

# DECISION INFORMATION MEETS LARGE LANGUAGE MODELS: THE FUTURE OF EXPLAINABLE OPERATIONS RESEARCH

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Operations Research (OR) is vital for decision-making in many industries. While recent OR methods have seen significant improvements in automation and efficiency through integrating Large Language Models (LLMs), they still struggle to produce meaningful explanations. This lack of clarity raises concerns about transparency and trustworthiness in OR applications. To address these challenges, we propose a comprehensive framework, Explainable Operations Research (EOR), emphasizing actionable and understandable explanations accompanying optimization. The core of EOR is the concept of *Decision Information*, which emerges from what-if analysis and focuses on evaluating the impact of complex constraints (or parameters) changes on decision-making. Specifically, we utilize bipartite graphs to quantify the changes in the OR model and adopt LLMs to improve the explanation capabilities. Additionally, we introduce the first industrial benchmark to rigorously evaluate the effectiveness of explanations and analyses in OR, establishing a new standard for transparency and clarity in the field.

## 1 INTRODUCTION

Operations Research (OR) has a long history of optimizing complex decision-making processes, such as in logistics, finance, investment, transportation, and healthcare, etc., where even small improvements can lead to significant operational profits. As these optimization algorithms increasingly contribute to daily life, it is essential to ensure their trustworthiness and reliability through explanations, which build user confidence (Faulhaber et al., 2021). Governments are also responding to this need by enacting laws like the General Data Protection Regulation (GDPR) of the European Union (Goodman & Flaxman, 2017), emphasize the “right to explanation” for algorithmic decisions (Selbst & Powles, 2018) in automated systems.

In recent years, Large Language Models (LLMs) have emerged as powerful tools in the OR domain, offering new opportunities to automate and enhance the modeling process. Current research of LLMs in OR, such as works (Xiao et al., 2023; AhmadiTeshnizi et al., 2024; Tang et al., 2024; Zhang et al., 2024), explore the potential to streamline the formulation and solutions of complex OR problems. However, LLMs in OR have primarily focused on improving efficiency and accuracy by generating codes for external solvers to obtain OR solutions, with less attention to enhancing solution explainability, especially in real-time collaborative automated systems.

Meanwhile, several studies (Čyras et al., 2019; Li et al., 2023; Erwig & Kumar, 2024; De Bock et al., 2024) have explored explainable optimization related to OR, but there are still limitations. For example, (Erwig & Kumar, 2024) focuses specifically on combinatorial optimization problems, which limits its applicability across the broader OR landscape and fails to leverage the advantages of LLMs in real-time modeling and explanations. Another work, OptiGuide (Li et al., 2023), emphasizes what-if analysis, which, while useful for specific easy scenarios, lacks the robust modeling capability to address more complex cases like deleting or combining constraints. For example, if a warehouse closes, the OR model must remove the related storage capacity constraint and adjust the distribution network. Current methods struggle to achieve this level of flexibility, yet such adaptability is crucial for accurately reflecting real-world changes. Most critically, the explanations these methods provide are often superficial, merely summarizing the outcomes without exploring the un-

derlying reasons behind the results, thus lacking the quantitative analysis, depth, and clarity required to fully understand and trust the decision-making process.

Given the limitations of existing approaches to explainable optimization, we are motivated to develop a more comprehensive framework, EOR, for explaining OR models. As shown in Figure 1, our framework addresses the critical need for transparency in OR by shifting from purely modeling a problem (a natural language description) to providing clear, actionable explanations for a user query, such as “What if transportation costs increase by 15%?”. First, we formulate the problem of explainable OR within the context of LLMs. This formulation is essential for laying a foundation for future research in this emerging area. Second, our framework emphasizes two critical types of explanations, 1) *Explanation of correctness*: the reasons for code updates during the modeling process, and 2) *Explanation of the Results*: the rationale for generating specific solutions. Unlike traditional methods that provide superficial explanations and analyses, our approach incorporates more sophisticated what-if analysis, quantifying the effects of changes prompted by user queries and providing deeper insights into the decision-making process. A significant component of our contribution is the introduction of the concept of “Decision Information”. By leveraging bipartite graphs, we measure the importance of different decision factors in response to user queries, particularly when complex changes in constraints arise. This approach enhances both the modeling capabilities and the explanatory power of OR models within LLMs. Our dual focus ensures that the framework not only yields accurate optimal solutions but also effectively communicates the underlying rationale for these solutions. Finally, recognizing the need for standardized evaluation of explainable OR methods, we design a new industrial benchmark specifically tailored to assess the effectiveness of explanations in OR. This benchmark fills the gap of current approaches and sets a new standard for evaluating the transparency and comprehensibility of OR models.

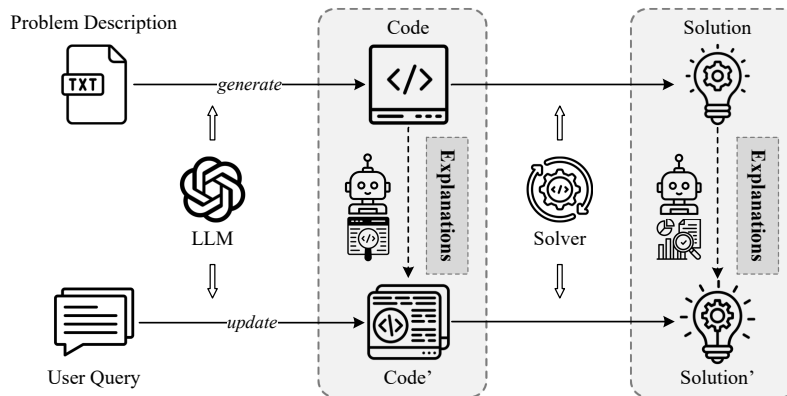


Figure 1: The framework of EOR.

**Contributions.** 1) We formulate the problem of the explainable OR problem within the context of LLMs, laying a foundation for future research in this area. 2) We introduce the concept of “Decision Information” and utilize bipartite graphs in conjunction with LLMs to quantify its importance in response to user queries, enhancing both the modeling capabilities and the explanation of complex what-if analysis within OR. 3) We develop a new benchmark specifically designed to evaluate the effectiveness of explanations in OR, setting a new standard for explainability in the field.

## 2 RELATED WORK

### 2.1 LLMs FOR OR

LLMs show great promise for OR, offering innovative approaches to optimize and automate modeling processes (Xiao et al., 2023; AhmadiTeshnizi et al., 2024; Tang et al., 2024; Zhang et al., 2024; Huang et al., 2024; Mostajabdaveh et al., 2024). Although LLMs have shown potential in various OR tasks, their application has primarily focused on automating modeling processes and enhancing computational efficiency. In contrast, our approach distinguishes itself by using LLMs to provide detailed, context-aware explanations of OR solutions, addressing the gap in explainability.

## 2.2 EXPLANATIONS FOR OR

Explanation in OR is essential for clarity and transparency, helping various stakeholders understand complex decision-making processes. Despite its significance, there is a notable lack of research, with only a few related works, such as (Thuy & Benoit, 2024; Erwig & Kumar, 2024; De Bock et al., 2024) and OptiGuide (Li et al., 2023), which only focus on the easy what-if analysis. Our approach focuses on more complex what-if analysis, seeking a broad range of methods that can uniformly quantify constraint changes or parameter changes (sensitivity analysis). We leverage LLMs to embed detailed, context-aware explanations directly within OR solutions, filling this gap. The comparison between EOR and OptiGuide is shown in Table 1. More details about the differences between the concepts are discussed in Appendix A.1.

Table 1: The comparison of OptiGuide and EOR.

	What-if Analysis		Decision Information Analysis	Sensitivity Analysis
	Easy	Complex		
OptiGuide	✓	✗	✗	✗
EOR	✓	✓	✓	✓

## 3 METHODOLOGY

### 3.1 PROBLEM FORMULATION

Given an OR problem  $\hat{p}$ , along with a user query  $\hat{q}$  related to the problem, our goal is to utilize an LLM to generate comprehensive explanations of solutions for the queries in real time. The LLM will provide two types of explanations: 1) *Attribution Explanation*, which outlines the general attributes and structure of the problem, and 2) *Justification Explanation*, which elucidates the correctness and derivation of the solutions. The mathematical formulation is as follows:

**Input:** The origin problem description  $\hat{p} = \langle \hat{d}, \hat{o}, \hat{c} \rangle$ , and a user query  $\hat{q} = \langle \hat{o}', \hat{c}' \rangle$ . The updated problem description, incorporating the user query, is denoted as  $\hat{p}' = \langle \hat{d}, \hat{o}', \hat{c}' \rangle$ . Here,  $\hat{d}$  represents the set of decision variables,  $\hat{o}$  and  $\hat{o}'$  are the objective functions to be maximized or minimized,  $\hat{c}$  and  $\hat{c}'$  denote the original and modified constraints in  $\hat{p}$  and  $\hat{p}'$ , respectively, that the decision variables must satisfy. In our setting, we assume the decision variables remain unchanged.

**Output:** We denote the problem solutions for  $\hat{p}$  and  $\hat{p}'$  as  $\hat{s}$  and  $\hat{s}'$ , respectively. The output comprises two types of explanations, 1) *Attribution Explanation*: A detailed description of the elements  $\hat{d}$ ,  $\hat{o}$ ,  $\hat{o}'$ ,  $\hat{c}$ ,  $\hat{c}'$ ,  $\hat{s}$ , and  $\hat{s}'$  within the context of the problem. 2) *Justification Explanation*: A rationale for the correctness of  $\hat{s}$  and  $\hat{s}'$ , and clarifies how  $\hat{s}'$  is derived from  $\hat{s}$ .

### 3.2 THE EOR FRAMEWORK

Our proposed framework, EOR, is an end-to-end solution designed to enhance OR model transparency using LLMs. Unlike current methods that provide limited and shallow explanations, primarily the form of attribution explanations, our framework emphasizes delivering clear, actionable insights for diverse stakeholders. As illustrated in Figure 1, we focus on the justification explanations. We will offer two critical explanations: 1) justifications for code updates during the modeling process, and 2) the rationale behind specific solutions. By doing so, we strive to make OR solutions more accessible, understandable, and applicable to a broader audience, thereby improving decision-making quality and fostering user trust.

