of the constraints in the instance are inactive, making it a good candidate for sifting.

## 5. Experiments

In this section, we conduct a comprehensive evaluation of OptiMUS-0.3. We showcase the superior performance of OptiMUS-0.3 on both the NLP4LP dataset and two existing benchmark datasets for optimization modeling highlighting its strengths and weaknesses. An ablation study demonstrates the impact of different system components on our results, and a sensitivity analysis probes the internal dynamics of OptiMUS-0.3. We conclude this section by identifying failure cases and potential areas for further improvement.

### 5.1. Overall Performance

**Baselines** To evaluate the overall performance of OptiMUS, we compare it with standard prompting, Reflexion, and Chain-of-Experts (CoE) (Shinn et al. 2023, Xiao et al. 2023). Reflexion is the highest-performing general-purpose framework, and CoE is a recently proposed agentic framework for natural-language optimization modeling. We include self-reported results for two state-of-the-art auto-formulation methods based on fine-tuning LLMs specifically for optimization modeling: ORLM (Tang et al. 2024), and

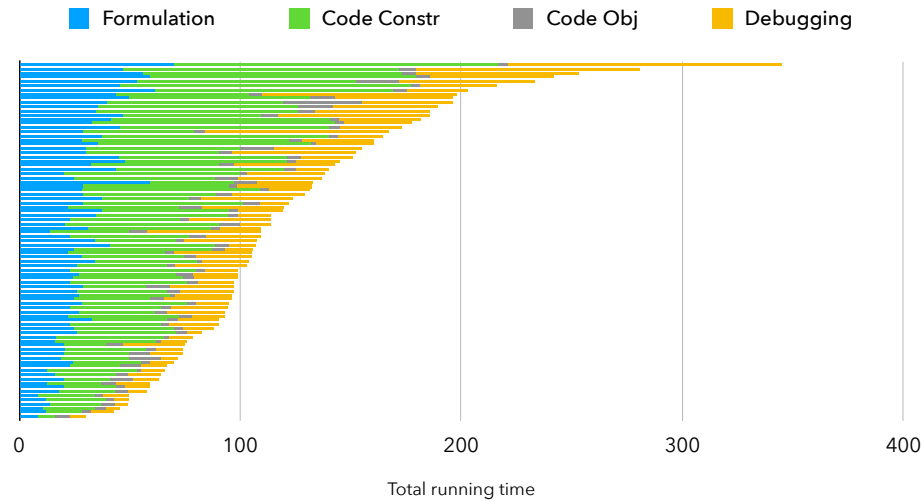
**Table 2 Accuracy of auto-formulation methods on benchmark datasets. \* indicates performance taken directly from original paper, and blanks indicate no results reported on the dataset.**

	LLM	NL4OPT	NLP4LP	IndustryOR
<i>Methods based on direct prompting</i>				
Standard	GPT-4o	47.3%	33.2%	28.0%
Standard	o1	> 95%	68.8%	44.0%
Reflexion	GPT-4o	53.0%	42.6%	–
<i>Methods based on fine-tuning LLMs</i>				
LLMOPT	Qwen1.5-14B	93.0%*	83.8%*	46.0%*
ORLM	Deepseek-Math	86.5%*	72.9%*	38.0%*
<i>Methods based on agentic frameworks</i>				
CoE	GPT-4o	64.2%	49.2%	–
OptiMUS-0.2	GPT-4o	78.8%	68.0%	–
OptiMUS-0.3	GPT-4o	86.6%	73.7%	37.0%
OptiMUS-0.3	o1	–	80.6%	46.0%

LLMOPT (Jiang et al. 2024). We also include results from the initial conference publication of this work, which we denote as OptiMUS-0.2. We report results for the OptiMUS-0.3 framework with both GPT-4o (OpenAI 2023) and o1 (Jaech et al. 2024), an LLM that automatically includes more advanced reasoning at the expense of longer inference times. Due to the cost of running o1, we report performance with standard prompting from Lu et al. (2025) and only run the OptiMUS framework on datasets where o1 alone does not already achieve near-perfect accuracy. In our experiments, we use the default API parameters for all LLM calls.

**Evaluation Metrics and Datasets** Three main metrics have been used to assess the accuracy of LP modeling tools in the literature: accuracy, compilation error (CE) rate, and runtime error (RE) rate. However, a method can generate an irrelevant short code that runs, or fix runtime and compilation errors by completely removing relevant sections of the code. Hence, we only compare the models’ accuracy. Accuracy is defined as the number of instances correctly solved. An instance is considered as correctly solved only if 1) the code runs successfully, 2) the optimal value is correct, and 3) the optimal solution is correct. Optimal values and solutions are obtained from the dataset or by solving the problems manually. Since the output of the code does not necessarily match the format of the optimal solution in the dataset, we use an LLM to determine whether the solution generated by OptiMUS matches the optimal solution in the dataset. In addition to the NLP4LP dataset introduced in this paper, we include results for two existing datasets for optimization modeling: NL4OPT (Ramamonjison et al 2023) and IndustryOR (Tang et al. 2024). The latter dataset includes 100 real-world case studies from eight different industries. In the tests, we assume that the parameter names are given and identified in the problem description.

Results are presented in Table 2. OptiMUS-0.3 outperforms all other direct prompting and agentic methods in all datasets by a large margin, beating the standard GPT-4o baseline by over 40% on the NLP4LP



**Figure 7** Distribution of running time (s) for different stages in the OptiMUS framework.

dataset. This improvement persists even with more advanced reasoning models such as o1. This impressive performance improvement highlights the importance of modularity and structure compared to a single prompt to solve complex problems using LLMs.

Remarkably, the OptiMUS framework also remains competitive with finetuning approaches, beating the performance of ORLM and remaining within 4% of LLMOPT. Note that fine-tuning based frameworks — which must be trained and deployed by the user, and updated via expensive retraining to take advantage of a new base model — represent a much more logistically complex approach compared to agentic frameworks that can be run with simple API calls to LLM providers and instantly improved by access to a new base model. We see that modular design of an agentic system can close the performance gap between direct prompting and finetuning, giving practitioners a simpler alternative to state-of-the-art systems. The improved performance of the OptiMUS-0.3 framework with o1 over GPT-4o also underscores a key benefit of an agentic framework: as more advanced models are released, the framework can be easily adjusted to run with the new models by substituting a simple API call whereas finetuning methods require repeating the costly finetuning procedure.

*On Computation Time* In Figure 7 we investigate the running time of the OptiMUS framework across 85 randomly selected instances of NLP4LP using GPT-4o and along different stages of the pipeline. The OptiMUS-0.3 framework takes no more than 350 seconds (median time 108 seconds) to model and solve any problem instance. In contrast, recent work highlighted that human subject matter experts take on average approximately 150 *minutes* for similar tasks (Tang et al. 2024). Coding and debugging represent the most time-consuming parts of the OptiMUS pipeline. Smaller LLM agents specialized to coding (e.g., codex) could further improve the computation time of the framework.

*On Stochasticity in LLM Outputs* LLM outputs are well-known to be probabilistic, meaning that calling an LLM multiple times can lead to different outputs. To evaluate the impact of this stochasticity on the performance of OptiMUS-0.3 we re-ran 10 hard instances of the NLP4LP dataset over five different random seeds. Across all instances, the result (i.e., success or failure) remained consistent across all 5 random seeds, providing strong evidence that OptiMUS is robust to this source of stochasticity. We hypothesize that the EC modules help mitigate performance variability arising from stochastic LLM outputs. Conversely, we do not observe improved performance for multiple samples of the LLM output. One interpretation is that for hard instances, new strategies, and not simply repeated samples, are needed to improve performance.

*On Calibration* A key feature of the OptiMUS-0.3 system is the confidence-based LLM feedback detailed in Section 4.3.2. To evaluate the calibration of the LLM-based confidence scores we sampled 40 confidence scores from the clause modeling module and enlisted a doctoral student in operations research to manually annotate whether each clause was correct. Of the 40 clauses, the LLM responded with a confidence score of 5/5 in 28/40 of the instances and was correct 100% of the time, and a score of less than 5 in 12/40 instances and was correct 91.7% of the time. This provides promising evidence that the LLM-based confidence scores can effectively guide a user to identify errors. Exploring more advanced mechanisms for eliciting calibrated confidence scores (Xiong et al. 2023) is a promising direction for future work.

## 5.2. Ablation Study

Table 3 shows the impact of the choice of LLM on the performance of OptiMUS, broken down by its performance on hard and easy instances of the NLP4LP dataset. We evaluate the OptiMUS system with a leading closed-source LLM, GPT-4o, a reasoning model, o1, as well as a weaker open-source alternative, LLaMa-3.1-70B-Instruct (Meta AI 2024). As expected, the performance of the OptiMUS-0.3 framework degrades when less capable LLMs, such as LLaMa-3.1-70B-Instruct, are used instead of GPT-4o or o1, especially on hard modeling tasks. Optimization tasks need complex reasoning, and smaller models may be better suited for simpler tasks. However, the OptiMUS-0.3 framework with LLaMa-3.1-70B-Instruct outperforms GPT-4o with previous state-of-the-art prompting strategies, including Reflexion, on easy problems. We see that OptiMUS-0.3 is more robust to the quality of the underlying LLM, allowing weaker cheaper models to match the performance of stronger models on simple optimization modeling tasks.

