AUTOFORMULATION OF MATHEMATICAL OPTIMIZA-TION MODELS USING LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Mathematical optimization is fundamental to decision-making across diverse domains, from operations research to healthcare. Yet, translating real-world problems into optimization models remains a formidable challenge, often demanding specialized expertise. This paper formally introduces the concept of *autoformu*lation—an automated approach to creating optimization models from natural language descriptions for commercial solvers. We identify the three core challenges of autoformulation: (1) defining the vast, problem-dependent hypothesis space, (2) efficiently searching this space under uncertainty, and (3) evaluating formulation correctness (ensuring a formulation accurately represents the problem). To address these challenges, we introduce a novel method leveraging Large Language Models (LLMs) within a Monte-Carlo Tree Search framework. This approach systematically explores the space of possible formulations by exploiting the hierarchical nature of optimization modeling. LLMs serve two key roles: as dynamic formulation hypothesis generators and as evaluators of formulation correctness. To enhance search efficiency, we introduce a pruning technique to remove trivially equivalent formulations. Empirical evaluations across benchmarks containing linear and mixed-integer programming problems demonstrate our method's superior performance. Additionally, we observe significant efficiency gains from employing LLMs for correctness evaluation and from our pruning techniques.

028 029

031 032

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

1 INTRODUCTION

033 Mathematical optimization has long been a cornerstone of decision-making processes across various domains, from supply chain management (Bramel & Simchi-Levi, 1997) and healthcare resource allocation (Delgado et al., 2022) to portfolio optimization (Mokhtar et al., 2014). These problems are characterized by maximizing an objective function subject to constraints (Williams, 2013). Traditionally, optimization modeling follows a three-step process:
 gathering problem requirements, 037 typically expressed in unstructured formats and domain terminology; ► formulating these requirements into a formal mathematical model, including variables, constraints, and objective functions; implementing the model computationally using specialized modeling language for solution using 040 commercial solvers. These solvers (e.g. Gurobi (Gurobi Optimization, LLC, 2024), CPLEX (Cplex, 041 2009)), leverage sophisticated algorithms to tractably optimize a wide array of problems to global 042 optimality, including convex problems, (e.g. linear and quadratic programs), and certain non-convex 043 problems (e.g. mixed-integer linear programs) (Boyd & Vandenberghe, 2004). 044

Autoformulation. Despite major advances in solving algorithms over the past decades, the process of formulating optimization models still relies largely on human expertise (Karmarkar, 1984). Aut-046 oformulation aims to address this bottleneck by automating the formulation process, bridging the 047 gap between problem descriptions and formal mathematical models. This approach enhances time 048 and cost efficiency while enabling access for users without deep optimization expertise. At its core, autoformulation can be conceptualized as a search for an optimal formulation within a vast hypothesis space. However, this process is complicated by several challenges. For one, the hypothesis space 051 is highly problem-dependent, making it difficult to manually specify. Second, efficiently searching through this hypothesis space requires methods that can balance exploitation and exploration, par-052 ticularly given the uncertainty in the optimal formulation. Lastly, the search should be guided by a signal of formulation correctness, which is itself an ill-posed problem. While solvers can indi-

085

087

089

090 091

cate optimality gap and computational efficiency, evaluating that the formulated model accurately represents the intended real-world problem remains a distinct challenge.

A keen insight of this work is that recent advances in *Large Language Models* (LLMs) (Brown, 057 2020; Chowdhery et al., 2023) have opened new avenues for autoformulation. LLMs contribute 058 several crucial capabilities to this process: contextual understanding for nuanced interpretation of 059 problem descriptions, vast domain knowledge to incorporate relevant modeling techniques, and rea-060 soning capabilities to support approximate evaluation of formulation correctness. Recent works 061 (Ramamonjison et al., 2023; Xiao et al., 2023; AhmadiTeshnizi et al., 2024) have demonstrated the 062 promising potential of LLMs in autoformulation, laying important groundwork in this field. Build-063 ing upon these contributions, our work focuses on developing techniques for efficient, systematic 064 exploration and introducing novel methods for evaluating formulation correctness.