#### 3.2.1 THE WORKFLOW OF EOR

As shown in Figure 2, EOR framework comprises three key agents: Commander, Writer, and Safe-guard, each serving a distinct role to ensure an efficient, accurate, and secure optimization process.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

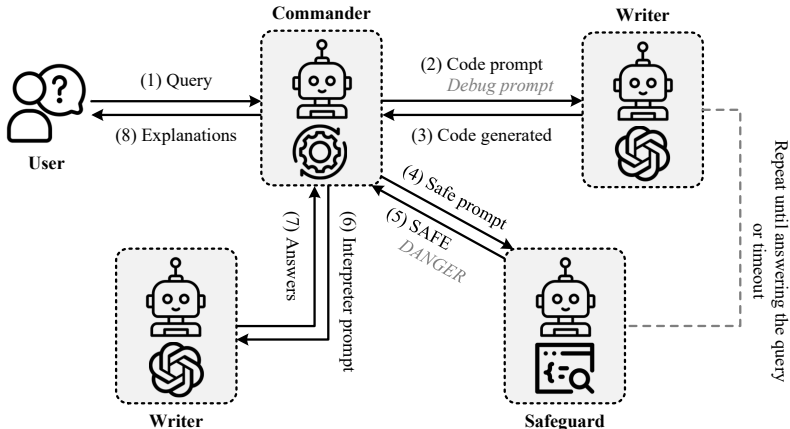


Figure 2: The overall workflow of EOR.

**Commander Agent:** The Commander acts as the central hub or “data bus” in the system, responsible for receiving user queries and managing data flow between the agents. When an end-user submits a new user query, the Commander first interprets the query’s context, identifies the intent, and then forwards the message to the appropriate agents.

**Writer Agent:** Upon receiving the processed query from the Commander, the Writer initially assumes the role of analyzing and modifying the code. Based on the query’s requirements, the Writer determines whether to add, delete, or update specific constraints and parameters. By leveraging LLMs, the Writer guarantees that the generated code accurately reflects the intended changes. Subsequently, the updated code is sent to the Safeguard for verification. Once the Safeguard provides a SAFE confirmation, the Writer transitions into an interpreter role, generating detailed explanations for the modifications and the rationale behind the decision-making process.

**Safeguard Agent:** The Safeguard is responsible for ensuring the safety and correctness of the generated code. It conducts thorough checks to verify whether the code adheres to predefined safety standards and is free from logical or syntactical errors that could compromise the optimization process. If the code passes these safety checks, the Safeguard approves it for execution; otherwise, it triggers a debugging process where the Writer regenerates and corrects the code as necessary.

**The overall workflow:** The EOR workflow starts when a user submits a query to the Commander (1), who relays it to the Writer with a code prompt (2). The Writer analyzes the query, determining whether to add, delete, or update code, and returns the updated code to the Commander (3). This step ensures that the generated code aligns with the updated problem requirements. The Commander then sends the code to the Safeguard for verification (4). If the Safeguard determines the code is safe (5), the process moves to (6), where the Commander sends an interpreter prompt to the Writer. If the code is deemed dangerous, the process loops back to (2), with the Commander sending a debug code prompt to the Writer. Once the code is safe, the Writer generates answers about explanations for the modifications and results (7) and sends them to the Commander. Finally, the Commander sends these explanations to the user (8). This iterative process, which repeats until a satisfactory answer or timeout, ensures robust, explainable solutions tailored to the user’s query. The detailed design of the prompt template is presented in Appendix A.2.

### 3.2.2 JUSTIFICATION EXPLANATION GENERATION

In our framework, explanations are generated to ensure transparency and trustworthiness in the decision-making process. These explanations are divided into two main categories:

**Explanation of Correctness:** This type of explanation serves to validate the code modifications introduced by the Writer, offering a detailed rationale for the changes made. It clarifies the necessity of these modifications in addressing the problem’s requirements and ensures that they adhere to safety standards and logical constraints. Through this process, the reliability of the generated code is substantiated, thereby enhancing the accuracy of the model prior to execution.

**Explanation of the Results:** Once the code is executed and the results are obtained, this type of explanation focuses on interpreting the outcome. It breaks down the results into understandable terms, illustrating the impact of the code changes on the final solution. This explanation connects the modifications to their direct effects, providing users with a clear understanding of how the new solution addresses their initial query and any resulting trade-offs or benefits.

Before formalizing the concept of “Decision Information”, consider the following user query from an OR problem in flight operations: *How should the aircraft configuration be adjusted if the company limits Type A aircraft to 15 and Type B aircraft to 30?* In this context, “Decision Information” refers to the new constraints, limiting the number of Type A and Type B aircraft to 15 and 30, respectively, which directly modify the optimization model. These changes reshape the solution space, requiring adjustments to meet demand within the newly imposed limits. Thus, “Decision Information” captures the key query elements that drive changes in the problem.

**Definition 1 Decision Information:** *Decision Information encompasses the parameters and constraints specified in a user’s query for an OR problem, capturing the essential details needed to reconfigure the problem according to the user’s intent.*

Following this definition, we turn to a quantitative evaluation of “Decision Information” to assess its impact on decision-making processes. However, existing methods lack a measure for assessing changes in decision information caused by constraint modifications, focusing only on sensitivity analysis in parameter changes. Inspired by (Xing et al., 2024), we convert both the updated and original programs into a standardized format and then calculate their differences to determine the importance of information changes. This process can be outlined in three steps:

**Conversion to Linear Programs (LPs) in General Form:** Both the updated code (resulting from user queries) and the original code are first parsed and translated into a standardized LP format, that is widely adopted by various LP solvers, including CPLEX (Cplex, 2009), Gurobi (Bixby, 2007), COPT (Ge et al., 2023) and so on. This conversion includes expressing all decision variables, objective functions, and constraints uniformly to allow a direct comparison. The general LP form captures the essence of both the initial and modified decision scenarios, providing a clear basis for analyzing how changes in input data or constraints affect the outcome.

Formally, an LP with  $n$  variables and  $m$  constraints can be represented as:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & \mathbf{l}^s \leq \mathbf{A}\mathbf{x} \leq \mathbf{u}^s \\ & \mathbf{l}^x \leq \mathbf{x} \leq \mathbf{u}^x, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is the constraint matrix,  $\mathbf{c} \in \mathbb{R}^n$  is the cost vector,  $\mathbf{x} \in \mathbb{R}^n$  is the decision variables,  $\mathbf{l}^x \in \mathbb{R}^n$  and  $\mathbf{u}^x \in \mathbb{R}^n$  are lower/upper bound of the decision variables  $\mathbf{x}$ , and  $\mathbf{l}^s \in \mathbb{R}^n$  and  $\mathbf{u}^s \in \mathbb{R}^n$  are lower/upper bound of the constraint.

**Graph Representation Conversion:** The standardized LPs can be transformed into bipartite graphs, which consist of two distinct sets of nodes: decision variables and constraints (Gasse et al., 2019; Fan et al., 2023; Xing et al., 2024). In this representation, edges between nodes signify dependencies or relationships, with edge weights indicating the strength or nature of these connections. We define the bipartite graph as  $\mathcal{G} = (S \cup X, E)$ , where  $S = \{s_i | i \in [m]\}$  represents the set of constraint nodes,  $X = \{x_j | j \in [n]\}$  represents the set of decision variable nodes, and  $E = \{e_{i,j} | i \in [m], j \in [n]\}$  represents the edges between them. Here,  $[\cdot]$  denotes a set of consecutive numbers. The attribute of a constraint vertex  $s_i$  is expressed as  $\text{attr}(s_i) = [l_i^s, u_i^s]^\top$ , indicating its lower and upper bounds. Similarly, the attribute of a decision variable vertex  $x_j$  is given by  $\text{attr}(x_j) = [l_j^x, u_j^x, c_j]^\top$ , which includes its lower bound  $l_j^x$ , upper bound  $u_j^x$ , and objective coefficient  $c_j$ . This graph-based approach facilitates a structural analysis of the decision-making framework, allowing us to visualize the interactions among variables under different conditions and to compare updated and original programs on a structural level.

**Graph Edit Distance (GED) Calculation:** To quantify the impact of changes in “Decision Information”, we compute the GED between the two bipartite graphs derived from the updated and original LPs. GED represents the minimum cost necessary to transform one graph into another through a sequence of operations, such as inserting, deleting, or substituting vertices and edges. This metric

effectively captures the types of modifications made to the code in response to a query. This metric quantifies the minimal number of modifications (such as adding, deleting, or substituting nodes or edges) required to transform one graph into the other. A smaller edit distance indicates fewer changes, suggesting the updated program closely aligns with the original decision-making context. Conversely, a larger edit distance highlights significant alterations, signaling a substantial impact of the updated information on the decision-making process.