Table 3 also shows the impact of different components of the OptiMUS framework on its performance. Debugging (i.e., running error correction on the output of the Coding stage) has the largest impact on the performance of OptiMUS, increasing its accuracy on hard problems by almost 50%. The debugging step, is particularly important for harder problems, given LLMs often make a few small mistakes when they generate code. Error correction also leads to modest improvements in the performance of OptiMUS. Running error correction after extracting clauses is particularly important for difficult problems. Table 4 further breaks

down the impact of error correction methods on correct and incorrect items (i.e., constraints). In both constraint extraction and modeling, EC modules are able to fix a large fraction of errors without modifying most correct items. The LLM feedback step also marginally improves the performance of the system. Note that we do not include the impact of the structure detection agent and advanced optimization coding agent as they do not impact the accuracy of the framework, but just the efficiency (i.e., runtime). For examples showcasing the impact of these components, see Appendix C.

Fig. 8 (left) shows the relation between the performance and the number of debugging steps for OptiMUS with both GPT-4o and LLaMa-3.1-70B-Instruct. In both models, the first few debugging steps improve performance, but further steps do not result in improvement.

### 5.3. Failure Cases

We next analyze the most common reasons why OptiMUS fails. We categorized failure cases manually via a qualitative coding process based in grounded theory (Charmaz 2006). We followed an open-coding process to analyze instances on which OptiMUS failed and tag the type of error. These initial codes were then grouped together into the following coarser categories:

- Extraction errors: OptiMUS extracts the wrong natural language constraints or objective (e.g., price is non-negative where price is a parameter), or fails to extract all of the constraints from the description.
- Formulation errors: OptiMUS formulates a mathematical model that does not match the natural language constraint (e.g., it uses the wrong parameter to lower-bound a variable).
- Coding errors: OptiMUS does not always generate error-free code even after debugging. Coding errors most often occur when the LLM is confused by the language used (e.g. a runtime error due to accessing indices that do not exist in a parameter array).

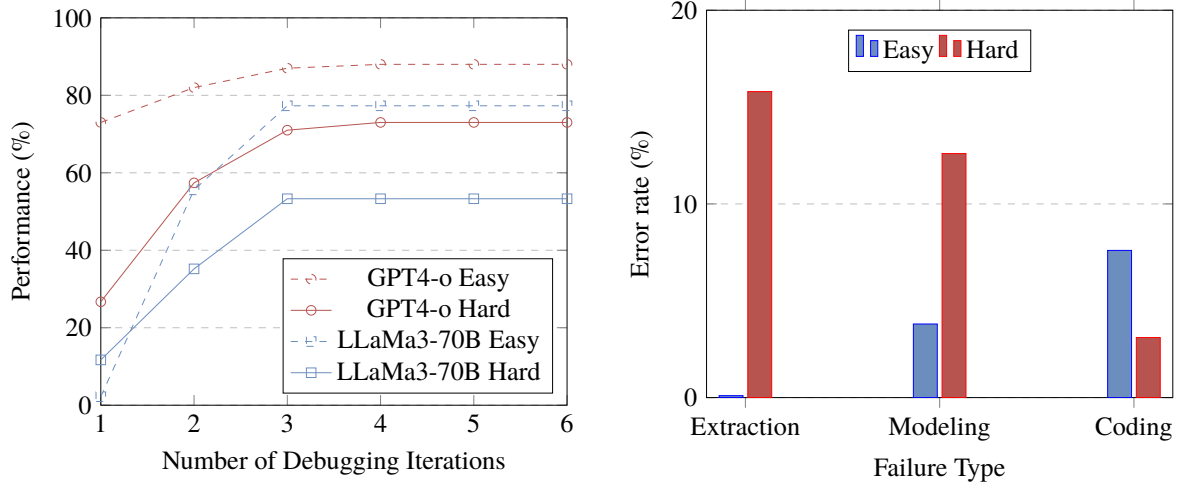
As shown in Fig. 8, OptiMUS successfully extracts almost all the correct clauses for the easy dataset. This task is easy on the easy dataset because the easy instances are all LPs and involve only scalar values. However, for the hard instances — comprising MILP problems and multi-dimensional variables — parameter extraction and modeling are considerably more challenging than coding the resulting model. Note that the coding error rate is lower for the hard dataset because OptiMUS sometimes fails in modeling or formulating hard clauses, producing clauses that cannot be coded, which we do not include in the coding statistics.

## 6. Web Application

As discussed in 1, the OptiMUS project aims to 1) help optimization experts develop and maintain their models with ease and 2) empower domain experts to make better decisions faster and reduce costs. To further these goals, we designed and built a user-centric web app with an intuitive interface that enables users to interact seamlessly with the system and leverage the advantages of large language models (LLMs). The webapp allows users to follow the overall process of modeling an optimization model including extracting parameters and clauses, modeling each clause, and finally coding it up. Importantly, users can correct any output of the

**Table 3** Ablation studies on OptiMUS-0.3

	NL4OPT	NLP4LP
<b>Importance of Different Components</b>		
w/o Debugging	73.2%	26.7%
w/o Extraction EC	86.7%	60.5%
w/o Modeling EC	83.8%	65.7%
w/o LLM Feedback	86.6%	68.4%
<b>OptiMUS-0.3 (GPT-4o)</b>	<b>86.6%</b>	<b>73.7%</b>
<b>Performance with Different LLMs</b>		
LLaMa3.1-70B-Instruct	70.4%	31.5%
<b>GPT-4o</b>	<b>86.6%</b>	<b>73.7%</b>
<b>o1</b>	—	<b>80.6%</b>

**Figure 8** Left) Further debugging iterations improve performance. Right) For harder problems, most failures arise from clause extraction mistakes. For easier problems, most failures are due to coding errors.**Table 4** Error correction methods can find and fix a large fraction of errors in constraint extraction (left) and constraint modeling (right), without modifying most correct items. (Perfect performance is diagonal.)

	Not Modified	Modified
Right	219	7
Wrong	9	41

	Not Modified	Modified
Right	231	2
Wrong	4	22

framework, including natural language clauses,  $\text{\LaTeX}$  model, and code, to provide better supervision over the entire process. By allowing users to observe and provide input throughout each step, the web app brings significant speed and convenience to the modeling process while reducing the risk of errors. Figure 9 shows a sample model in the Coding stage. Further details of all the different stages and features incorporated in the webapp can be found in Appendix E. The webapp is publicly available at <https://optimus-solver.com/>.



**Figure 9** Sample view of the OptiMUS webapp during the coding phase in which individual clauses are translated to Python code. Users have the ability to see the confidence scores, correct the code directly in the panes to the right, or ask OptiMUS to regenerate code.

## 7. Conclusion

How can we leverage LLMs to achieve complex goals? This paper interrogates this question in the domain of optimization and showcases the importance of modular structure. We develop OptiMUS-0.3, a modular LLM-based agent designed to formulate and solve optimization problems from natural language descriptions. Our research illustrates the potential to automate aspects of the optimization process by combining LLMs with traditional solvers. To assess the performance of OptiMUS-0.3, we released NLP4LP, a comprehensive dataset of long and challenging optimization problems. OptiMUS-0.3 achieves SOTA performance and scales to problems with large amounts of data and long descriptions.

Real-world optimization problems are often complex and multifaceted. Developing LLM-based solutions for these problems requires domain-specific considerations, including integrating existing optimization techniques to leverage problem structure. We are at the early stages of this research, but anticipate significant developments that will enable these systems to address more complex, industrial-level problems. It is interesting to notice that the challenge of using AI for an applied domain is much larger in safety-critical domains such as self-driving, which demand extremely high accuracy, than in domains where AI can function as an assistant and where answers can be checked by humans, as in theorem-proving or optimization. Here, AI systems with moderate accuracy can still usefully augment human effort.

**Future directions.** The OptiMUS system presented here serves as a prototype for more powerful AI-assisted optimization modeling. Much more work remains to make such a system reliable and to facilitate the use of classical optimization tools for domain users who know little about optimization. While some

improvements can be expected as language models grow more powerful, we expect that substantial breakthroughs by the optimization community will be required to address many of the following challenges:

- **Reliability:** Optimization algorithms generally offer guaranteed reliability. How can a system like OptiMUS reliably solve large-scale problems using unreliable LLMs that suffer hallucinations, overconfidence, and instability? What kinds of guarantees can be made, and how expensive might these be? How should developers of such a system assess whether the system has succeeded in producing a correct model or not?
- **Trust:** A system like OptiMUS should help domain experts understand the model that was selected and decide whether it meets their needs. How should the system communicate a model back to the human user to facilitate effective feedback about the deficiencies in the model, or trust in a model that is correct? How should the system integrate this feedback? How can the system developer assess the effectiveness of this feedback and improve the system to enable quicker and easier feedback cycles?
- **Ambiguity:** Most optimization problems specified in natural language do not correspond unambiguously to a single optimization model. Rather, the process of formalizing the objective and constraints reduces the vagueness so that it is possible to distinguish a solution from any other proposal. A system like OptiMUS should help decision-makers think through their goals (which may be vague) to choose a model formulation by clarifying implicit requirements, eliciting the relative importance of different objectives, and searching internal or external databases for relevant data.
- **Fast Solvers:** Rapid decisions require a fast solver: possibly a well-chosen solver for the problem, or possibly even a custom solver written for the problem. A system like OptiMUS should be able to choose between available heuristics and parameters for provably optimal solvers using all information available about the problem to determine a solution method that meets runtime and accuracy requirements, or even to write a custom solver when it is warranted. To choose or write the best solver, we expect that the system will likely make use of natural-language information about the problem setting (using an LLM) as well as structured information about problem parameters (using a graph neural network).
- **Larger datasets.** All of these research directions depend on the availability of appropriate datasets to develop, hone, and test new ideas. Most of these directions require a dataset of optimization problems consisting of a natural-language description coupled with problem data and solution code. While there are many excellent datasets for developing optimization algorithms, these are generally not associated with natural-language descriptions; and problems associated with natural language descriptions are generally either small (e.g., textbook problems) or lack associated problem data and solutions (e.g., published papers on applications of MIPs) as these elements are often proprietary. Semi-synthetic approaches to generating larger datasets are promising (Tang et al. 2024), but creating MIP problems that are large-scale, realistic, and feasible is a challenge. Curating a larger dataset of optimization problems could be a key enabler driving further research and development in this area.