065 **Key considerations.** We conceptualize autoformulation as a search problem, leveraging optimiza-066 tion modeling's inherent hierarchical structure to efficiently explore the vast hypothesis space, 067 guided by feedback on formulation correctness. We decompose optimization modeling into hi-068 erarchical components and introduce a Monte-Carlo Tree Search (MCTS) method to incrementally 069 explore each component's formulation space (Coulom, 2006). We employ LLMs for two specialized roles: (1) as context-dependent hypothesis generators to produce component formulations at each level of the tree search; (2) as *evaluators* of formulation correctness, which is combined with solver-071 returned data to obtain a reward signal to guide exploration. To further improve search efficiency, 072 we introduce a pruning technique using Satisfiability Modulo Theories solvers, eliminating redun-073 dant hypotheses that are syntactically different yet functionally equivalent (Barrett & Tinelli, 2018). 074 Empirically, we observed that this significantly reduces search efforts expended on equivalent for-075 mulations. Our algorithm systematically explores the hypothesis space through multiple iterations, 076 producing a set of functionally distinct models, each scored using our evaluation function. 077

Contributions. Our main contributions are: A① We formally introduce *autoformulation* of mathematical optimization models, framing it as a search problem and identifying its core challenges. ②
We propose a novel approach integrating LLMs as hypothesis generators and evaluators within an MCTS framework, enabling efficient systematic exploration of the optimization model space. ③
Across two benchmarks containing linear and mixed-integer programming problems, we demonstrate our method's superior performance in formulating correct models, observing efficiency gains from pruning and LLM evaluation of formulation correctness.

2 AUTOFORMULATION: TOWARDS AUTOMATED OPTIMIZATION MODELING

Optimization modeling seeks to minimize an objective function subject to specific constraints on decision variables (Dantzig, 1990). The mathematical model can be expressed in a **general form**:

....

Minimize	$f(\mathbf{x})$			
subject to	$g_i(\mathbf{x}) \leq 0,$	$i=1,\ldots,I,$	((1)
	$h_j(\mathbf{x}) = 0,$	$j=1,\ldots,J.$		

Here $\mathbf{x} \in \mathcal{X}$ represents the vector of decision variables, and $\mathcal{X} \subseteq \mathbb{R}^{\ell} \times \mathbb{Z}^{k}$ is the domain of the problem for which the objective and constraints functions are all defined. Furthermore, $f : \mathcal{X} \to \mathbb{R}$ is the objective function to be minimized, $g_i : \mathcal{X} \to \mathbb{R}$ are inequality constraints, $h_j : \mathcal{X} \to \mathbb{R}$ are equality constraints, and I and J are the numbers of inequality and equality constraints respectively. The feasible region is the set of all possible points that satisfy the problem constraints: $\{\mathbf{x} \in \mathcal{X} \mid g_i(\mathbf{x}) \leq 0, \forall i \in [I], h_j(\mathbf{x}) = 0, \forall j \in [J]\}$.

Convex problems. An optimization problem is **convex** if f and $g_i \forall i \in [I]$ are convex, and $h_j \forall j \in [J]$ are affine. Convexity is significant as any local optimum of a convex problem is globally optimal, and specialized solvers can efficiently solve convex problems to global optimality using advanced algorithms (e.g., Gurobi (Gurobi Optimization, LLC, 2024), CVXPY (Diamond & Boyd, 2016)). For this reason, convexity is the widely accepted watershed between "easy" and "hard" problems. Before utilizing these solvers, the mathematical models are first represented in code as computational models, which are then passed to the solvers for optimization.



Figure 1: Illustration of autoformulation and its challenges. Autoformulation involves translating a problem description $d \in D$ into a mathematical model $m \in M$. This model is then transformed into a program that can be executed by solvers (e.g., Gurobi).