There are many well-established GED algorithms (Gao et al., 2010; Stauffer et al., 2017; Abu-Aisheh et al., 2015; Xing et al., 2024), we follow a straightforward principle provided by (Xing et al., 2024): each operation on an attribute in graph incurs a unit cost of 1. Formally, given the graph of the original program  $\mathcal{G}^p = (S^p \cup X^p, E^p)$  and updated program  $\mathcal{G}^{p'} = (S^{p'} \cup X^{p'}, E^{p'})$ , we define the vertex cost matrix as follows,

$$\begin{array}{c|ccc}
 & s_i^p \in S^p & x_i^p \in X^p & \epsilon \\
 \hline
 s_i^{p'} \in S^{p'} & \#\text{msm}(s_i^{p'}, s_i^p) & \infty & \#\text{attr}(s_i^{p'}) \\
 x_i^{p'} \in X^{p'} & \infty & \#\text{msm}(x_i^{p'}, x_i^p) & \#\text{attr}(x_i^{p'}) \\
 \epsilon & \#\text{attr}(s_i^p) & \#\text{attr}(x_i^p) & \infty
 \end{array}$$

In this context, all operations are treated as matching processes; for example, deleting a vertex is conceptualized as matching the vertex to an empty vertex, denoted by  $\epsilon$ . Here,  $\#\text{attr}(\cdot)$  denotes the total number of vertex attributes, while  $\#\text{msm}(\cdot)$  represents the count of mismatched attributes between two vertices. Similarly, the edge cost matrix is defined as follows,

$$\begin{array}{c|cc}
 & e_{i,j}^p \in E^p & \epsilon \\
 \hline
 e_{i,j}^{p'} \in E^{p'} & \#\text{msm}(e_{i,j}^{p'}, e_{i,j}^p) & \#\text{attr}(e_{i,j}^{p'}) \\
 \epsilon & \#\text{attr}(e_{i,j}^p) & \infty
 \end{array}$$

Given the varying scales of changes in parameters and constraints, we further normalize  $\text{GED}(\mathcal{G}^{p'}, \mathcal{G}^p)$  by the graph size. The normalized GED (NGED) is calculated as:  $\text{NGED}(\mathcal{G}^{p'}, \mathcal{G}^p) = \frac{\text{GED}(\mathcal{G}^{p'}, \mathcal{G}^p)}{\max(|\mathcal{G}^{p'}|, |\mathcal{G}^p|)}$ , where  $|\mathcal{G}|$  is defined as the sum of the number of attributes for all edges and vertices in the graph, specifically:  $|\mathcal{G}| = \sum_{e \in E} \#\text{attr}(e) + \sum_{v \in S \cup X} \#\text{attr}(v)$ . This normalization accounts for graph size variations, allowing for a more consistent comparison of the impact of changes across different scenarios. By quantifying discrepancies between the updated and original graphs through an analysis of their attribute changes, our approach offers a more precise measurement of changes in the decision-making framework.

By employing these three steps, LP conversion, graph representation, and graph edit distance calculation, we provide a rigorous and systematic approach to quantifying ‘‘Decision Information’’. Since LLMs cannot directly perform this quantification, we utilize them to sense these processes and generate explanatory insights. This approach fills the gap in previous methods that could not quantitatively analyze the impact of changes in constraints. By offering a unified methodology for measuring changes in constraints and parameters, it enables a precise evaluation of information changes. This framework provides a measurable and comparable basis for assessing the impact of these changes on decision-making. It thus anchors the concept of ‘‘Decision Information’’ in both practical and theoretical contexts, offering more profound insight into how updates shape decisions.

## 4 EXPERIMENTS

### 4.1 EVALUATION BENCHMARK

Despite the availability of numerous open-source datasets in OR, such as NL4OPT (Ramamonjison et al., 2023), ComplexOR (Xiao et al., 2023), NLP4LP (Holzer et al., 2024), and IndustryOR (Tang et al., 2024), these datasets are limited to problem descriptions and are primarily suited for OR modeling needs. There remains a significant gap in datasets specifically tailored for explainable OR, which are crucial for advancing transparency and interpretability in decision-making processes.

To address this issue, we developed a novel benchmark dataset based on the open-source commercial IndustryOR, specifically designed to evaluate explainability in OR tasks. The benchmark

324 includes 30 categorized problems across various domains (e.g., supply chain management, financial  
325 investment, logistics management, etc.). Each problem is paired with 10 unique queries that involve  
326 diverse or combined modifications to parameters or constraints (e.g., deleting, adding, or updating  
327 constraints and parameters). These queries were crafted by OR experts with significant industry  
328 experience to ensure both diversity and practical relevance. The dataset’s quality and comprehen-  
329 siveness are validated through iterative expert reviews and comparisons with real-world cases. For  
330 every problem, we provide corresponding Python code and employ the Gurobi optimization solver  
331 (Bixby, 2007) to determine optimal solutions. Additionally, we also include the ground truth labels  
332 for each query for each problem to ensure accurate evaluation. Notably, the question sets and the  
333 queries in this benchmark are developed from scratch and managed in-house, guaranteeing that they  
334 have not been part of LLM training data, enhancing the robustness of the benchmark for assessing  
335 model performance. We provide a complete example in Appendix A.3.

## 336 4.2 EVALUATION METHODOLOGY

337  
338 Our evaluation focuses on two aspects: **Modeling Accuracy** and **Explanation Quality**. For the  
339 modeling accuracy assessment, we recognize that different code implementations can produce the  
340 same optimization results. For instance, two programs solving a linear programming problem may  
341 employ different formulations or techniques (e.g., distinct constraint orderings or variable namings),  
342 but both can still yield identical optimal solutions. Therefore, rather than directly comparing the  
343 generated code to a reference, we evaluate accuracy by comparing the optimization outcomes to  
344 ensure correctness and consistency across implementations.

345 Regarding the explanation quality, although our task is highly specialized and requires expert-level  
346 interpretability in OR, we aim to develop an automated evaluation method. Drawing inspiration  
347 from (Kondapaneni et al., 2024) and utilizing the capabilities of LLMs for text quality evaluation  
348 (Chen et al., 2023; Chiang & Lee, 2023; Hu et al., 2024; Chu et al., 2024; Zytek et al., 2024),  
349 we establish expert-crafted templates and use LLMs to assess explanation quality, aiming for a  
350 human-level standard. However, this evaluation approach has certain limitations. To address this,  
351 we propose two methods in this paper. First, we establish a structured template that specifies key  
352 criteria for effective explanations, such as clarity, relevance, and logical coherence. For instance,  
353 explanations should explicitly describe the rationale for modifying specific parameters or constraints  
354 and explain how these changes affect the optimization results. Second, we conduct a blind review  
355 process where OR experts anonymously score the explanations generated by different methods.  
356 This approach helps minimize bias, providing a reliable measure of how effectively the explanations  
357 convey meaningful insights to users. We will assess whether the proposed automated evaluation  
358 method aligns with human evaluation. This dual approach enables us to assess the consistency  
359 between the proposed automated method and human evaluation.

## 360 4.3 BASELINES

361  
362 We employ two baselines to ensure a comprehensive evaluation: Standard and OptiGuide (Li et al.,  
363 2023). The Standard involves a proprietary LLM (e.g., GPT-4, GPT-4-Turbo, etc.) that gener-  
364 ates updated programs and explanations, serving as a basic comparison point for assessing overall  
365 performance. The OptiGuide represents a specialized method used in supply chain optimization,  
366 providing a domain-specific comparison that evaluates the adaptability and effectiveness of LLMs  
367 in industry-relevant scenarios.

## 368 4.4 MODEL SETUP

369  
370 For our experiments, we evaluate the performance of the proposed baselines using four LLMs, GPT-  
371 4 (Achiam et al., 2023), GPT-4-1106-preview, GPT-4-0125-preview, and GPT-4-Turbo, under both  
372 zero-shot and one-shot learning settings. The zero-shot setting requires the models to generate up-  
373 dated programs and explanations without any prior examples, testing their inherent understanding  
374 and generalization capabilities. The one-shot setting provides a single example to guide the models,  
375 allowing us to assess the impact of minimal contextual information on their ability to perform the  
376 tasks accurately and coherently. To ensure fairness and reproducibility, we fix the hyperparameter  
377 *temperature* at 0 and apply the same examples for all models in the one-shot setting, and all models  
are implemented under the framework AutoGen (Wu et al., 2024). This paper focuses on fully au-

Table 2: Accuracy across different models under different LLMs with zero/one-shot setting. The bold scores are the best in each row.

Setting	Model	Standard	OptiGuide	EOR
Zero-shot	GPT-4	18.67%	30.33%	<b>75.67%</b>
	GPT-4-1106-preview	75.33%	36.00%	<b>81.67%</b>
	GPT-4-0125-preview	68.33%	47.33%	<b>87.33%</b>
	GPT-4-Turbo	63.00%	30.33%	<b>88.33%</b>
One-shot	GPT-4	71.67%	75.00%	<b>90.33%</b>
	GPT-4-1106-preview	69.67%	55.33%	<b>87.67%</b>
	GPT-4-0125-preview	76.00%	69.33%	<b>92.00%</b>
	GPT-4-Turbo	88.00%	69.33%	<b>95.33%</b>

tomating all processes, excluding user involvement. However, as our implementation is built on the AutoGen, it inherently supports seamless integration of user feedback. A detailed hyperparameter sensitivity analysis is provided in the Appendix A.4.

## 4.5 QUANTITATIVE PERFORMANCE

### 4.5.1 COMPARISON OF MODELING ACCURACY

Table 2 shows the accuracy results across different models. We have the following observations. In the zero-shot setting, EOR consistently outperforms both Standard and OptiGuide across all LLM versions. These results emphasize the robustness of the EOR in zero-shot tasks, which substantially improves performance compared to other methods. For instance, GPT-4-Turbo achieves 88.33% accuracy with EOR, while Standard and OptiGuide methods yield only 63.00% and 30.33%, respectively. Moreover, for all LLM models except GPT-4, Standard outperforms OptiGuide, indicating that the LLM’s modeling capabilities are quite strong.

In the one-shot setting, EOR continues to outperform Standard and OptiGuide in all LLM versions, achieving an average of 90.00% accuracy and even reaching 95.33% accuracy on the GPT-4-Turbo, the highest among all results. Additionally, we find that providing an example can significantly improve modeling accuracy, and nearly all methods perform better in the one-shot setting than in the zero-shot setting. Specifically, the accuracy of OptiGuide on GPT-4 improves from 30.33% to 75.00%. However, OptiGuide still produces the worst result except on GPT-4.

Overall, EOR consistently outperforms other methods in both zero-shot and one-shot settings, particularly with models like GPT-4-Turbo and GPT-4-0125-preview. The accuracy gains observed when transitioning from zero-shot to one-shot highlight the importance of using examples in improving model performance. In terms of model comparisons, GPT-4-Turbo demonstrates the highest adaptability across both settings, achieving the best overall accuracy. While OptiGuide provides modest improvements, it does not match the performance of Standard and EOR. These results underscore the value of carefully selecting both models and example strategies to maximize accuracy, with EOR emerging as the most effective for high-accuracy tasks.