- **Beyond MILPs.** OptiMUS currently uses gurobipy as a modeling language. We made this choice as the wealth of online examples allows LLMs to successfully model many problems in gurobipy. However, many optimization problems are best expressed in a different framework: for example, Concorde for traveling salesman problems (Cook et al. 2011), cvxpy for control problems (Diamond and Boyd 2016), or minizinc for constraint programming problems (Nethercote et al. 2007). Integrating these alternative modeling languages and identifying, for a given problem, which modeling language to use (or whether the problem is likely to be successfully modeled in the OptiMUS framework at all) is an important future direction.

## Acknowledgments

AA, WG, CL, and MU gratefully acknowledge support from the National Science Foundation (NSF) Award IIS-2233762, the Office of Naval Research (ONR) Awards N000142212825, N000142412306, and N000142312203, the Alfred P. Sloan Foundation, and from IBM Research as a founding member of Stanford Institute for Human-centered Artificial Intelligence (HAI). HB thanks the Aker Scholarship Foundation for financial support. ST gratefully acknowledges support from the Air Force Office of Scientific Research (AFOSR) under Grant FA9550-23-1-0251 and in part by the Office of Naval Research under Grant N00014-24-1-2164. This manuscript was improved by feedback from anonymous referees.

## References

- Jesper Aastrup and Herbert Kotzab. Forty years of out-of-stock research — and shelves are still empty. *The International Review of Retail, Distribution and Consumer Research*, 20(1):147–164, 2010.
- Tobias Achterberg. What’s new in gurobi 9.0. *Webinar Talk url: <https://www.gurobi.com/wp-content/uploads/2019/12/Gurobi-90-Overview-Webinar-Slides-1.pdf>*, 2019.
- Ali AhmadiTeshnizi, Wenzhi Gao, and Madeleine Udell. OptiMUS: Scalable optimization modeling with (MI)LP solvers and large language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 577–596. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/ahmaditeshnizi24a.html>.
- Alibaba Cloud. Alibaba cloud mindopt copilot, 2022. URL <https://opt.alibabacloud.com/chat>.
- Nicolás Astorga, Tennison Liu, Yuanzhang Xiao, and Mihaela van der Schaar. Autoformulation of mathematical optimization models using llms. *arXiv preprint arXiv:2411.01679*, 2024.
- EML Beale and John JH Forrest. Global optimization using special ordered sets. *Mathematical Programming*, 10: 52–69, 1976.
- Dimitris Bertsimas and John N Tsitsiklis. *Introduction to linear optimization*, volume 6. Athena scientific Belmont, MA, 1997.

- Robert E Bixby, John W Gregory, Irvin J Lustig, Roy E Marsten, and David F Shanno. Very large-scale linear programming: A case study in combining interior point and simplex methods. *Operations research*, 40(5):885–897, 1992.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR, 2022.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, March 2004.
- Kathy Charmaz. *Constructing grounded theory: A practical guide through qualitative analysis*. Sage, New York, NY, 2006.
- Hao Chen, Gonzalo E Constante-Flores, and Can Li. Diagnosing infeasible optimization problems using large language models. *arXiv preprint arXiv:2308.12923*, 2023.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language modeling with pathways, 2022.
- William J Cook, David L Applegate, Robert E Bixby, and Vasek Chvatal. *The traveling salesman problem: a computational study*. Princeton university press, 2011.
- CPLEX User’s Manual. IBM ILOG CPLEX optimization studio. *Version*, 12(1987-2018):1, 1987.
- Steven Diamond and Stephen Boyd. CVXPY: A python-embedded modeling language for convex optimization. *The Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context learning, 2024. URL <https://arxiv.org/abs/2301.00234>.
- Gerald Gamrath, Timo Berthold, Stefan Heinz, and Michael Winkler. *Structure-based primal heuristics for mixed integer programming*. Springer, 2016.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR, 2023.

- Ambros Gleixner, Gregor Hendel, Gerald Gamrath, Tobias Achterberg, Michael Bastubbe, Timo Berthold, Philipp Christophel, Kati Jarck, Thorsten Koch, Jeff Linderoth, et al. Miplib 2017: data-driven compilation of the 6th mixed-integer programming library. *Mathematical Programming Computation*, 13(3):443–490, 2021.
- Google. Google Gemini AI: Advancements and Performance, 2024. URL <https://blog.google/technology/ai/google-gemini-ai/#performance>. [Accessed: October 27, 2024].
- Gurobi Optimization. 2023 state of mathematical optimization report, 2023. URL <https://www.gurobi.com/resources/report-state-of-mathematical-optimization-2023/>.
- Jesse Holzer, Carleton Coffrin, Christopher DeMarco, Ray Duthu, Stephen Elbert, Brent Eldridge, Tarek Elgindy, Scott Greene, Nongchao Guo, Elaine Hale, Bernard Lesieutre, Terrence Mak, Colin McMillan, Hans Mittelmann, Hyungseon Oh, Richard O’Neill, Thomas Overbye, Bryan Palmintier, Farnaz Safdarian, Ahmad Tbaileh, Pascal Van Hentenryck, Arun Veeramany, and Jessica Wert. Grid optimization competition challenge 3 problem formulation. [https://gocompetition.energy.gov/sites/default/files/Challenge3\\_Problem\\_Formulation\\_20230126.pdf](https://gocompetition.energy.gov/sites/default/files/Challenge3_Problem_Formulation_20230126.pdf), 2023. Accessed: Access Date.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Caigao Jiang, Xiang Shu, Hong Qian, Xingyu Lu, Jun Zhou, Aimin Zhou, and Yang Yu. Llmopt: Learning to define and solve general optimization problems from scratch. *arXiv preprint arXiv:2410.13213*, 2024.
- Jakob Krarup and Peter Mark Pruzan. The simple plant location problem: Survey and synthesis. *European journal of operational research*, 12(36-81):41, 1983.
- Connor Lawless, Jakob Schoeffler, Lindy Le, Kael Rowan, Shilad Sen, Cristina St. Hill, Jina Suh, and Bahareh Sarrafzadeh. “I want it that way”: Enabling interactive decision support using large language models and constraint programming. *ACM Transactions on Interactive Intelligent Systems*, 14(3):1–33, 2024.
- Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length on the reasoning performance of large language models, 2024.
- Beibin Li, Konstantina Mellou, Bo Zhang, Jeevan Pathuri, and Ishai Menache. Large language models for supply chain optimization. *arXiv preprint arXiv:2307.03875*, 2023.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts, 2023.
- Steven H. Low. Convex relaxation of optimal power flow: A tutorial. In *2013 IREP Symposium Bulk Power System Dynamics and Control - IX Optimization, Security and Control of the Emerging Power Grid*, pages 1–15, 2013. doi: 10.1109/IREP.2013.6629391.
- Hongliang Lu, Zhonglin Xie, Yaoyu Wu, Can Ren, Yuxuan Chen, and Zaiwen Wen. Optmath: A scalable bidirectional data synthesis framework for optimization modeling. *arXiv preprint arXiv:2502.11102*, 2025.
- Zhi-Quan Luo, Wing-Kin Ma, Anthony Man-Cho So, Yinyu Ye, and Shuzhong Zhang. Semidefinite relaxation of quadratic optimization problems. *IEEE Signal Processing Magazine*, 27(3):20–34, 2010.