2.1 AUTOFORMULATION: PROBLEM DEFINITION

Boyd & Vandenberghe (2004) aptly recognized that "the challenge, and art, in using convex optimization is in recognizing and formulating the problem. Once this formulation is done, solving the problem is ... (almost) technology". While solver technology has significantly matured, the process of formulating optimization models remains largely human expertise driven. Building on this insight, we define *autoformulation* as the automated process of transforming natural language descriptions of real-world problems into formal optimization models. This process aims to bridge the gap between human-readable problem statements and computational models suitable for optimization solvers, thus automating the "challenge and art" of problem formulation.

Formal Definition

Let \mathcal{D} represent the space of natural language problem descriptions. Additionally, let \mathcal{M} and \mathcal{C} represent the space of all possible mathematical formulations and space of all possible computational models of optimization problems. Specifically, each $c \in \mathcal{C}$ is a computational representation (e.g. Python code), which includes choice of solving algorithm and its configurations. Given a problem $d \in \mathcal{D}$, autoformulation involves two transformations:

- 1. Mathematical Formulation $p_{\phi} : \mathcal{D} \to P(\mathcal{M})$: Transforming problem description into a precise mathematical formulation. Here, $P(\cdot)$ represents the space of probability distributions.
- 2. Computational Representation $p_{\psi} : \mathcal{M} \to P(\mathcal{C})$: Converting the mathematical formulation into computational formats suitable for solvers. This includes encoding the model in a programming framework and specifying or configuring an appropriate solving algorithm.

Autoformulator. Here, p_{ϕ} and p_{ψ} are models of each transformation, with ϕ and ψ their respective parameters. The complete autoformulation process can thus be represented as inferring the joint distribution $p_{\phi,\psi}(m, c \mid d) = p_{\psi}(c \mid m) \cdot p_{\phi}(m \mid d)$. We refer to any algorithm designed for the autoformulation problem as an *autoformulator*.

Objective. For a given problem d, autoformulation aims to find optimal mathematical and computational formulations that maximize an evaluation measure $Q(\cdot)$:

$$(m^*, c^*) \in \arg\max_{m \in \mathcal{M}, c \in \mathcal{C}} Q(m, c; d)$$
 (2)

Evaluation criteria. Here, Q assesses the quality of (m, c) relative to d, considering factors such as correctness, solvability, and efficiency. The quality of a formulation is primarily evaluated based on **formulation correctness**—how accurately it reflects the problem description. Given that a formulation is correct, two additional criteria come into play. **Optimality gap:** the gap between the value of the objective function at the solution and the optimal value. For example, convex problems and certain non-convex problems like mixed-integer linear programs (MILPs) can be solved efficiently to global optimality (i.e. zero optimality gap). **Computational efficiency:** evaluating the resource requirements and solution time of a particular model. These latter two criteria are only meaningful for formulations that correctly capture the problem description.

162 Challenges. A closer scrutiny of Eq. (2) reveals a few key challenges: 163 164 [C1] Problem-dependent hypothesis space: Given a problem d, we define the problem*dependent* hypothesis space of an autoformulator as $\mathcal{H}_{\phi,\psi}(d)$.^{*a*} This space encompasses 165 all plausible formulations, including correct, incorrect, or trivially redundant formulations. 166 Manually specifying this hypothesis is highly challenging, due to complex interdependen-167 cies, the vast number of formulations, and the problem-specific nature of the space. 168 **[C2]** Efficient search under uncertainty: Efficiently navigating the problem-dependent hypoth-169 esis space is challenging, as performant formulations (e.g. efficient and globally optimal) 170 can be sparse. This search involves managing two key sources of uncertainty: a) modeling 171 decisions: uncertainty in the optimal way to formulate the problem, b) problem ambigu-172 ity: uncertainty due to ambiguous requirements such as implicit or 'common-sense' con-173 straints (e.g. non-negativity or integer constraints for individual resources). An additional 174 complexity is **trivial model equivalence**—formulations that are identical but have minor 175 syntactic differences (e.g. functions 2x + 3y and 3y + 2x). Note, we term these 'trivial' to distinguish this type of equivalence from mathematically equivalent reformulation such 176 as converting a non-convex constraint to convex. Exploring these trivial variations ineffi-177 ciently can lead to overlooking more diverse and potentially valuable formulations. **[C3]** Model evaluation: While solvers can assess solvability and computational efficiency, they 179 cannot evaluate formulation correctness-whether the model accurately represents all re-180 quirements in the problem description. This lack of feedback on semantic correctness sig-181 nificantly complicates the search process, as an efficient solution to an incorrectly formu-182 lated problem is ultimately invalid. 183 ^{*a*}For notational simplicity, we omit the subscript (ϕ, ψ) , when the meaning is clear. 185 A few observations. Our definition of both steps as probabilistic distributions includes deterministic 187