Table 3: Quality scores (0-10) across different models with zero/one-shot setting.

Method		EC $\uparrow$		ER $\uparrow$		Overall $\uparrow$	
		Auto	Expert	Auto	Expert	Auto	Expert
Zero-shot	Standard	0.12	-	6.98	6.98	5.36	5.63
	OptiGuide	0.69	-	7.93	7.90	6.31	6.59
	EOR	9.76	9.86	9.41	9.47	9.47	9.47
One-shot	Standard	0.34	-	7.10	6.90	5.54	5.37
	OptiGuide	0.20	-	7.60	7.61	5.96	6.05
	EOR	9.61	9.72	9.30	9.35	9.33	9.35



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

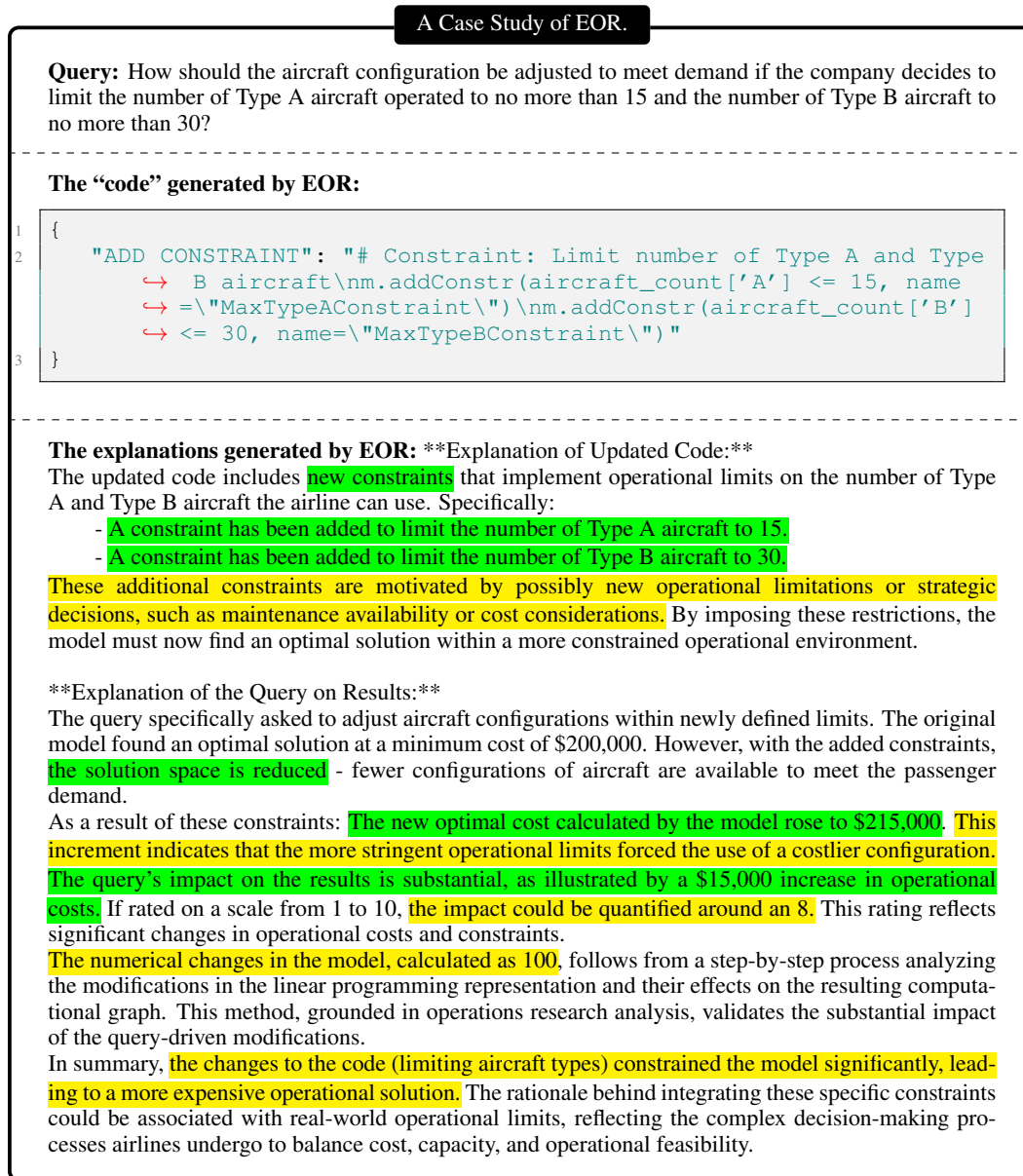


Figure 3: An example illustrating the code and explanations generated by EOR.

#### 4.5.2 COMPARISON OF THE QUALITY OF EXPLANATIONS

As discussed in Sec. 4.2, a critical part of our evaluation is assessing the quality of the explanations generated by the models. It is important to note that explanations based on incorrect results are irrelevant, as explaining failure cases offers no meaningful insights. Therefore, we first filter out all incorrect modeling cases, ensuring that only correct outputs are evaluated. We adopt two evaluation methods: an automated evaluation (Auto) and expert evaluation (Expert), both introduced in Sec. 4.2. These evaluations focus on two main aspects: Explanation of Correctness (EC) and Explanation of Results (ER). EC assesses the clarity and correctness of the explanation regarding the model’s output, while ER evaluates how well the explanation conveys the reasoning behind the optimization results. Additionally, we also measure the overall explanation quality (Overall). The explanations are scored on a 0-10 scale, with 0 indicating poor quality and 10 indicating excellent quality. The prompt template for evaluating the explanation quality is provided in Appendix A.5. The results of these evaluations using GPT-4o are summarized in Table 3.

As shown in Table 3, we have the following observations. Firstly, Auto is nearly as effective as Expert, with minimal differences in scores across all metrics, indicating that Auto is a reliable and valid approach. It is worth noting that Auto scores individual queries on EC, although it does not explicitly evaluate explanations. Secondly, EOR consistently outperforms both Standard and OptiGuide in explanation quality, achieving the highest scores in both zero-shot and one-shot settings. This demonstrates the superiority of our method in providing clearer and more accurate interpretations compared to other models. Finally, we find that Auto scores slightly higher on Standard than Expert, while for other models, the opposite is observed, with Auto scoring slightly lower than Expert. One possible reason is that LLMs may be biased toward the results they generate.

#### 4.6 CASE STUDY

In this section, we provide an example of the code and explanations generated by EOR, as illustrated in Figure 3, where simple explanations are highlighted in green to represent attribute explanations, while deeper ones in yellow indicate justification explanations with quantitative analysis. Results from other models are provided in Appendix A.6. The EOR approach stands out by not only providing a clear explanation of the updated code but also offering detailed insights into how the query affects the outcomes. It explicitly highlights the newly added constraints that limit the number of Type A and Type B aircraft, demonstrating transparency in addressing the operational limits specified in the query. Unlike other methods that merely report result differences, EOR incorporates a quantitative analysis, explaining the \$15,000 increase in operational costs due to the restricted solution space caused by the newly added constraints. This step-by-step explanation not only justifies the modifications but also clarifies the broader impact of the constraints, making EOR more comprehensive and contextually aware than other models.

#### 4.7 FAILURE CASES

To understand the strengths and weaknesses of EOR, we analyzed common failure cases, summarized in Table 4. These failures can be grouped into three categories: JSON format errors, where the LLMs fail to generate correct JSON outputs; correct execution, where code runs but produces incorrect results due to modeling logic errors or incomplete modeling; and runtime errors, such as variable name errors, syntax errors, or indentation errors. Table 4 shows a 60.00% reduction in total errors from zero-shot to one-shot, demonstrating a substantial improvement in the model’s performance. No syntax errors in the one-shot results highlights the effectiveness of EOR. However, ongoing attention to modeling logic and runtime errors is still crucial for further improvement.

Table 4: Failure cases on GPT-4-Turbo with zero/one-shot setting.

Failure Types		Zero-shot (Total 35)		One-shot (Total 14)	
		#Number	Percentage	#Number	Percentage
JSON Format Errors		8	22.86%	2	14.29%
Correct Execution	Modeling Logic Errors	<b>13</b>	<b>37.14%</b>	<b>4</b>	<b>28.57%</b>
	Incomplete Modeling	4	11.43%	2	14.29%
Runtime Errors	Variable Name Errors	2	5.71%	3	21.43%
	Syntax Errors	6	17.14%	0	0.00%
	Indent Errors	2	5.71%	3	21.43%

## 5 CONCLUSION AND FUTURE WORK

In this paper, we present EOR, a novel framework that addresses transparency and interpretability challenges in OR. We introduce the concept of “Decision Information” through what-if analysis and use bipartite graphs to quantify changes in OR models. By leveraging LLM capabilities, our framework enhances the modeling process and generates actionable, comprehensible explanations. Furthermore, we establish the first industrial benchmark for evaluating explanation quality in OR. In the future, we aim to develop more accurate modeling methods to improve the accuracy while exploring effective automated evaluation techniques for explanations.