- Meta AI. Meta LLaMA 3 Announcement, 2024. URL <https://ai.meta.com/blog/meta-llama-3/>. [Accessed: October 27, 2024].
- Dritan Nace. Lecture notes in linear programming modeling, 2020. URL [https://www.hds.utc.fr/~dnace/dokuwiki/\\_media/fr/lp-modelling\\_upt\\_p2021.pdf](https://www.hds.utc.fr/~dnace/dokuwiki/_media/fr/lp-modelling_upt_p2021.pdf).
- Nicholas Nethercote, Peter J Stuckey, Ralph Becket, Sebastian Brand, Gregory J Duck, and Guido Tack. Minizinc: Towards a standard cp modelling language. In *International Conference on Principles and Practice of Constraint Programming*, pages 529–543. Springer, 2007.
- OpenAI. Gpt-4 technical report, 2023.
- Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B. Hashimoto. Proving test set contamination in black box language models, 2023. URL <https://arxiv.org/abs/2310.17623>.
- PapersOwl. Free Plagiarism Checker, 2024. URL <https://papersowl.com/free-plagiarism-checker>. [Accessed: Date of Access, e.g., October 27, 2024].
- Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. ART: Automatic multi-step reasoning and tool-use for large language models. *arXiv preprint arXiv:2303.09014*, 2023.
- . Ramamonjison et al. Augmenting operations research with auto-formulation of optimization models from problem descriptions. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 29–62, Abu Dhabi, UAE, December 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-industry.4>.
- . Ramamonjison et al. NL4Opt competition: Formulating optimization problems based on their natural language descriptions, 2023. URL <https://arxiv.org/abs/2303.08233>.
- Tomáš Roubíček. *Relaxation in optimization theory and variational calculus*, volume 4. Walter de Gruyter GmbH & Co KG, 2020.
- Soroush Saghaian, Garrett Austin, and Stephen J Traub. Operations research/management contributions to emergency department patient flow optimization: Review and research prospects. *IIE Transactions on Healthcare Systems Engineering*, 5(2):101–123, 2015.
- Rabia Shakoar, Mohammad Yusri Hassan, Abdur Raheem, and Yuan-Kang Wu. Wake effect modeling: A review of wind farm layout optimization using Jensen’s model. *Renewable and Sustainable Energy Reviews*, 58:1048–1059, 2016.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- Zhengyang Tang, Chenyu Huang, Xin Zheng, Shixi Hu, Zizhuo Wang, Dongdong Ge, and Benyou Wang. Orlm: Training large language models for optimization modeling, 2024. URL <https://arxiv.org/abs/2405.17743>.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Yuvraj Virk, Premkumar Devanbu, and Toufique Ahmed. Enhancing trust in llm-generated code summaries with calibrated confidence scores. *arXiv preprint arXiv:2404.19318*, 2024.
- Segev Wasserkrug, Leonard Boussiou, Dick den Hertog, Farzaneh Mirzazadeh, Ilker Birbil, Jannis Kurtz, and Donato Maragno. From large language models and optimization to decision optimization copilot: A research manifesto. *arXiv preprint arXiv:2402.16269*, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-Thought prompting elicits reasoning in large language models, 2023.
- H Paul Williams. *Model building in mathematical programming*. John Wiley & Sons, 2013.
- Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin Fu, Tao Zhong, Jia Zeng, Mingli Song, et al. Chain-of-Experts: When LLMs meet complex operations research problems. In *The Twelfth International Conference on Learning Representations*, 2023.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*, 2023.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers, 2023.
- Zhicheng Yang, Yiwei Wang, Yinya Huang, Zhijiang Guo, Wei Shi, Xiongwei Han, Liang Feng, Linqi Song, Xiaodan Liang, and Jing Tang. Optibench meets resocratic: Measure and improve llms for optimization modeling. *arXiv preprint arXiv:2407.09887*, 2024.
- Enjian Yao, Tong Liu, Tianwei Lu, and Yang Yang. Optimization of electric vehicle scheduling with multiple vehicle types in public transport. *Sustainable Cities and Society*, 52:101862, 2020.
- Fariba Zohrizadeh, Cedric Jozs, Ming Jin, Ramtin Madani, Javad Lavaei, and Somayeh Sojoudi. A survey on conic relaxations of optimal power flow problem. *European journal of operational research*, 287(2):391–409, 2020.

## Appendix

### A. Terms and Definitions

Table 5: Definitions of Terms Used in OptiMUS-0.3

Term: Clause	
Definition:	A constraint or objective within the optimization problem.
How We Use Them:	We refer to constraints and objectives as clauses to simplify terminology.
Example:	The constraint $x + y \leq 10$ is a clause.

**Term: Connection graph**

Definition:	A graph that records which variables and parameters appear in each constraint.
How We Use Them:	Used to ensure consistency of formulations and focus the LLM on relevant context.
Example:	Constraint $C1$ connects to variables $x$ and $y$ in the connection graph.

**Term: State**

Definition:	The collection of all parameters, clauses, variables, and background information managed and modified during the solution process.
How We Use Them:	Saved and updated in JSON format throughout the problem-solving steps.
Example:	The state includes parameters like demand, variables like production quantity, and clauses like capacity constraints.

**Term: Reflective prompts**

Definition:	Prompts designed to encourage the LLM to reflect on and correct its own mistakes.
How We Use Them:	Used to reduce modeling errors by having the LLM check and fix its outputs.
Example:	Asking "Are units the same for both sides of constraint $C$ ?"

**Term: Confidence-based user feedback**

Definition:	A method where the LLM assesses its confidence and, if low, requests help from the user or a stronger LLM.
How We Use Them:	To improve the accuracy of the model when the LLM is unsure about its outputs.
Example:	The LLM says "I am not confident about constraint $C$ " and asks the user to verify.

**Term: Parameters**

Definition:	Known quantities in the optimization problem, each with a symbol, shape, and definition.
How We Use Them:	Extracted from the problem description and used in clause formulations.
Example:	The cost per unit, denoted as $c$ , is a parameter.

**Term: Variables**

Definition:	Unknown quantities to be determined in the optimization problem, each with a symbol, shape, definition, and type.
How We Use Them:	Defined during clause formulation and used in constraints and objectives.
Example:	The production quantity $x$ is a variable.

**Term: Background**

Definition:	A short string that explains the real-world context of the problem.
How We Use Them:	Included in every prompt to improve common sense reasoning.
Example:	"This problem involves optimizing factory production."

**Term: Error Correction**

Definition:	Techniques used to mitigate the impact of LLM hallucinations and correct errors.
How We Use Them:	Using reflective prompts and confidence-based user feedback to improve reliability.

Example: Correcting a misidentified parameter by prompting "Is the value of  $P$  known or not?"

---

**Term: Formulate Clauses**

---

Definition: A process step where clauses are mathematically formulated, and variables and auxiliary constraints are defined.

How We Use Them: To generate the mathematical representation of constraints and objectives.

Example: Formulating "Total production must meet demand" as  $x \geq d$ .

---

**Term: Extract Parameters**

---

Definition: A process step where parameters are extracted from the problem description.

How We Use Them: To identify all known quantities needed for the optimization problem.

Example: Extracting "demand" as a parameter from the description.

---

**Term: Correct Errors**

---

Definition: A process step where errors in the parameters, clauses, and variables are corrected.

How We Use Them: To ensure accuracy of the extracted and formulated components.

Example: Correcting a parameter's shape if it was misidentified.

---

**Term: Extract Targets**

---

Definition: A process step where constraints, objective, and connection graph are extracted.

How We Use Them: To identify the clauses and their relationships to parameters and variables.

Example: Extracting the constraint "Production capacity cannot exceed limit."

---

**Term: Code Targets**

---

Definition: A process step where code snippets for all parameters, clauses, and variables are generated.

How We Use Them: To assemble these snippets into a single runnable code file.

Example: Generating code to define variables in a solver.

---

**Term: Assemble Code**

---

Definition: A process step where all code snippets are combined into a single code file.

How We Use Them: To produce executable code for solving the optimization problem.

Example: Combining parameter definitions, variable declarations, and constraints into one script.

---

**Term: Debug**

---

Definition: A process step where code is iteratively executed and errors are fixed until the desired results are achieved.

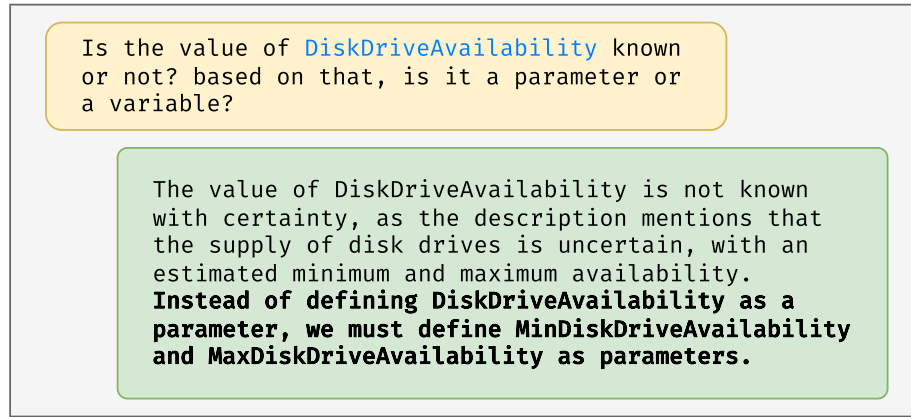
How We Use Them: To correct syntax or runtime errors in the generated code.

Example: Fixing a variable name mismatch that caused a runtime error.

---

## B. Reflexive Prompts

The reflective prompts used by OptiMUS-0.3 are as follows:

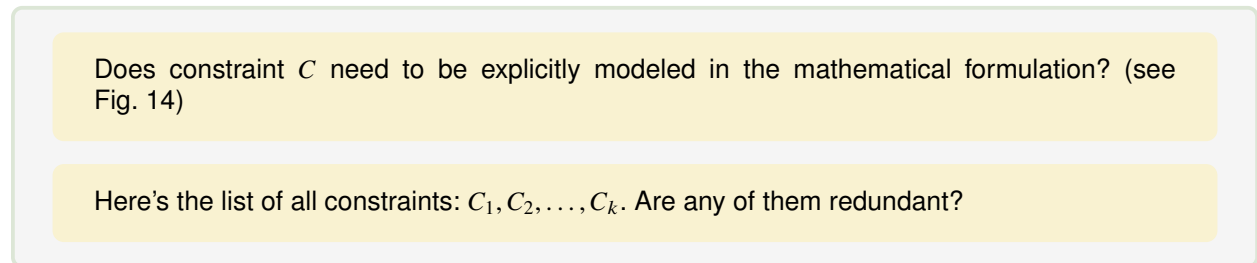


**Figure 10** OptiMUS-0.3 can fix parameter identification errors when prompted “*Is the value of  $P$  known or not?*”

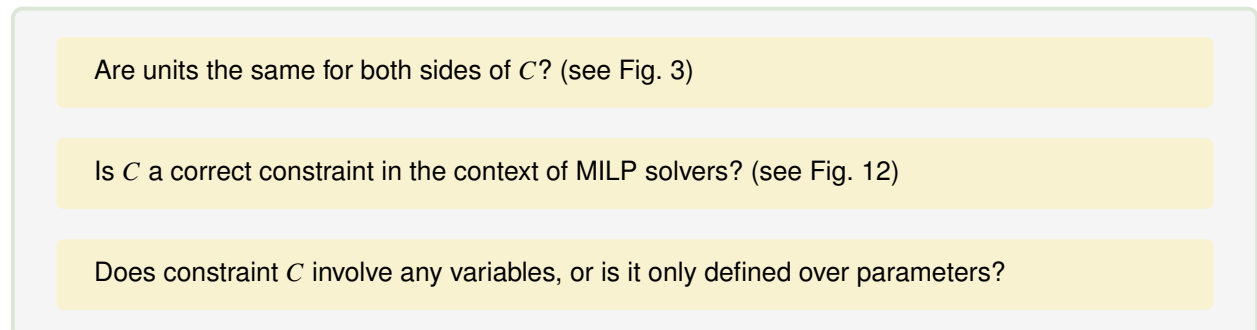
*Parameter Extraction.* The model often confuses parameters with variables, misidentifies parameter shapes or misses some parameters. We use the following reflective prompt to correct these errors:

- “*Is the value of  $P$  known or not? Based on that, is it a parameter or a variable?*” (see Fig. 10)

*Constraint Extraction.* During constraint extraction, the model sometimes extracts trivial or vague constraints and may regenerate a previously generated constraint in different language. We use the following reflective prompts to correct these errors:



*Clause Modeling.* Once the clauses are extracted, OptiMUS-0.3 models them, defining variables and auxiliary constraints if needed. Often, the model generates incorrect mathematical formulations for constraints. We use the following reflective prompts to detect and correct these errors:





Are units the same for both sides of this constraint?

$$(p_a + x_a) \cdot d_a \cdot (1 + e_a \cdot \frac{x_a}{p_a}) \leq m_a, \forall a \in A?$$

... Left-hand side (LHS):

- $(p_a + x_a)$  represents the new price for article  $a$ , which is in euros (€).
- $d_a$  represents the sales forecast (demand) for article  $a$  for the next twelve months at the current price, which is in units of the article.
- $(1 + e_a \cdot \frac{x_a}{p_a})$  is a unitless factor ...

Therefore, the unit of the left-hand side is: **euros (€) × units of the article**

Right-hand side (RHS):

- $m_a$  represents the maximum production volume for article  $a$ , which is in units of the article (e.g., number of shirts, pants, etc.).

The unit of the right-hand side is: **units of the article**

... this inconsistency suggests an error in the formulation of Constraint 5. To correct this, we should ... here is the corrected constraint:

$$d_a \cdot (1 + e_a \cdot \frac{x_a}{p_a}) \leq m_a, \forall a \in A$$

**Figure 11** OptiMUS-0.3 can fix its constraint modeling errors when prompted “Are units the same for both sides of  $C$ ?”

OptiMUS-0.3 uses similar reflective prompts to perform error checking for the optimization objective. For variables, we use the following prompt:

Is the value of  $V$  known or not? Based on that, is it a parameter or a variable? (see Fig. 13)

### C. Additional details on Structure Detection

Table 6 shows the distribution of special problem structures in the NLP4LP dataset.

Is  $p_i - \text{InitialPrice}_i \equiv 0 \pmod{1}$  a correct constraint in the context of MILP solvers?

The constraint  $p_i - \text{InitialPrice}_i \equiv 0 \pmod{1}$  is not a correct constraint in the context of MILP solvers. In MILP, constraints must be linear, and the  $\equiv$  symbol is not a standard operator in linear programming. The mod operator is not directly supported in MILP solvers.

To model the constraint that price changes must be in multiples of 1 euro, we can introduce an additional integer variable  $d_i$  to represent the price change in multiples of 1 euro:

$$p_i = \text{InitialPrice}_i + d_i$$

This constraint ensures that the new price  $p_i$  is the initial price plus a multiple of 1 euro, represented by the integer variable  $d_i$ .

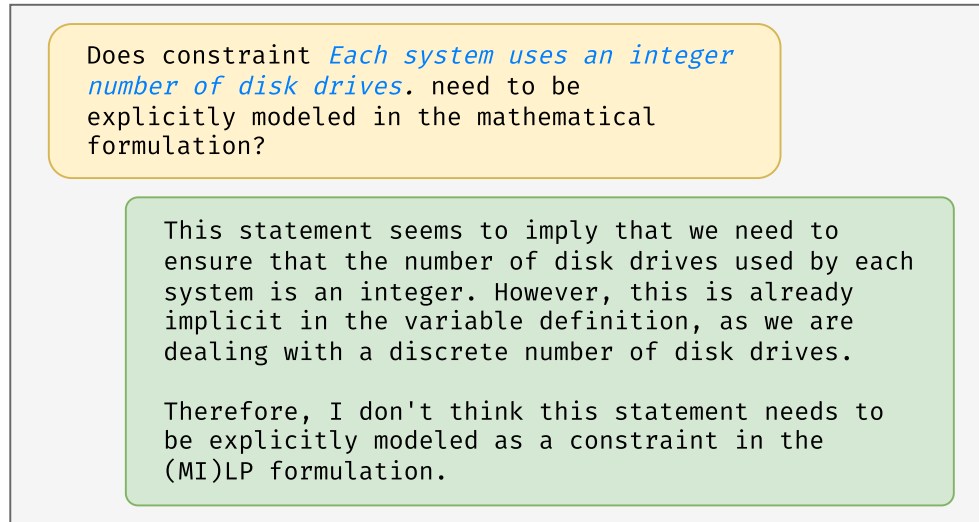
Figure 12 Through introspection, OptiMUS-0.3 can fix its errors in constraint modeling.

is the value of  $e_i$  known or not? based on that, is it a parameter or a variable?

Based on the problem statement, it seems that the value of  $e_i$  (price elasticity of article  $i$ ) is known. It is mentioned that "past elasticities" are given, implying that the elasticities are known constants.

Apologies for my mistake. In this case,  $e_i$  would be a parameter, not a variable. It would be a known value that is used in the calculation of the objective function, but it would not be a decision variable that is optimized by the solver. Here's the fixed formulation:

Figure 13 OptiMUS-0.3 can fix its parameter and variable identification errors when prompted "Is the value of  $P$  known or not?"



**Figure 14** OptiMUS-0.3 can fix its parameter and variable identification errors when prompted “*Is the value of  $P$  known or not?*”

**Table 6** Distribution of problem structures in NLP4LP

Structure	Number of instances
SOS	2
Indicator	6
Piecewise-linear	17

#### D. SCUC problem

**Description.** The Security Constrained Unit Commitment (SCUC) is a critical optimization problem in the operation of electrical power systems. Its goal is to schedule generation units in a cost-effective manner while ensuring reliable operation over a specified time period (typically 24 hours), considering the system’s physical and operational constraints. This appendix presents a detailed description of the SCUC problem, including its objectives and constraints:

**Objective.** The objective of the SCUC problem is to minimize the total operational cost of the power system. This cost generally includes:

- **Fuel Cost:** The cost associated with consuming fuel to generate electricity.
- **Start-up and Shut-down Cost:** Cost incurred from starting up or shutting down generation units.
- **Emission Cost:** Cost related to the emissions produced by the generating units, if applicable.

**Constraints.** The SCUC must satisfy several constraints to ensure the safe and reliable operation of the power system:

- **Power Balance Constraint:** Ensures that the total power generation meets the total demand plus system losses at every interval. Mathematically, this is expressed as the sum of the outputs of all online generators being equal to the sum of the demand and the transmission losses.
- **Minimum and Maximum Output Limits:** Each generator has a minimum and maximum generation capacity when it is online. The SCUC ensures that the output of each unit stays within these bounds.

### Example: Indicator Constraint

#### Definition

An indicator constraint over a binary variable  $z$  and a constraint  $C$  states that  $z = 1 \Rightarrow C$ .

#### Relevant description

A retail company wants to open new stores in a set of  $K$  potential locations. A store at location  $k$  must have a minimum staff of  $n_k$  employees (to ensure adequate customer service) and a minimum inventory level of  $l_k$  units (to avoid stockout). The maximum inventory level for store  $k$  is  $u_k$ , and there is a fixed cost associated with opening each store.

#### Relevant Constraints

If a store at location  $k$  is open, it must have at least  $n_k$  employees. Moreover, let the inventory level of store  $k$  be  $y_k$ .

#### Formulations

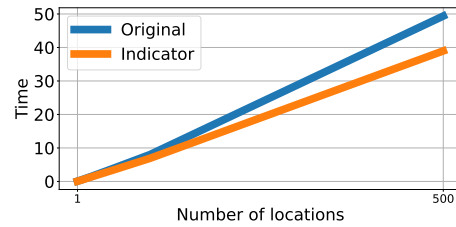
Original:

$$x_k \in \{0, 1\}, y_k \leq u_k x_k, \\ z_k \geq l_k - M(1 - x_k)$$

OptiMUS:

$$x_k \Rightarrow y_k \leq u_k \\ x_k \Rightarrow z_k \geq l_k$$

#### Performance plot



- **Ramp-up and Ramp-down Limits:** These limits specify the maximum rate at which a generator can increase or decrease its output. They are critical for handling load changes throughout the day.
- **Unit Start-up and Shut-down Constraints:** These constraints handle the logistics of turning units on or off. Start-up constraints may include minimum down times (the minimum time a unit must remain off before it can be restarted) and minimum up times (the minimum time a unit must remain on once started).
- **Minimum Up and Down Time Constraints:** Ensures that once a generator is turned on, it stays on for at least its minimum up time, and similarly, once turned off, it remains off for at least its minimum down time. These constraints are crucial for the mechanical integrity of the generation units.
- **Reserve Requirements:** The system operator must ensure that sufficient spinning and non-spinning reserves are available. These reserves are needed to handle sudden increases in demand or unexpected generator failures.