mappings as a special case (i.e. Dirac delta distributions). This formalism recognizes the inherent uncertainty in autoformulation, stemming from language ambiguity, domain knowledge limitations, 188 and variability in modeling decisions. Additionally, the second transformation is assumed to be in-189 dependent of the problem description d. In general, the first step (formulating mathematical models, 190 p_{ϕ}) presents a significantly greater challenge than the second step (creating computational models, 191 p_{ψ}). The former requires deep domain understanding, the ability to abstract real-world complexities 192 into mathematical constructs, and creativity in choosing effective problem formulations. In contrast, 193 formulating computational models often follows more standardized patterns, with several commer-194 cial packages already offering automation in translating mathematical models into solver-compatible 195 code (Fourer et al., 1990). However, the second transformation can present distinct challenges that 196 introduce uncertainty, most notably through problem-specific choices and configurations for solvers.

197

199

200

201

202

203

Types of optimization problems. Additionally, we note that the challenges faced by an autoformulator fundamentally depends on the nature of the problem *d*, particularly its their convexity properties and reformulation possibilities. While some problems are naturally convex and enable direct solution for global optimality, others begin as non-convex and require the autoformulator to either identify equivalent convex reformulations or develop appropriate relaxation strategies that balance optimality with computational efficiency. In the interest of completeness, we provide a detailed categorization of different challenges presented by different problem types in App. E

- 204 205
- 206 207

3 LLM-ENHANCED MCTS FRAMEWORK FOR AUTOFORMULATION

208 Building upon our analysis of key challenges in autoformulation, we present a novel approach that 209 takes initial steps towards addressing them. In the scope of this work, our method primarily focuses 210 on evaluating formulation correctness—ensuring that autoformulated models accurately represent 211 the problem description. Moreover, we mainly focus on the first transformation, as we observed 212 that the second transformation was relatively straightforward for problems in available benchmarks. 213 Instead, we implement the second transformation using a custom deterministic parser (represented as parser, and detailed in App. A), which successfully parsed all available problems. By con-214 centrating on these fundamental aspects, we aim to catalyze the development of more advanced and 215 computationally efficient autoformulators that address more challenging problems of each type.

216 **Overview.** At a high level, our method integrates Large Language Models (LLMs) within a Monte-217 Carlo Tree Search (MCTS) framework, founded on two key insights: 1) leveraging the inherent 218 hierarchical structure of optimization modeling by decomposing the autoformulation process into 219 distinct stages, enabling systematic exploration using a tailored MCTS algorithm, and 2) employing 220 LLMs in two specialized roles: as dynamic hypothesis generators in each formulation stage, and as approximate evaluators of formulation correctness. This approach harnesses LLMs' extensive do-221 main knowledge and contextual understanding to implicitly create and search a problem-dependent 222 hypothesis space. This strategy bypasses the need for manual specification of the search space, which would be intractable due to the vast number of possible formulations. Furthermore, it utilizes 224 LLMs' reasoning capability to evaluate formulation correctness against problem descriptions, pro-225 viding an approximate yet meaningful signal to guide the search process. Note: In the interest of 226 space, we present detailed information about all prompts used in the algorithm in App. A, providing 227 only high-level details in the following subsections. 228