## REFERENCES

- 540  
541  
542 Zeina Abu-Aisheh, Romain Raveaux, Jean-Yves Ramel, and Patrick Martineau. An exact graph edit  
543 distance algorithm for solving pattern recognition problems. In *4th International Conference on*  
544 *Pattern Recognition Applications and Methods 2015*, 2015.
- 545  
546 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
547 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
548 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 549  
550 Ali AhmadiTeshnizi, Wenzhi Gao, and Madeleine Udell. Optimus: Scalable optimization modeling  
551 with (mi) Ip solvers and large language models. *arXiv preprint arXiv:2402.10172*, 2024.
- 552  
553 Bob Bixby. The gurobi optimizer. *Transfp. Re-search Part B*, 41(2):159–178, 2007.
- 554  
555 Yi Chen, Rui Wang, Haiyun Jiang, Shuming Shi, and Ruifeng Xu. Exploring the use of large  
556 language models for reference-free text quality evaluation: An empirical study. *arXiv preprint*  
557 *arXiv:2304.00723*, 2023.
- 558  
559 Cheng-Han Chiang and Hung-yi Lee. A closer look into using large language models for automatic  
560 evaluation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp.  
561 8928–8942, 2023.
- 562  
563 KuanChao Chu, Yi-Pei Chen, and Hideki Nakayama. A better llm evaluator for text generation: The  
564 impact of prompt output sequencing and optimization. *arXiv preprint arXiv:2406.09972*, 2024.
- 565  
566 IBM ILOG Cplex. V12. 1: User’s manual for cplex. *International Business Machines Corporation*,  
567 46(53):157, 2009.
- 568  
569 Kristijonas Čygas, Dimitrios Letsios, Ruth Misener, and Francesca Toni. Argumentation for explain-  
570 able scheduling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 2752–2759,  
571 2019.
- 572  
573 Koen W De Bock, Kristof Coussement, Arno De Caigny, Roman Słowiński, Bart Baesens, Robert N  
574 Boute, Tsan-Ming Choi, Dursun Delen, Mathias Kraus, Stefan Lessmann, et al. Explainable ai  
575 for operational research: A defining framework, methods, applications, and a research agenda.  
576 *European Journal of Operational Research*, 317(2):249–272, 2024.
- 577  
578 Martin Erwig and Prashant Kumar. Explanations for combinatorial optimization problems. *Journal*  
579 *of Computer Languages*, 79:101272, 2024.
- 580  
581 Zhenan Fan, Xinglu Wang, Oleksandr Yakovenko, Abdullah Ali Sivas, Owen Ren, Yong Zhang,  
582 and Zirui Zhou. Smart initial basis selection for linear programs. In *International Conference on*  
583 *Machine Learning*, pp. 9650–9664. PMLR, 2023.
- 584  
585 Anja K Faulhaber, Ina Ni, and Ludger Schmidt. The effect of explanations on trust in an assistance  
586 system for public transport users and the role of the propensity to trust. In *Proceedings of Mensch*  
587 *und Computer 2021*, pp. 303–310. Association for Computing Machinery, 2021.
- 588  
589 Xinbo Gao, Bing Xiao, Dacheng Tao, and Xuelong Li. A survey of graph edit distance. *Pattern*  
590 *Analysis and applications*, 13:113–129, 2010.
- 591  
592 Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. Exact combi-  
593 natorial optimization with graph convolutional neural networks. *Advances in neural information*  
processing systems, 32, 2019.
- Dongdong Ge, Qi Huangfu, Zizhuo Wang, Jian Wu, and Yinyu Ye. Cardinal Optimizer (COPT) user  
guide. <https://guide.coap.online/copt/en-doc>, 2023.
- Bryce Goodman and Seth Flaxman. European union regulations on algorithmic decision-making  
and a “right to explanation”. *AI magazine*, 38(3):50–57, 2017.

- 594 Jesse T Holzer, Carleton J Coffrin, Christopher DeMarco, Ray Duthu, Stephen T Elbert, Brent C  
595 Eldridge, Tarek Elgindy, Manuel Garcia, Scott L Greene, Nongchao Guo, et al. Grid optimization  
596 competition challenge 3 problem formulation. Technical report, Pacific Northwest National  
597 Laboratory (PNNL), Richland, WA (United States), 2024.
- 598 Xinyu Hu, Mingqi Gao, Sen Hu, Yang Zhang, Yicheng Chen, Teng Xu, and Xiaojun Wan. Are  
599 llm-based evaluators confusing nlg quality criteria? *arXiv preprint arXiv:2402.12055*, 2024.
- 600 Sen Huang, Kaixiang Yang, Sheng Qi, and Rui Wang. When large language model meets optimization.  
601 *arXiv preprint arXiv:2405.10098*, 2024.
- 602 Neehar Kondapaneni, Markus Marks, Oisín Mac Aodha, and Pietro Perona. Less is more: Discovering  
603 concise network explanations. In *ICLR 2024 Workshop on Representational Alignment*,  
604 2024.
- 605 Beibin Li, Konstantina Mellou, Bo Zhang, Jeevan Pathuri, and Ishai Menache. Large language  
606 models for supply chain optimization. *arXiv preprint arXiv:2307.03875*, 2023.
- 607 Mahdi Mostajabdaveh, Timothy T Yu, Rindranirina Ramamonjison, Giuseppe Carenini, Zirui Zhou,  
608 and Yong Zhang. Optimization modeling and verification from problem specifications using a  
609 multi-agent multi-stage llm framework. *INFOR: Information Systems and Operational Research*,  
610 pp. 1–19, 2024.
- 611 Rindranirina Ramamonjison, Timothy Yu, Raymond Li, Haley Li, Giuseppe Carenini, Bissan Ghad-  
612 dar, Shiqi He, Mahdi Mostajabdaveh, Amin Banitalebi-Dehkordi, Zirui Zhou, et al. N14opt com-  
613 petition: Formulating optimization problems based on their natural language descriptions. In  
614 *NeurIPS 2022 Competition Track*, pp. 189–203. PMLR, 2023.
- 615 Andrew Selbst and Julia Powles. “meaningful information” and the right to explanation. In *confer-  
616 ence on fairness, accountability and transparency*, pp. 48–48. PMLR, 2018.
- 617 Michael Stauffer, Thomas Tschachtli, Andreas Fischer, and Kaspar Riesen. A survey on applica-  
618 tions of bipartite graph edit distance. In *Graph-Based Representations in Pattern Recognition:  
619 11th IAPR-TC-15 International Workshop, GbRPR 2017, Anacapri, Italy, May 16–18, 2017, Pro-  
620 ceedings 11*, pp. 242–252. Springer, 2017.
- 621 Zhengyang Tang, Chenyu Huang, Xin Zheng, Shixi Hu, Zizhuo Wang, Dongdong Ge, and Benyou  
622 Wang. Orlm: Training large language models for optimization modeling. *arXiv preprint  
623 arXiv:2405.17743*, 2024.
- 624 Arthur Thuy and Dries F Benoit. Explainability through uncertainty: Trustworthy decision-making  
625 with neural networks. *European Journal of Operational Research*, 317(2):330–340, 2024.
- 626 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun  
627 Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and  
628 Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework.  
629 In *COLM*, 2024.
- 630 Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin  
631 Fu, Tao Zhong, Jia Zeng, Mingli Song, et al. Chain-of-experts: When llms meet complex opera-  
632 tions research problems. In *The Twelfth International Conference on Learning Representations*,  
633 2023.
- 634 Linzi Xing, Xinglu Wang, Yuxi Feng, Zhenan Fan, Jing Xiong, Zhijiang Guo, Xiaojin Fu, Rindra  
635 Ramamonjison, Mahdi Mostajabdaveh, Xiongwei Han, et al. Towards human-aligned evaluation  
636 for linear programming word problems. In *Proceedings of the 2024 Joint International Confer-  
637 ence on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*,  
638 pp. 16550–16556, 2024.
- 639 Jihai Zhang, Wei Wang, Siyan Guo, Li Wang, Fangquan Lin, Cheng Yang, and Wotao Yin. Solving  
640 general natural-language-description optimization problems with large language models. *arXiv  
641 preprint arXiv:2407.07924*, 2024.
- 642 Alexandra Zyteck, Sara Pidò, and Kalyan Veeramachaneni. Llms for xai: Future directions for  
643 explaining explanations. *arXiv preprint arXiv:2405.06064*, 2024.
- 644
- 645
- 646
- 647

## 648 A APPENDIX

### 649 A.1 THE COMPARISON BETWEEN WHAT-IF ANALYSIS, SENSITIVITY ANALYSIS, AND 650 DECISION INFORMATION ANALYSIS

651 **What-if Analysis:** What-if analysis explores the impact of changing parameters or conditions on  
652 overall outcomes, allowing users to adjust inputs to observe potential scenarios manually. This  
653 approach emphasizes scenario exploration at a macro level, helping users understand how different  
654 “what if” situations might influence results without delving into specific details.

655 **Sensitivity Analysis:** Sensitivity analysis measures how minor variations in input parameters affect  
656 the model’s output, identifying which inputs have the most significant influence on results. This  
657 method provides detailed quantitative insights into the effects of parameter changes, making it par-  
658 ticularly effective in optimization contexts where analytical solutions can be derived.

659 **Decision Information Analysis:** Decision information analysis extends sensitivity analysis by fo-  
660 cusing on how changes in constraints impact decision-making and outcomes. This approach exam-  
661 ines the sensitivity of the model to variations in constraints, identifying key decision factors while  
662 also capturing the broader implications of both constraint and parameter changes, particularly in  
663 complex optimization problems without analytical solutions.

664 In summary, decision information and sensitivity analysis are more detailed, with the former center-  
665 ing on decision-related changes (constraints or parameters) and the latter on parameter sensitivity,  
666 while both can be part of a broader what-if analysis framework. Previous research has predominantly  
667 focused on sensitivity analysis through the simplex method, which is recognized for its efficiency,  
668 precision, and ability to provide analytical solutions. However, the simplex method has limitations  
669 in evaluating the effects of changes to constraints. In contrast, our proposed approach leverages  
670 bipartite graphs to assess the impact of constraint modifications, addressing this gap and providing  
671 a more generalized solution.