## E. Web Application

As discussed in 1, the OptiMUS project aims to 1) help optimization experts develop and maintain their models with ease and 2) empower domain experts to make better decisions faster and reduce costs. To further these goals, we designed and built a user-centric web app with an intuitive interface that enables users to interact seamlessly with the system and leverage the advantages of large language models (LLMs).

## Example: Special Ordered Set

## Definition

Given a set of variables  $\mathcal{X} = \{x_1, \dots, x_n\}$ ,

- $\{x_1, \dots, x_n\} \in \text{SOS}_1$  if at most one element of  $\mathcal{X}$  can be nonzero
- $\{x_1, \dots, x_n\} \in \text{SOS}_2$  if at most two elements of  $\mathcal{X}$  can be nonzero

## Relevant Description

An airline needs to assign  $N$  crew members to a set of  $K$  flights for a day. Each crew member can only be assigned to at most one flight at a time, and there must be a minimum rest period of 2 hours between consecutive flights to comply with aviation regulations and ensure crew well-being. Each crew member has a maximum number of flying hours per day, set at 8 hours, to comply with aviation regulations and prevent fatigue. There are minimum and maximum layover times between flights: minimum layover of 2 hours (to allow for crew rest and flight preparations) and maximum layover of 6 hours (to optimize crew utilization and reduce idle time). Assigning a crew member to a specific flight has some cost.

## Relevant Constraints

Each crew member can work on only one flight. Let  $x_{ij} \in \{0, 1\}$  denote whether crew  $i$  works on flight  $j$  and  $T_{ij}$  denote the working time of crew  $i$  on flight  $j$ .

## Formulations

Original:

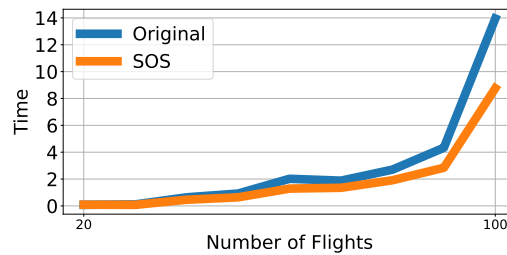
$$\sum_{j=1}^K x_{ij} \leq 1, \quad \forall i \in [N],$$

$$T_{ij} \leq 24x_{ij}, \quad \forall i \in [N], j \in [K]$$

OptiMUS:

$$\{T_{i1}, \dots, T_{iK}\} \in \text{SOS}_1$$

## Performance plot



We had the following key characteristics in mind when designing the web app:

- **User-centered design:** The design centers around the user's modeling flow rather than requiring users to adapt their workflow to the system.
- **Observability and intervention:** Given that LLM-based systems often are prone to errors, user supervision is critical. The interface is designed to facilitate this in every step.

The process includes several key steps:

1. **Parameter Extraction:** Users can input text in any format, and the model automatically extracts parameters. Users have the flexibility to edit and update parameters if necessary (Fig. 15).
2. **Clause Detection:** The model detects the objective and constraints given the provided information (Fig. 16).
3. **Clause Formulation:** Each clause is formulated by the model (Fig. 17).

**OptiMUS**

1 Description  
2 **Parameters**  
3 Constraints & Objective  
4 Mathematical Formulation  
5 Coding  
6 Data Processing  
7 Testing

**Formatted Description**

Consider a production problem. Given a number of products  $\backslash\text{param}\{P\}$ , each product is produced at a specific rate  $\backslash\text{param}\{\text{ProductionRate}\}$  (in tons per hour). There are  $\backslash\text{param}\{\text{HoursAvailable}\}$  hours available in a week. A ton of each product results in a known profit  $\backslash\text{param}\{\text{ProfitPerTon}\}$ . For each product, there is a lower limit  $\backslash\text{param}\{\text{LowerLimit}\}$  and an upper limit  $\backslash\text{param}\{\text{UpperLimit}\}$  on the tons of that product sold in a week. The problem aims to maximize the total profit from selling all products. The total number of hours used by all products may not exceed  $\backslash\text{param}\{\text{HoursAvailable}\}$ . How to decide the tons of each product to be produce?

**Parameters**

Symbol	Shape	Definition	Action
LowerLimit	[P]	Lower limit on the tons of each product sold in a week	Remove
P	[I]	Number of different products	Remove
ProfitPerTon	[P]	Profit obtained by selling a ton of each product	Remove
HoursAvailable	[I]	Number of hours available in a week	Remove
UpperLimit	[P]	Upper limit on the tons of each product sold in a week	Remove
ProductionRate	[P]	Production rate of each product in tons per hour	Remove

**Reset**  
[Have Feedback?](#)  
Made with ❤ at Udell Lab

**Next →**

**Figure 15 Parameter Identification:** The parameters are automatically extracted, and users can add, remove, or modify them as needed.

**OptiMUS**

1 Description  
2 Parameters  
3 **Constraints & Objective**  
4 Mathematical Formulation  
5 Coding  
6 Data Processing  
7 Testing

**Extract Constraints and Objective**

**Objective**

The company aims to maximize its total profit from selling all products

**Background**

A company produces a variety of products, each with specific production rates, profit margins, and time constraints within a given week.

**Constraints**

Description	Action
Each product can only be produced in non-negative quantities	Remove
The production volume for each product must be an integer value if products cannot be fractionally produced	Remove
Each product has a minimum production limit of LowerLimit tons per week	Remove
Each product has a maximum production limit of UpperLimit tons per week	Remove
The total number of production hours for all products must not exceed HoursAvailable hours per week	Remove
The production of each product is constrained by its respective ProductionRate in tons per hour	Remove

**Reset**  
[Have Feedback?](#)  
Made with ❤ at Udell Lab

**Next →**

**Figure 16 Clause Detection:** In this example, one identified constraint is incorrectly marked as a constraint (The production limit is automatically enforced and does not need to be an explicit constraint). The interface allows the user to easily remove it.

4. **Code Generation:** The model generates code for each formulated clause (Fig. 18).
5. **Data File Structure Inference:** Based on the parameter information, a data file structure is inferred. Users can upload data in this format or use randomly generated data to continue (Fig. 19).

**OptiMUS**

1 Description  
2 Parameters  
3 Constraints & Objective  
4 Mathematical Formulation  
5 Coding  
6 Data Processing  
7 Testing

**Formulate All**

**Objective**

Maximize the total profit from selling all products **Formulate**

$$\max \sum_{p=1}^P (\text{ProfitPerTon}_p \times \text{TonsSold}_p)$$

**Constraints**

Non-negative production constraint for each product **Formulate**

$$\text{TonsSold}_p \geq \text{LowerLimit}_p \quad \forall p \in \{1, 2, \dots, P\}$$

Every product must meet a minimum production **Formulate**

$$\text{TonsSold}_p \geq \text{LowerLimit}_p \quad \forall p \in \{1, 2, \dots, P\}$$

Each product must not be produced in quantities **Formulate**

$$0 \leq \text{TonsSold}_p \leq \text{UpperLimit}_p \quad \forall p \in \{1, 2, \dots, P\}$$

**Reset**  
Have Feedback?  
Made with ❤ at Udell Lab

**Figure 17** Clause Formulation: An error is identified where one constraint is incorrect, allowing the user to correct it.

**OptiMUS**

1 Description  
2 Parameters  
3 Constraints & Objective  
4 Mathematical Formulation  
5 Coding  
6 Data Processing  
7 Testing

**Objective**

$$\max \sum_{p=1}^P (\text{ProfitPerTon}_p \times \text{TonsSold}_p)$$

**Generate Code**

```
1 # Set objective
2 model.setObjective(gp.quicksum(ProfitPerTon[p] *
  TonsSold[p] for p in range(P)), gp.GRE_MAXIMIZE)
```

**Constraints**

$\text{TonsSold}_p \geq \text{LowerLimit}_p \quad \forall p \in \{1, 2, \dots, P\}$  **Generate Code**

```
1 # Add non-negative production constraints for each product
2 for p in range(P):
3     model.addConstr(TonsSold[p] >= LowerLimit[p], name
      =f"non_neg_prod_constr_{p}")
```

$\text{TonsSold}_p \geq \text{LowerLimit}_p \quad \forall p \in \{1, 2, \dots, P\}$  **Generate Code**

```
1 # Add constraints to ensure every product meets the minimum
  production limit per week
2 for p in range(P):
3     model.addConstr(TonsSold[p] >= LowerLimit[p], name
      =f"min_prod_limit_{p}")
```

$0 \leq \text{TonsSold}_p \leq \text{UpperLimit}_p \quad \forall p \in \{1, 2, \dots, P\}$  **Generate Code**

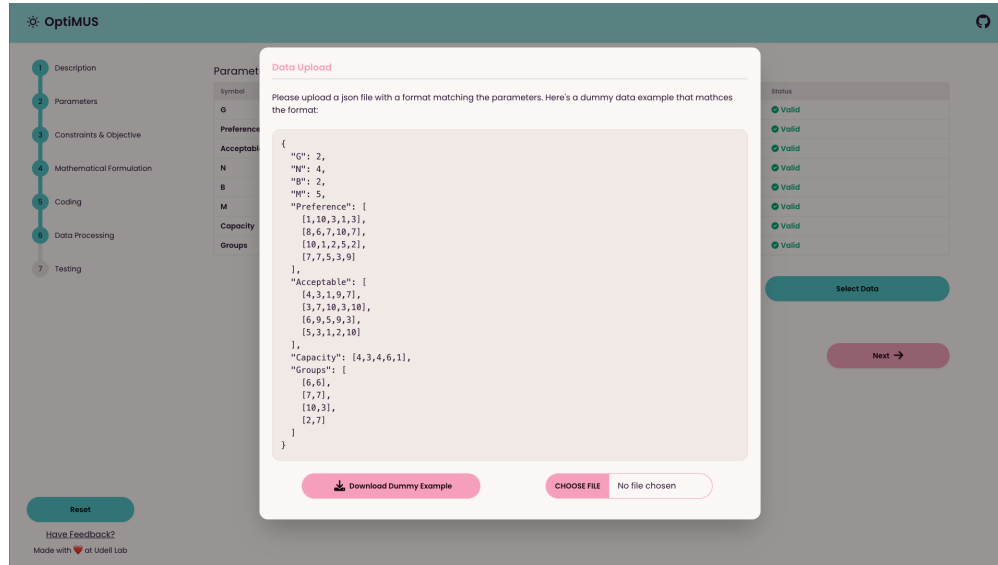
```
1 # Add upper limit constraints for products sold in a week
2 for p in range(P):
3     model.addConstr(TonsSold[p] <= UpperLimit[p], name
      =f"upper_limit_{p}")
```