229 230

3.1 STRUCTURED DECOMPOSITION OF AUTOFORMULATION

Optimization modeling is inherently complex, involving multiple interconnected components. To manage this complexity and improve search efficiency, we propose a decomposition of the formulation process. This approach allows us to sequentially explore each model component rather than searching for entire formulations at once, potentially leading to more efficient search.

Specifically, we structurally decompose the autoformulation process into four distinct stages, each represented by m_i . The complete mathematical formulation is defined as $m = \bigoplus_{i=1}^{4} m_i$, where \bigoplus denotes the composition of model components: m_1 —parameters and decision variables, m_2 objective function, m_3 —equality constraints, and m_4 —inequality constraints. Given a problem description d, the joint distribution $p_{\phi,\psi}(c, m \mid d)$ is decomposed hierarchically:

240 241

246

$$\psi(c,m \mid d) = p_{\psi}(c \mid m) \prod_{i=1}^{4} p_{\phi}(m_i \mid m_{\leq i}, d)$$
(3)

Here, $p_{\phi}(m_i \mid m_{<i}, d)$ represents the sequential nature of mathematical formulation, where each component m_i depends on the *partial formulation* $m_{<i} = \bigoplus_{j=0}^{i-1} m_j$ (with $m_0 = \emptyset$) and the problem description d. Additionally, the term $p_{\psi}(c \mid m)$ represents the computational model's dependency on the completed mathematical formulation.

247 3.2 MCTS-BASED AUTOFORMULATOR

 p_{ϕ}

Having established a structured decomposition of the autoformulation process, we now address the challenge of efficiently navigating this hierarchical space. We employ an MCTS-based algorithm, which is particularly well-suited for exploring complex, hierarchical search spaces (Coulom, 2006). Our MCTS constructs a search tree of depth 4 to explore possible formulations, where each of the four levels corresponds to a component in our structured decomposition (m_1 to m_4). Nodes in this tree contain component formulations, and a complete formulation is represented by a path from the root to a terminal node, with each path yielding a unique formulation.

The MCTS algorithm iteratively builds the search tree through four key steps: \blacktriangleright selection, \triangleright expansion, \triangleright evaluation, and \triangleright backpropagation. For notational clarity, we denote a tree node as n and any of its child nodes as $n_{child} \in Child(n)$, where Child(n) is the set of all child nodes of n. We use \vec{n} to represent the *partial* formulation contained in the path from root to node n. For instance, \vec{n} for a node of depth 2 is the partial formulation containing the parameters, decision variables, and the objective function. Terminal nodes are denoted as n_t .

261 262

3.2.1 SELECTION

The selection step guides the search towards promising regions of the tree. Starting from the root, the algorithm recursively selects child nodes using the Upper Confidence Bound for Trees (UCT): $n_{child}^* = \arg \max_{n_{child} \in Child(n)} \left(V(n_{child}) + \omega \sqrt{\frac{\ln N(n)}{N(n_{child})}} \right)$ (Kocsis & Szepesvári, 2006). This process continues until reaching an unexpanded node. Here, n_{child}^* is the selected child node, $V(n_{child})$ is its estimated value, N(n) and $N(n_{child})$ are visit counts for the parent and child nodes respectively and ω is an exploration constant. This formula balances exploitation (first term, favoring high-value nodes) with exploration (second term, favoring less-visited nodes).