### 672 A.2 PROMPT TEMPLATE DESIGN

#### 673 A.2.1 PROMPT TEMPLATE FOR WRITER AGENT

674 The prompt template for Writer with system message for the ChatCompletion inference:

```
680 1 WRITER_SYSTEM_MSG = """
681 2
682 3 **Role:** You are a chatbot tasked with:
683 4 (1) Writing Python code for operations research-related projects.
684 5 (2) Explaining solutions using the Gurobi Python solver.
685 6
686 7 --- Problem Description: ---
687 8 {description}
688 9
689 10 --- Source Code: ---
690 11 {source_code}
691 12
692 13 --- Documentation: ---
693 14 {doc_str}
694 15
695 16 --- Example Q&A: ---
696 17 {example_qa}
697 18
698 19 --- Original Execution Result: ---
699 20 {execution_result}
700 21
701 22 **Task:**
702 23 You are provided with the original problem description and the correct
703 24 ↪ code for an operations research problem. Based on the user’s query,
704 ↪ update the code accordingly. Your task may involve either deleting
705 ↪ constraints or adding new data or constraints.
```

```

702 25
703 26
704 27 **Steps:**
705 28
706 29 1. Determine the Required Operations:
707 30
708 31 - **Add Operation:** Generate new code for data or constraints to be
709 32 ↪ added.
710 33 - Insert the new data between the markers:
711 34 "# ORExplainer DATA CODE GOES HERE" and "# ORExplainer DATA CODE
712 35 ↪ ENDS HERE".
713 36 - Insert the new constraints between the markers:
714 37 "# ORExplainer CONSTRAINT CODE GOES HERE" and "# ORExplainer
715 38 ↪ CONSTRAINT CODE ENDS HERE".
716 39
717 40 - **Delete Operation:** Identify the relevant block within the
718 41 ↪ constraints section that needs to be deleted.
719 42 - The code to be deleted will be between the markers:
720 43 "# ORExplainer CONSTRAINT CODE GOES HERE" and "# ORExplainer
721 44 ↪ CONSTRAINT CODE ENDS HERE".
722 45
723 46 2. Return the Changes in JSON Format:
724 47
725 48 - Use the keys "DELETE CONSTRAINT", "ADD CONSTRAINT", or "ADD DATA".
726 49 ↪ Only these keys are allowed.
727 50 - The values should be the Python code snippet to be deleted (exactly
728 51 ↪ as it appears in the original code) or the new code to be added.
729 52 - Include comments within the code snippets using the prefix "#".
730 53 - Ensure the line breaks and indents of the code are correct.
731 54
732 55 3. Output Requirements:
733 56
734 57 - Return only the JSON object with the changes.
735 58 - Do not include any additional information in the response.
736 59 - Do not add new decision variables.
737 60 - Ensure the JSON is valid, properly formatted.
738 61
739 62 The above explained instructions are your guide to accomplish the task
740 63 ↪ effectively. Your user's success heavily relies upon your ability
741 64 ↪ to provide the precise and accurate Python code changes within the
742 65 ↪ existing operations research problem. Good Luck!
743 66
744 67 """
745 68
746 69
747 70
748 71
749 72
750 73
751 74
752 75
753 76
754 77
755 78

```

The prompt template for Code (2):

```

741 1 CODE_PROMPT = """
742 2
743 3 **Role:** You are a professional software developer tasked with handling
744 4 ↪ code modification requests. Your role is to interpret these
745 5 ↪ requests which describe changes needed in a source code.
746 6
747 7 **Task:** Your task is to return the changes to be made to the source
748 8 ↪ code based on the user's query. The modifications should be
749 9 ↪ returned in a JSON format containing only the necessary changes.
750 10
751 11 **JSON Format Requirements:**
752 12
753 13 - Use the keys "DELETE CONSTRAINT", "ADD CONSTRAINT", or "ADD DATA".
754 14 ↪ Only these keys are allowed.
755 15 - The values should be the Python code snippet to be deleted (exactly
756 16 ↪ as it appears in the original code) or the new code to be added.

```

```

756 13 - Include comments within the code snippets using the prefix "#".
757 14 - Ensure the line breaks and indents of the code are correct.
758 15
759 16
760 17 **Output Requirements:**
761 18
762 19 - Return only the JSON object with the changes.
763 20 - Do not include any additional information in the response.
764 21 - Do not add new decision variables.
765 22 - Ensure the JSON is valid, properly formatted.
766 23
767 24
768 25 The above explained instructions are your guide to accomplish the task
769 26 ↪ effectively. Your user's success heavily relies upon your ability
770 27 ↪ to provide the precise and accurate Python code changes within the
771 28 ↪ existing operations research problem. Good Luck!
772 29
773 30 --- Answer Code: ---
774 31 """

```

#### The prompt template for Debug (2):

```

777 1 DEBUG_PROMPT = """
778 2
779 3 **Role:** You are a professional code debugger.
780 4
781 5
782 6 **Task:** Identify and fix the error in the code, ensuring the corrected
783 7 ↪ version runs smoothly and error-free.
784 8
785 9 **Details:**
786 10 - Error Type: {error_type}
787 11 - Error Message: {error_message}
788 12
789 13 **Instructions:**
790 14 Please analyze the error details, resolve the bug based on the type and
791 15 ↪ message provided, and rewrite the corrected version of the code
792 16 ↪ snippet below.
793 17
794 18 **Corrected Code:**
795 19 --- NEW CODE ---
796 20
797 21 """

```

#### The prompt template for Interpreter (6):

```

800 1 INTERPRETER_PROMPT = """
801 2
802 3 **Role:** You are a skilled interpreter with expertise in analyzing and
803 4 ↪ explaining changes in computational model code and their effects on
804 5 ↪ results.
805 6
806 7 **Task:** Present a clear and thorough explanation of the code updates
807 8 ↪ and their effects on the results. Structure your explanation in two
808 9 ↪ key parts: Explanation of the updated code and Explanation of the
809 10 ↪ Query on Results.
810 11
811 12 """

```

```

810 9  **Inputs:**
811 10  - Original code: {source_code}
812 11  - Updated code: {new_code}
813 12  - Code changes induced by the query: {json_data}
814 13  - Original execution results: {original_execution_result}
815 14  - New execution results: {execution_rst}
816 15  - Measure of numerical changes in the model induced by the query: {
817 16    ↪ different_model}
818 17
819 18  **Key Points to Understand:**
820 19
821 20  1. Explanation of Updated code:
822 21
823 22  - Explain the rationale behind each specific change to the code, such
824 23    ↪ as why certain constraints or data were added, deleted, or
825 24    ↪ modified.
826 25
827 26  2. Explanation of the Query on Results:
828 27
829 28  - Clarify why the specific results were produced in response to the
830 29    ↪ query.
831 30  - Assess the query's impact on the results by comparing the new
832 31    ↪ execution results with the original ones and the corresponding
833 32    ↪ numerical changes in the model.
834 33  - Use a scale from 1 to 10 to quantify the query's impact on the
835 34    ↪ results, with 1 indicating minimal impact and 10 indicating
836 35    ↪ significant impact.
837 36
838 37  **Background for Numerical Changes Calculation:**
839 38
840 39  The impact of the query is measured using a three-step process:
841 40  1. LP Conversion: The problem is converted into a linear programming (
842 41    ↪ LP) format to identify key components and constraints.
843 42  2. Graph Representation: The LP model is then represented as a
844 43    ↪ bipartite graph, where nodes and edges correspond to variables,
845 44    ↪ constraints, and relationships.
846 45  3. Graph Edit Distance Calculation: The difference between the
847 46    ↪ original and modified graphs is computed by measuring the graph
848 47    ↪ edit distance, which involves operations like insertion,
849 48    ↪ deletion, and substitution of nodes and edges, each with a unit
850 49    ↪ cost of 1.
851 50
852 51  **Output:**
853 52
854 53  Provide the explanations in two distinct parts:
855 54  (1) Explanation of the Updated code
856 55  (2) Explanation of the Query on Results
857 56
858 57  **Requirements:**
859 58
860 59  - Ensure the explanations are detailed and comprehensive, covering all
861 60    ↪ relevant aspects of the code updates and their impact on the
862 61    ↪ results.
863 62  - Ensure that explanations are delivered in a narrative style, suited
864 63    ↪ for a non-technical audience, avoiding jargon or direct
865 64    ↪ references to specific variable names.
866 65  - Offer clear, precise, easy-to-understand descriptions that
867 66    ↪ effectively bridge complex information with clarity and insight.

```



```

864 53 The above explained instructions are your guide to accomplish the task
865      ↪ effectively. Your user's success heavily relies upon your ability
866      ↪ to provide the explanations within the existing operations research
867      ↪ problem. Good Luck!
868 54
869 55 --- HUMAN READABLE ANSWER ---
870 56
871 57 ""
872 58

```

## 873 A.2.2 PROMPT TEMPLATE FOR SAFEGUARD AGENT

874 The prompt template for Safeguard with system message for the ChatCompletion inference:

```

877 1 SAFEGUARD_SYSTEM_MSG = ""
878 2
879 3 **Role:** You are a code safety evaluator.
880 4
881 5
882 6 **Task:** Review the provided source code to determine if it is safe to
883      ↪ execute, ensuring it does not contain any malicious code that could
884      ↪ compromise security or privacy.
885 7
886 8
887 9 **Instructions:**
888 10 --- Source Code: ---
889 11 {source_code}
890 12
891 13
892 14 **Question:**
893 15 Is the code safe to run?
894 16
895 17
896 18 **Answer:**
897 19 Respond with one word:
898 20     `SAFE` if the code is secure.
899 21     `DANGER` if the code poses any risk.
900 22
901 23 ""

```

902 The prompt template for Safe (4):

```

903 1 SAFEGUARD_PROMPT = ""
904 2
905 3 **Role:** You are a professional code safety evaluator.
906 4
907 5
908 6 **Task:** Examine the safety of each code snippet contained within a
909      ↪ provided JSON file.
910 7
911 8
912 9 **Details:**
913 10
914 11 - **Code Structured as JSON:** Each value in the JSON represents a
915      ↪ code snippet intended for review. These snippets may be newly
916      ↪ generated or under consideration for deletion.
917 12
918 13
919 14 **Instructions:**
920 15
921 16 - Thoroughly analyze each code snippet found in the JSON.
922 17 - For each snippet, determine its safety for execution.
923 18 - Provide your assessment as a single word for each snippet: either "
924 19     ↪ SAFE" or "DANGER".