**Reset**  
Have Feedback?  
Made with ❤ at Udell Lab

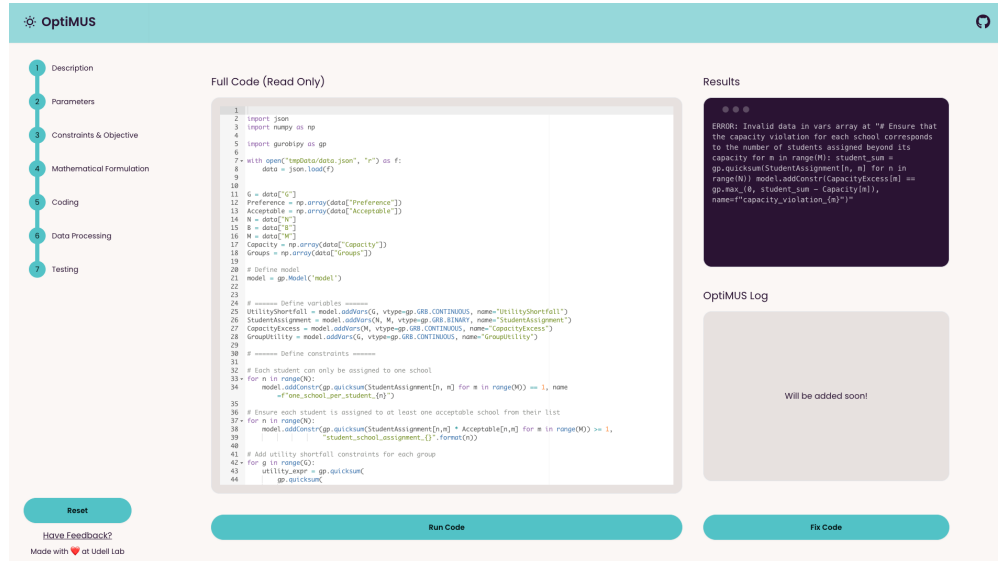
**Figure 18** Code Generation for each clause based on the formulated constraints and objectives.

6. **Code Editing and Testing:** Users can edit the generated code, run it on the dataset, and debug it with assistance from the LLM (Figs. 20 and 21).

By allowing users to observe and provide input throughout each step, the web app brings significant speed and convenience to the modeling process while minimizing error risks.



**Figure 19** Data File Structure: Inferred from the extracted parameters, with options for data upload or random generation.



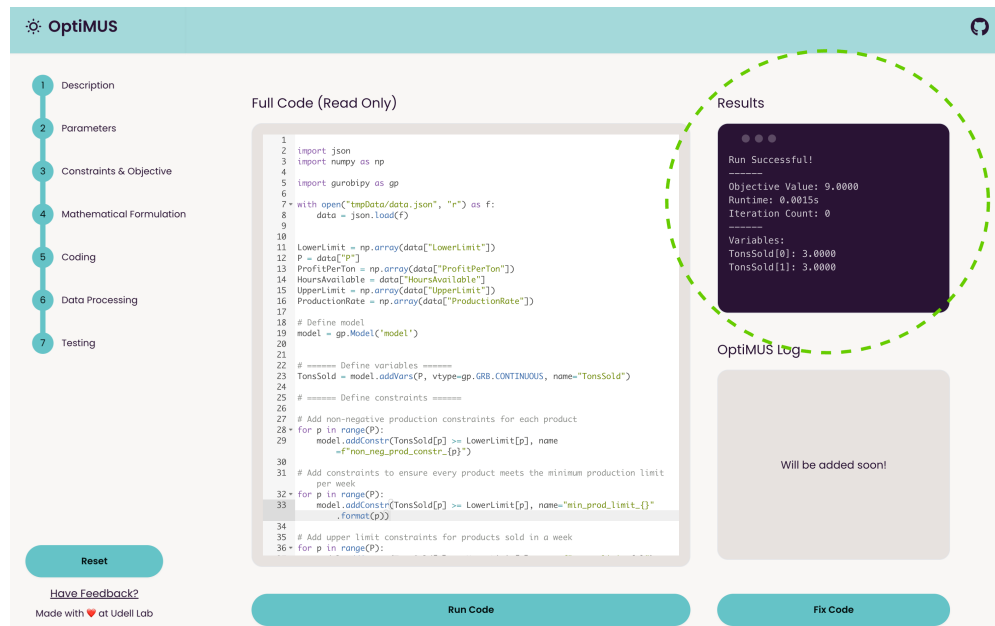
**Figure 20** Testing Phase: The code is debugged and revised as errors are detected.

## F. Dataset Figures

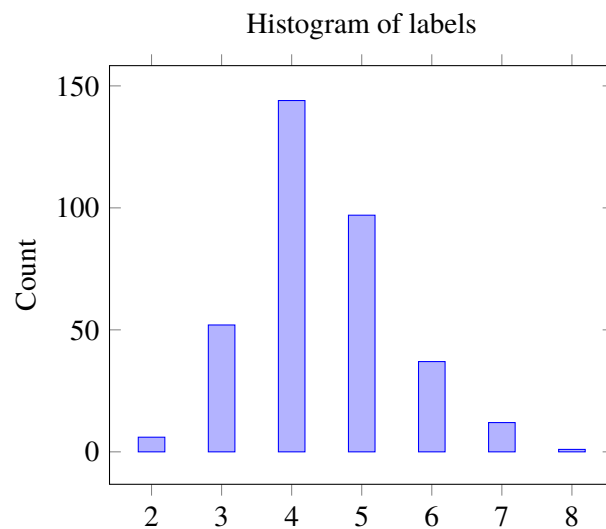
Fig. 23 and Fig. 24 represent the number of instances in the dataset per industry sectors and operational areas. Resource allocation is the most common area and manufacturing is the most common sector. Figure 22 shows the distribution of instances by the number of associated labels with the problem. Most instances have around 4 labels, with some having as few as 2 and as many as 8

For plagiarism tests we use Papers Owl online tool PapersOwl (2024). See Fig. 25 for an example and Section F for a histogram of originality scores for all instances.