270 3.2.2 EXPANSION

272

273

274

Upon reaching an unexpanded node n, we generate its child nodes Child(n) through an expansion process. Unlike traditional MCTS, which typically operates within a predefined space, our expansion step explores an *undefined* hypothesis space of component formulations. To address this challenge, we employ LLMs as *dynamic hypothesis generators*, conditioned on the partial formulation constructed up to this level, to propose potential formulations for the next component. Our process involves two key steps: first, generating a diverse set of candidate formulations, and then pruning this set to remove trivial equivalences, thereby ensuring a manageable and meaningful search space.

282 Context-aware hypothesis generator. At node n, the LLM generates potential child 283 nodes (next component formulations) by conditioning on the partial formulation and prob-284 When expanding nodes at depth i, this is equivlem description: $p_{\phi}(n_{child} \mid \vec{n}, d).$ 285 alent to $p_{\phi}(m_i)$ $| m_{\langle i}, d\rangle.$ We represent formulations using JSON format, where keys are descriptive labels and values are mathematical expressions. For example, when generating possible inequality constraints, the LLM might return the formulation: 287 {"material_balance": $x_1 + x_2 \le 100$, "quality_requirement": $0.8x_1 + 0.6x_2 \ge 75$ }. The LLM is 288 queried through a structured prompt with three elements: (1) Problem description: the original 289 natural language problem description d; (2) Partial formulation: the current partial formulation 290 $m_{<i}$ in JSON format; (3) Level-specific instructions: guidelines for the current modeling stage, 291 including output format and relevant considerations. Note that we also request the LLM return po-292 tential formulations using the same dictionary format. For each node expansion, we sample $H \in \mathbb{N}$ 293 hypotheses from the LLM's distribution: $\widehat{Child}(n) = \{\widetilde{n}_{child}^{(h)} | \widetilde{n}_{child}^{(h)} \sim p_{\phi}(\cdot | \vec{n}, d), \forall h \in [H]\},\$ where $\tilde{n}_{child}^{(h)}$ represents the *h*-th candidate component formulation. 295

296 Approximate pruning. To ensure diversity in our search 297 space and avoid redundant explorations, we prune candidates 298 containing **trivially equivalent formulations**: Child(n)= 299 pruning(Child(n)). This process removes formulations with 300 minor syntactic differences that could lead to inefficient searches, 301 while preserving meaningful reformulations. For this purpose, we employ Satisfiability Modulo Theories (SMT) solvers to check pair-302 wise equivalence of formulations (Barrett & Tinelli, 2018). SMT 303 solvers determine the satisfiability of logical formulas with re-304 spect to background theories. We represent objective functions and 305 constraints as systems of equations or inequalities (where objec-306 tive functions form a single-equation system). For two such sys-307 tems S_1 and S_2 over variables x, we check the satisfiability of 308 $\neg(\forall x \ (S_1(x) \iff S_2(x))), \text{ where } \neg \text{ denotes negation. Unsat-}$ 309 isfiability of this formula proves equivalence, as it indicates there 310 exists no x such that the systems differ. Conversely, satisfiabil-311 ity indicates the systems are distinct. We apply this check to each 312 pair of candidate formulations in $\widehat{C}hild(n)$, pruning those deemed 313 trivially equivalent. This approach balances maintaining a diverse



Figure 2: **Expansion and node evaluation.** Expansion involves generating candidate formulations, which are then pruned based on functionality, and remaining formulations are assigned a normalized ranking score.

search space with computational efficiency. We detail the exact formulae used for SMT equivalencechecks in App. A.

316 We note that the satisfiability problem in SMT is not universally decidable (Monniaux, 2016). While 317 linear arithmetic over real and integer domains is generally decidable, mixed-integer domains or 318 non-linear functions may be undecidable, depending on specific problem properties. When the 319 solver cannot reach a conclusion, we assume the formulations are distinct. This approach trades 320 efficiency for thoroughness, potentially exploring some equivalent formulations while avoiding pre-321 mature pruning of unique components. Importantly, SMT solvers are only applicable to systems defined over the same variable domains. Consequently, we employ them solely for pruning re-322 dundant nodes in levels m_2-m_4 , where child nodes share decision variables. For level m_1 , which 323 involves defining different variable domains, we utilize LLMs as approximate checkers.