```

```

918 20
919 21 **Example of Expected Response:** For snippet_1: SAFE, for snippet_2:
920 22     ↪ DANGER
921 23
922 24 --- Answer: ---
923 25
924 26 ""
925

```

### 927 A.3 AN EXAMPLE OF THE BENCHMARK

928 The original problem description and the result of this problem:

```

930 1 Problem Description:
931 2 An airline operates two types of aircraft: large aircraft (Type A) and small aircraft (Type B). Each
932 3     ↪ type of aircraft has different operating costs and passenger capacities. The company needs
933 4     ↪ to determine the number of each type of aircraft to operate in order to meet the demand
934 5     ↪ of transporting at least 10,000 passengers. Specifically, one Type A aircraft can carry 500
935 6     ↪ passengers, and one Type B aircraft can carry 200 passengers. However, due to the use
936 7     ↪ and maintenance requirements of the aircraft, the total number of Type A and Type B
937 8     ↪ aircraft operated cannot exceed 50. The cost of operating one Type A aircraft is $10,000,
938 9     ↪ and the cost of operating one Type B aircraft is $5,000. Due to the indivisibility of
939 10    ↪ the aircraft, both types of aircraft must be operated in integer quantities. Under these
940 11    ↪ conditions, what is the minimum total cost that satisfies the passenger transportation
941 12    ↪ demand while adhering to the operational restrictions? Please round the answer to the
942 13    ↪ nearest integer.
943 14 Problem Result:
944 15 200000.0
945

```

946 The original Python code with Gurobi:

```

947 1 import gurobipy as gp
948 2 from gurobipy import GRB
949 3
950 4
951 5 # Parameters Section Begin
952 6 # Define model parameters
953 7 aircraft_types = ['A', 'B']
954 8
955 9 # Passenger capacity per aircraft type
956 10 passenger_capacity = {
957 11     'A': 500,
958 12     'B': 200
959 13 }
960 14
961 15 # Operating cost per aircraft type (dollars)
962 16 operating_cost_per_aircraft = {
963 17     'A': 10000,
964 18     'B': 5000
965 19 }
966 20
967 21 # Minimum passenger demand
968 22 min_passenger_demand = 10000
969 23
970 24 # Maximum number of aircraft
971 25 max_aircraft_count = 50
972 26 # Parameters Section End
973 27
974 28
975 29 # ORExplainer DATA CODE GOES HERE
976 30

```

```

972 31
973 32 # ORExplainer DATA CODE ENDS HERE
974 33
975 34
976 35 # Create a Gurobi model
977 36 m = gp.Model("AirlineOptimization")
978 37
979 38 # Decision Variables Section Begin
979 39 # Create decision variables a and b for the number of large and small
980 40 ↪ aircraft respectively
981 41 aircraft_count = {
982 42     'A': m.addVar(vtype=GRB.INTEGER, name="aircraft_A"),
983 43     'B': m.addVar(vtype=GRB.INTEGER, name="aircraft_B")
984 44 }
984 44 # Decision Variables Section End
985 45
986 46
987 47 # Objective Function Section Begin
988 48 # Set the objective function to minimize total operating cost
988 49 m.setObjective(
989 50     gp.quicksum(
990 51         operating_cost_per_aircraft[t] * aircraft_count[t]
991 52         for t in aircraft_types
992 53     ),
992 54     sense=GRB.MINIMIZE
993 55 )
994 56 # Objective Function Section End
995 57
996 58
997 59 # ORExplainer CONSTRAINTS CODE GOES HERE
998 60
999 61
999 62 # Constraints Section Begin
1000 63 # Constraint: Meet the passenger transportation demand
1001 64 m.addConstr(
1002 65     gp.quicksum(passenger_capacity[t] * aircraft_count[t]
1003 66                 for t in aircraft_types) >= min_passenger_demand,
1004 67     name="PassengerDemandConstraint"
1005 68 )
1006 69
1007 70 # Constraint: The total number of aircraft cannot exceed the maximum
1007 71 ↪ allowed
1008 72 m.addConstr(
1009 73     gp.quicksum(aircraft_count[t]
1010 74                 for t in aircraft_types) <= max_aircraft_count,
1011 75     name="OperationalConstraint"
1012 76 )
1012 76 # Constraints Section End
1013 77
1014 78
1014 79 # ORExplainer CONSTRAINTS CODE MIDDLE HERE
1015 80
1016 81
1017 82 # ORExplainer CONSTRAINTS CODE ENDS HERE
1018 83
1019 84
1019 85 # Solving the Model Section Begin
1020 86 # Solve the model
1021 87 m.optimize()
1022 88
1023 89 # Output the results
1024 90 if m.status == GRB.OPTIMAL:
1024 91     print(f"Minimum total cost: {round(m.ObjVal)} dollars")
1025 92     for t in aircraft_types:
1025 93         print(f"Number of Type {t} aircraft: {aircraft_count[t].X}")

```

```

1026 else:
1027     print("No optimal solution found.")
1028 # Solving the Model Section End

```

The queries and truth labels of these updated problems with new queries:

```

1032 Query 1: How should the number of aircraft be adjusted to maximize economic efficiency if the
1033     ↪ operating cost of the large aircraft (Type A) is reduced to $8,000?
1034 Truth Label 1: 160000.0
1035
1036 Query 2: How should the number of aircraft be reassessed to meet transportation demand if the
1037     ↪ passenger capacity of the small aircraft (Type B) increases to 250 passengers?
1038 Truth Label 2: 200000.0
1039
1040 Query 3: How should the aircraft configuration be adjusted to maximize profit if the operating cost
1041     ↪ of the Type B aircraft decreases to $4,000 and the passenger capacity of the Type A
1042     ↪ aircraft increases to 550 passengers?
1043 Truth Label 3: 184000.0
1044
1045 Query 4: How should the number of aircraft be adjusted to maintain minimum total cost if the
1046     ↪ operating costs of both Type A and Type B aircraft increase by 10%?
1047 Truth Label 4: 220000.0
1048
1049 Query 5: How should the aircraft configuration be adjusted to meet demand if the company decides
1050     ↪ to limit the number of Type A aircraft operated to no more than 15 and the number of
1051     ↪ Type B aircraft to no more than 30?
1052 Truth Label 5: 215000.0
1053
1054 Query 6: How should the number of aircraft be adjusted to maintain demand if the passenger
1055     ↪ capacity of the Type A aircraft decreases to 450 passengers and the operating cost of the
1056     ↪ Type B aircraft increases to $6,000?
1057 Truth Label 6: 226000.0
1058
1059 Query 7: How should the number of aircraft be adjusted to maximize transportation efficiency if
1060     ↪ the passenger capacity of the large aircraft increases to 600 passengers and the cost of the
1061     ↪ small aircraft increases to $5,500?
1062 Truth Label 7: 170000.0
1063
1064 Query 8: How should the number of aircraft be arranged to meet passenger demand and minimize
1065     ↪ costs if the airline needs to operate at least 10 Type B aircraft?
1066 Truth Label 8: 210000.0
1067
1068 Query 9: How should the number of aircraft be adjusted to maintain economic efficiency if the
1069     ↪ passenger capacity of the small aircraft (Type B) decreases to 150 passengers and the cost
1070     ↪ of the large aircraft (Type A) increases to $12,000?
1071 Truth Label 9: 240000.0
1072
1073 Query 10: How will the optimal aircraft configuration change if the constraint that the total number
1074     ↪ of Type A and Type B aircraft operated cannot exceed 50 is removed?
1075 Truth Label 10: 200000.0

```

#### A.4 HYPERPARAMETER SENSITIVITY ANALYSIS

This section analyzes the sensitivity of two hyperparameters, *temperature* and *debug times*, to evaluate their impact on the model's reliability and stability. The *temperature* reflects the reliability of the model's outputs, while *debug times* assess its performance stability. The experimental results are summarized in Tables 5 and 6.

Table 5 shows that our model maintains reliable outputs across temperature settings of 0, 0.5, and 1 in both zero-shot and one-shot scenarios. This demonstrates the model’s robustness in generating consistent and precise outputs under varying temperature configurations, which is particularly important for tasks demanding factual accuracy, such as OR modeling.

Table 6 reveals that increasing the number of debug attempts from 3 to 10 does not significantly improve performance. Additional iterations primarily consume resources without yielding better results, reflecting the model’s limited self-correction capabilities. This suggests that external interventions, such as user feedback, may be necessary to enhance performance further.

In summary, the model exhibits strong reliability across different temperature settings and consistent stability regardless of the number of debugging attempts. These findings underscore the robustness of our approach in delivering reliable and stable performance.

Table 5: Accuracy under different methods in zero/one-shot setting with different *temperature*. The bold scores are the best in each row.

Method		Temperature		
		0	0.5	1
Zero-shot	OptiGuide	<b>30.33%</b>	22.00%	26.33%
	EOR	88.33%	<b>89.67%</b>	88.00%
One-shot	OptiGuide	69.33%	<b>70.00%</b>	68.33%
	EOR	<b>95.33%</b>	93.67%	93.33%

Table 6: Accuracy under different LLMs in zero/one-shot setting with different *debug times*. The bold scores are the best in each row.

Method		Debug times	
		3	10
Zero-shot	GPT-4-1106-preview	<b>81.67%</b>	81.33%
	GPT-4-Turbo	88.33%	<b>89.00%</b>
One-shot	GPT-4-1106-preview	<b>87.67%</b>	<b>87.67%</b>
	GPT-4-Turbo	<b>95.33%</b>	<b>95.33%</b>

## A.5 PROMPT TEMPLATE FOR EVALUATING THE EXPLANATIONS

```

1116 EXPLANATIONS_EVALUATION = ""
1117
1118 **Role:** You are an expert in Operations Research evaluating
1119     ↪ explanations provided by three different models: 'EOR', 'OptiGuide
1120     ↪ ', and 'Standard'. Your role is to assess the quality of these
1121     ↪ explanations based on a user query.
1122
1123 **Input:**
1124     - User Query: {query}
1125     - Explanations by EOR: {eor}
1126     - Explanations by OptiGuide: {optiguide}
1127     - Explanations by Standard: {standard}
1128
1129 **Task:** For each model, provide three scores:
1130
1131     1). Score 1 - Explanation of Updated Code:
1132     How clearly does the explanation describe the code modifications made
1133     ↪ in response to the query? If no explanation of code updates is
1134     ↪ provided, score 0.