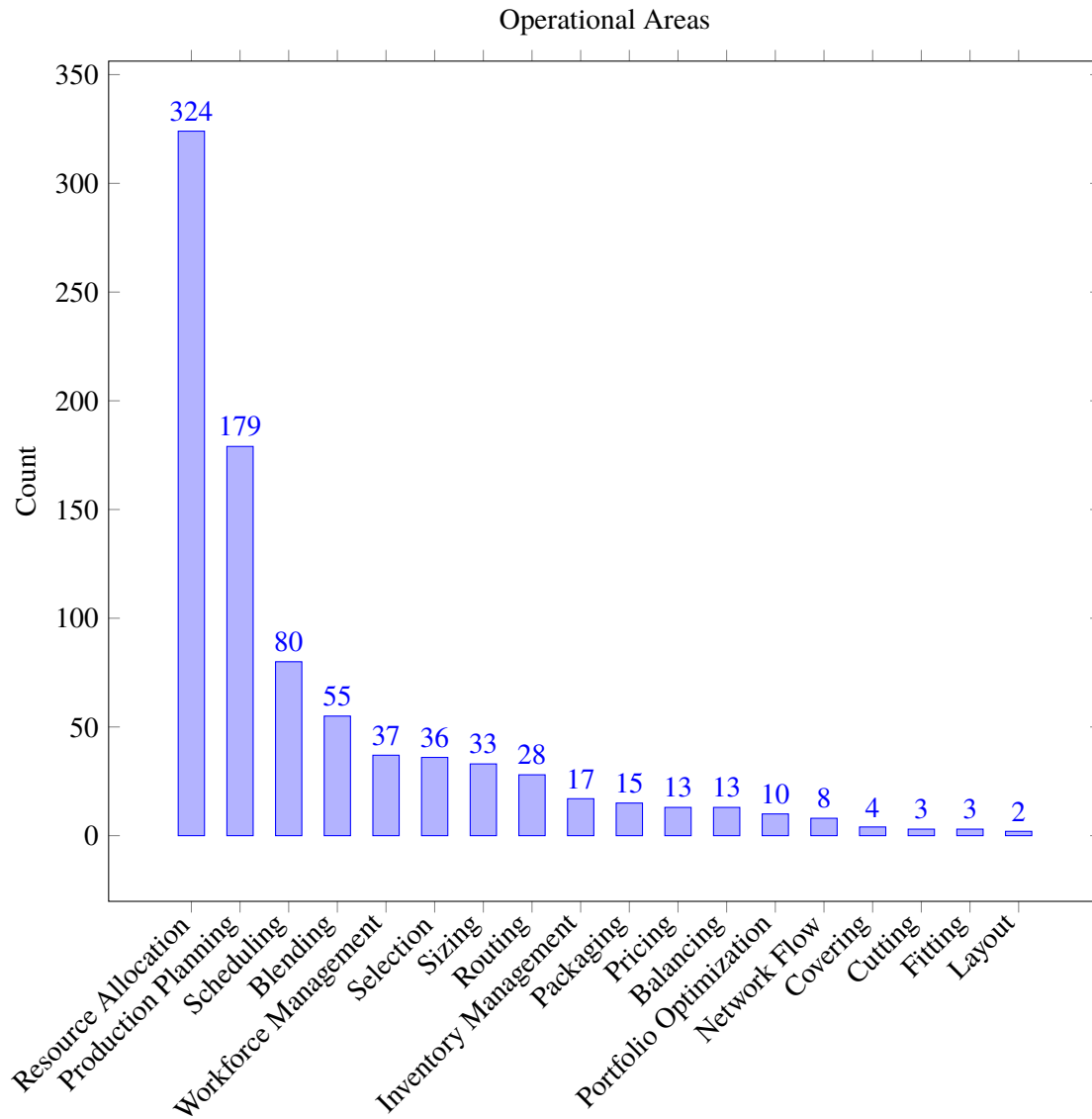




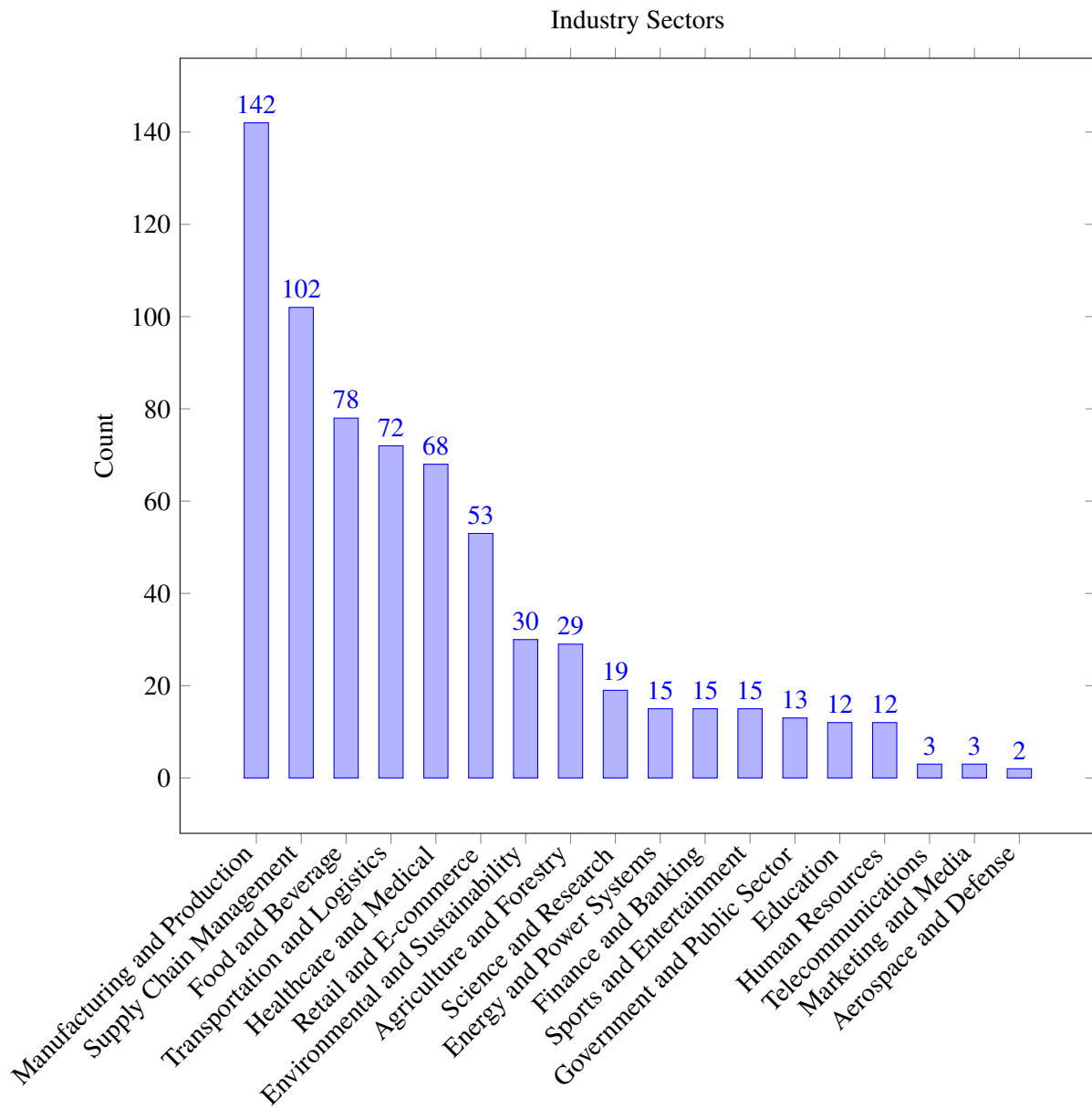
**Figure 21** Final Output: After debugging, the correct code is executed and the solution is displayed.



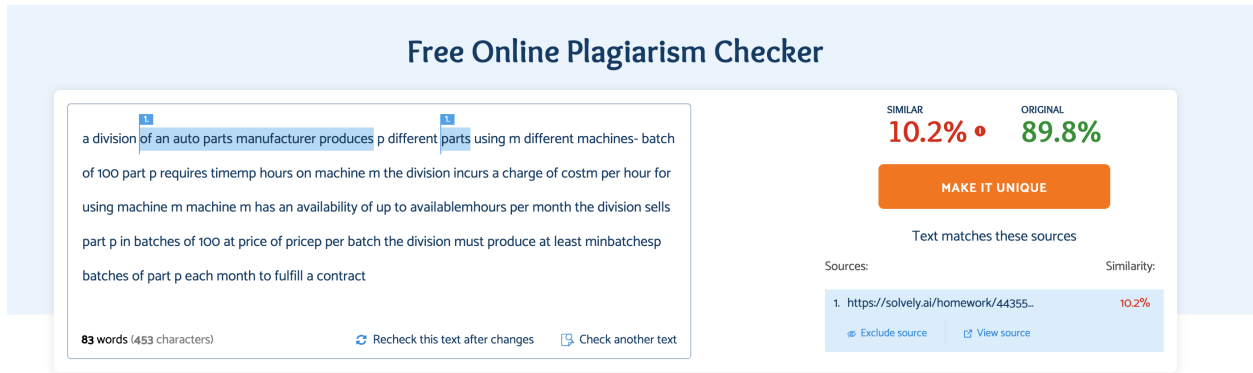
**Figure 22** Most problems have around 4 labels.



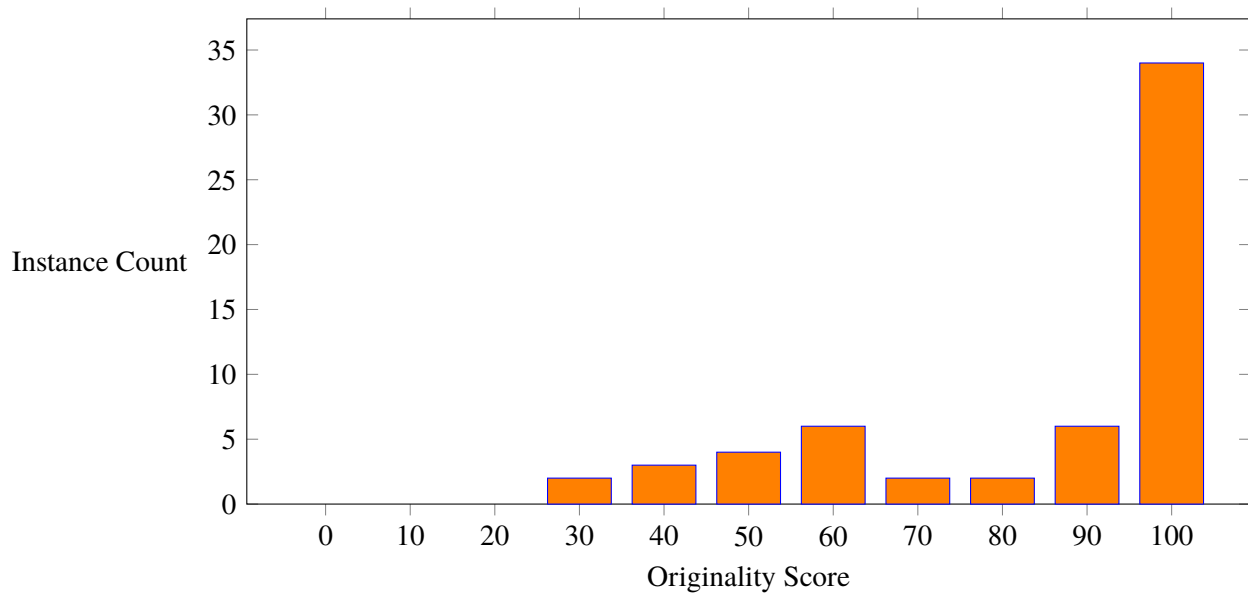
**Figure 23** NLP4LP covers common optimization problem types



**Figure 24** NLP4LP covers important optimization problem domains



**Figure 25** We used plagiarism detection tools to ensure that instances are novel and publicly available on the internet. Most of the similarities detected are at word level (e.g. from a news article about manufacturing, or from a different optimization problem in the same domain).



**Figure 26** Originality of the NLP4LP dataset. The instances in our dataset represent a really high originality score overall, minimizing the chances of them being using in internet-scraped training corpus used for training LLMs.