324 3.2.3 EVALUATION

326 After expanding a node, each newly created child node undergoes an initial evaluation to estimate 327 its value, guiding subsequent selection. This evaluation is non-trivial, as assessing the correctness of a partial formulation with respect to the original problem description is complex. Our method 328 employs an LLM-based ranking evaluation for each set of child nodes to provide more informed initial evaluations. Specifically, we rank the partial formulation from root to each child node, namely 330 $\{\vec{n}_{child} \mid n_{child} \in Child(n)\}$. These ranks are then center-normalized to [0, 1], with the middle 331 rank centered at 0.5. We denote this normalized score $s(\vec{n}_{child})$, which is used to initialize the 332 child node's value $V_{prior}(n_{child}) \leftarrow s(\vec{n}_{child})$. We note that this approach diverges from traditional 333 MCTS, which often uses uniform priors for expanded nodes. Here, the LLM assesses the formula-334 tion based on optimization principles and the specific problem context, potentially capturing aspects 335 such as formulation correctness, constraint feasibility, and alignment with the problem. 336

3.2.4 BACKPROPAGATION

337

338

345

346

Reward. Unlike conventional MCTS, which typically simulates the problem to a terminal state after expanding a child node, our approach continues expanding until a terminal node n_t is reached, where \vec{n}_t represents a complete formulation. This is computationally feasible due to our tree's limited depth of 4 levels, and provides more accurate rewards. We evaluate the complete formulation to obtain a *reward* $r(\vec{n}_t)$. We evaluate the complete formulation using a dual approach, combining assessments of both mathematical correctness and computational model's performance:

$$r(\vec{n}_t) = \mathbb{I}\left(E_{\text{solver}}^c(\text{parser}(\vec{n}_t)) = 1\right) \cdot E_{\text{LLM}}^m(\vec{n}_t; d) \tag{4}$$

where I is the indicator function. $E_{\text{LLM}}^m(\vec{n}_t; d)$ is the LLM's evaluation of the mathematical formulation's correctness, assessing how well it captures the problem requirements and constraints in *d*. $E_{\text{solver}}^c(\text{parser}(\vec{n}_t))$ is the solver's evaluation of the computational model's performance, providing a binary signal of whether the model is solved optimally. Note that $E_{\text{solver}}^c(\text{parser}(\vec{n}_t))$ is an imperfect signal, as an incorrectly formulated model could be solved to optimality despite not faithfully representing the original problem, highlighting the importance of our dual evaluation approach. The computational model *c* is derived using our custom deterministic parser: $c = \text{parser}(\vec{n}_t)$.

354 Here, it is inappropriate to use the same evaluation measure (based on ranking) as described be-355 fore, as the reward score would need to be comparable across all terminal nodes (across different subtrees). As such, we introduce a comparative evaluation method to obtain an LLM-evaluated cor-356 rectness score. This approach compares each formulation with a baseline, asking the LLM for its 357 preference. Specifically, the LLM returns a score $\in [0, 1]$, where < 0.5 values indicate preference 358 for the baseline formulation, and value greater than 0.5 favor the formulation in \vec{n}_t . Mathematically, 359 we represent this as $E_{\text{LLM}}^m(\vec{n}_t; d) \sim p_{\text{LLM}}(\cdot \mid \vec{n}_t, m_b; d)$, where m_b is the baseline formulation for 360 comparison. This baseline serves as a consistent reference point for all child nodes, enabling more 361 stationary and comparable evaluations of relative formulation correctness. In our implementation, 362 we generate the baseline formulation m_b through zero-shot prompting of the LLM.