```

```

1134 18 2). Score 2 - Explanation of Query on Results:
1135 19 How well does the explanation clarify why and how the results changed
1136      ↪ due to the query? Focus on the depth of the explanation,
1137      ↪ particularly the quantitative reasoning behind the changes, not
1138      ↪ just a description of the result.
1139 20
1140 21 3). Score 3 - Overall Score:
1141 22 Based on the previous scores, assign an overall score (0-10)
1142      ↪ reflecting the combined quality and effectiveness of the
1143      ↪ explanation.
1144 23
1145 24 **Scoring Criteria:**
1146 25
1147 26 - Scores should range from 0 (poor) to 10 (excellent) for each
1148      ↪ category.
1149 27 - Consider the overall clarity, conciseness, and structure. The
1150      ↪ explanation should be easy to follow and understand.
1151 28
1152 29 **Output JSON Format:**
1153 30
1154 31 - Only return the results in the JSON format.
1155 32 - The keys should be 'EOR', 'OptiGuide', and 'Standard'.
1156 33 - Each key should have a list with the three scores (0-10).
1157 34 - Do not include any additional information or comments in the
1158      ↪ response.
1159 35
1160 36 Your evaluation will help determine which model provides the most
1161      ↪ effective and clear explanations for the given query. Good Luck!
1162 37
1163 38 --- Answer: ---
1164 39
1165 40 ""
1166 41
1167 42
1168 43
1169 44

```

## 1166 A.6 CASE STUDY

1168 The comparison of code and explanations generated by different models are shown in Figure 4 and  
 1169 5, respectively.

1170 As shown in Figure 4, EOR returns a code snippet in JSON format, OptiGuide provides a code  
 1171 snippet, and Standard returns the complete code for the entire problem. Although Standard achieves  
 1172 better experimental results than OptiGuide, its modification and return of the entire code signifi-  
 1173 cantly increases maintenance costs, making it challenging to track changes. OptiGuide provides  
 1174 markers to guide code modification in larger models, but experiments have shown that these mark-  
 1175 ers are sometimes removed, complicating efforts to identify what and where changes were made. In  
 1176 contrast, the code format proposed in our paper enables targeted updates at specific locations while  
 1177 ensuring markers remain intact. This not only reduces maintenance costs but also, when paired with  
 1178 clear explanations, simplifies the identification of both the location and nature of code updates.

1179 As shown in Figure 5, both Standard and OptiGuide generate explanations that mainly describe  
 1180 results and provide a brief analysis of changes in solution outcomes. Still, they lack quantitative  
 1181 analysis, especially concerning the impact of implicit query changes on the model. These expla-  
 1182 nations align more closely with the attribution explanation, as defined in our paper. In contrast,  
 1183 our approach emphasizes the justification explanation, which not only explains and analyzes code  
 1184 updates but also provides a quantitative analysis of the impact of constraint changes on the model.

1185  
 1186  
 1187

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

### Comparison of Code Generated by Different Models.

**Query:** How should the aircraft configuration be adjusted to meet demand if the company decides to limit the number of Type A aircraft operated to no more than 15 and the number of Type B aircraft to no more than 30?

#### EOR:

```
1 {
2   "ADD CONSTRAINT": "# Constraint: Limit number of Type A and Type
   ↳ B aircraft\nm.addConstr(aircraft_count['A'] <= 15, name
   ↳ =\ "MaxTypeAConstraint\ ")\nm.addConstr(aircraft_count['B']
   ↳ <= 30, name=\ "MaxTypeBConstraint\ ") "
3 }
```

#### OptiGuide:

```
1 "# Additional constraints for the maximum number of Type A and Type
   ↳ B aircraft\nm.addConstr(aircraft_count['A'] <= 15, name=\ "
   ↳ MaxAircraft_A\ ")\nm.addConstr(aircraft_count['B'] <= 30, name
   ↳ =\ "MaxAircraft_B\ ") "
```

#### Standard:

```
1 "import gurobipy as gp\nfrom gurobipy import GRB\n\n# Parameters
   ↳ Section Begin\n# Define model parameters\naircraft_types = ['
   ↳ A', 'B']\n\n# Passenger capacity per aircraft type\n
   ↳ npassenger_capacity = {\n 'A': 500,\n 'B': 200\n}\n\n#
   ↳ Operating cost per aircraft type (dollars)\n
   ↳ noperating_cost_per_aircraft = {\n 'A': 10000,\n 'B': 5000\n
   ↳ }\n\n# Minimum passenger demand\nmin_passenger_demand =
   ↳ 10000\n\n# Maximum number of aircraft for Type A and Type B\n
   ↳ nmax_aircraft_count_A = 15\nmax_aircraft_count_B = 30\n\n#
   ↳ Parameters Section End\n\n# Create a new Gurobi model\nm = gp
   ↳ .Model(\ "AirlineRevisedOptimization\ ")\n\n# Decision
   ↳ Variables Section Begin\n# Create decision variables for the
   ↳ number of large and small aircraft respectively\n
   ↳ aircraft_count = {\n 'A': m.addVar(vtype=GRB.INTEGER, name
   ↳ =\ "aircraft_A\ ")\n 'B': m.addVar(vtype=GRB.INTEGER, name=\ "
   ↳ aircraft_B\ ")\n}\n\n# Decision Variables Section End\n\n#
   ↳ Objective Function Section Begin\n# Set the objective
   ↳ function to minimize total operating cost\nm.setObjective(\n
   ↳ gp.quicksum(\n operating_cost_per_aircraft[t] *
   ↳ aircraft_count[t]\n for t in aircraft_types\n ),\n sense=GRB.
   ↳ MINIMIZE\n)\n\n# Objective Function Section End\n\n#
   ↳ Constraints Section Begin\n# Constraint: Meet the passenger
   ↳ transportation demand\nm.addConstr(\n gp.quicksum(
   ↳ passenger_capacity[t] * aircraft_count[t]\n for t in
   ↳ aircraft_types) >= min_passenger_demand,\n name=\ "
   ↳ PassengerDemandConstraint\ ")\n\n# Constraint: The number of
   ↳ Type A aircraft cannot exceed 15\nm.addConstr(\n
   ↳ aircraft_count['A'] <= max_aircraft_count_A,\n name=\ "
   ↳ TypeAConstraint\ ")\n\n# Constraint: The number of Type B
   ↳ aircraft cannot exceed 30\nm.addConstr(\n aircraft_count['B']
   ↳ <= max_aircraft_count_B,\n name=\ "TypeBConstraint\ ")\n\n#
   ↳ Constraints Section End\n\n# Solving the Model Section Begin\n
   ↳ n# Solve the model\nm.optimize()\n\n# Output the results\nif
   ↳ m.status == GRB.OPTIMAL:\n print(f\ "Optimized minimal total
   ↳ cost with new configuration: {round(m.ObjVal)} dollars\ ")\n
   ↳ for t in aircraft_types:\n print(f\ "Number of Type {t}
   ↳ aircraft: {aircraft_count[t].X}\ ")\nelse:\n print(\ "No
   ↳ optimal solution found.\ ")\n\n# Solving the Model Section End"
```

Figure 4: An example illustrating codes generated by different models.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

#### Comparison of Explanations by Different Models.

**Query:** How should the aircraft configuration be adjusted to meet demand if the company decides to limit the number of Type A aircraft operated to no more than 15 and the number of Type B aircraft to no more than 30?

**EOR:** **\*\*Explanation of Updated Code:\*\***

The updated code includes new constraints that implement operational limits on the number of Type A and Type B aircraft the airline can use. Specifically:

- A constraint has been added to limit the number of Type A aircraft to 15.
- A constraint has been added to limit the number of Type B aircraft to 30.

These additional constraints are motivated by possibly new operational limitations or strategic decisions, such as maintenance availability or cost considerations. By imposing these restrictions, the model must now find an optimal solution within a more constrained operational environment.

**\*\*Explanation of the Query on Results:\*\***

The query specifically asked to adjust aircraft configurations within newly defined limits. The original model found an optimal solution at a minimum cost of \$200,000. However, with the added constraints, the solution space is reduced - fewer configurations of aircraft are available to meet the passenger demand.

As a result of these constraints: The new optimal cost calculated by the model rose to \$215,000. This increment indicates that the more stringent operational limits forced the use of a costlier configuration. The query's impact on the results is substantial, as illustrated by a \$15,000 increase in operational costs. If rated on a scale from 1 to 10, the impact could be quantified around an 8. This rating reflects significant changes in operational costs and constraints.

The numerical changes in the model, calculated as 100, follows from a step-by-step process analyzing the modifications in the linear programming representation and their effects on the resulting computational graph. This method, grounded in operations research analysis, validates the substantial impact of the query-driven modifications.

In summary, the changes to the code (limiting aircraft types) constrained the model significantly, leading to a more expensive operational solution. The rationale behind integrating these specific constraints could be associated with real-world operational limits, reflecting the complex decision-making processes airlines undergo to balance cost, capacity, and operational feasibility.

**OptiGuide:** The optimization problem was re-solved with additional constraints limiting the number of Type A aircraft to no more than 15 and Type B aircraft to no more than 30. The new configuration resulted in a total operating cost of \$215,000. When compared to the original solution, where the total operating cost was \$200,000, this new configuration increased the operating costs by \$15,000 due to the more restrictive constraints on the number and types of aircraft deployed. This demonstrates the trade-offs between operational flexibility and cost.

**Standard:** Based on the updated constraints where the maximum number of Type A aircraft is limited to 15 and Type B aircraft to 30, the revised optimization solution results in a total operating cost of \$215,000. This is an increase from the original scenario, due to the restrictions on the number of aircraft which limited the company's ability to opt for the least expensive configuration to meet passenger demand. In this setup, the solution ensures passenger requirements are met while respecting the new operational limits placed on the aircraft types.

Figure 5: An example illustrating explanations generated by different models.