Backpropagation. Following the reward calculation, we backpropagate this value to update the 364 statistics of all nodes along the trajectory. For each node in this path from root to terminal node n_t , 365 we apply the following updates: $V_{back}(n) \leftarrow \frac{V_{back}(n) \cdot N(n) + r(\vec{n}_t)}{N(n) + 1}$, $N(n) \leftarrow N(n) + 1$, $\forall n \in \vec{n}_t$. 366 Here, we increment the visit count N(n) by 1 and update the value V(n) with a weighted average 367 of its previous value and the new reward $r(\vec{n}_t)$. This backpropagation process ensures that the tree 368 gradually accumulates more accurate estimates of node values. These updated statistics then inform 369 the selection strategy in subsequent iterations. The value of the node used for selection is then 370 $V(n) = \lambda \cdot V_{prior}(n) + (1 - \lambda) \cdot V_{back}(n).$ 371

Summary. Our MCTS-based algorithm iterates through the aforementioned steps, progressively constructing and refining a tree of possible formulations. We execute this process for $T \in \mathbb{N}$ iterations, thoroughly exploring the space of potential models and identifying promising formulations. The final output is a set of $M \in \mathbb{N}$ functionally distinct optimization models, where $M \leq T$. Each model is defined by a unique trajectory through the tree. Formally, we express the overall algorithm as: $\{(m^{(i)}, c^{(i)}, r^{(i)})\}_{i=1}^{M} = \text{MCTS}_{\text{LLM}}(d)$. The superscript *i* indexes the functionally distinct formulation, and $r^{(i)}$ is the estimated value/reward of the corresponding terminal node.

378 4 **RELATED WORK** 379

Advances in LLMs. Recent works have demonstrated the substantial potential of Large Language 380 Models (LLMs) in solving complex reasoning tasks, including language understanding (Hendrycks 381 et al., 2021), commonsense reasoning (Brown, 2020), logical reasoning (Wei et al., 2022; Yao et al., 382 2024), mathematical problem-solving (Lewkowycz et al., 2022), and coding tasks (Chen et al., 383 2021). Of particular relevance are studies employing LLMs in optimization and search tasks, such 384 as Bayesian Optimization (Liu et al., 2024b), prompt optimization (Guo et al., 2023), evolutionary 385 optimization (Yang et al., 2024; Liu et al., 2024a), and symbolic program refinement (Madaan et al., 386 2024), as well as research exploring the integration of LLMs with planning algorithms (Huang et al., 387 2022; Zhao et al., 2024; Hao et al., 2023a; Zhou et al., 2024).

388 Autoformulation. Ramamonijson et al. (2023) introduced an early competition focused on translat-389 ing natural language descriptions of linear programming problems into mathematical formulations. 390 The competition involved two tasks: tagging problem entities and predicting formulations in prede-391 fined formats. Entries primarily used pre-LLM NLP models tailored for these tasks, which lacked 392 the ability to generalize beyond the given formats. More recently, Xiao et al. (2023) and Ahma-393 diTeshnizi et al. (2024) explored multi-agent LLM frameworks for optimization model formulation, 394 where LLM agents generated complete formulations in each iteration, refining them locally in sub-395 sequent steps. Our approach differs by breaking down the formulation process into key stages and 396 using MCTS to systematically explore the formulation space. Additionally, we guide the search with a composite reward function that combines solver feedback with LLM evaluation. 397

398 399

400 401

403

5 **EXPERIMENTS**

The experiments aim to evaluate the performance of the autoformulation model across a diverse set 402 of problems (Sec. 5.1). In addition, we study two key factors of our framework: (a) the use of LLMbased evaluation to ensure the correctness of formulations (Sec. 5.2), and (b) enhanced efficiency 404 by reducing search space in equivalent formulations (Sec. 5.3). Finally, the experiments examine 405 failure modes to better understand the limitations of the model (Sec. 5.4).

406 We use NL4OPT (Ramamonjison et al., 2023) and IndustryOR as benchmarks to evaluate our ap-407 proach. NL4OPT is a standard dataset for operations research tasks, primarily focusing on linear 408 programming problems. We used the filtered version presented in (Tang et al., 2024), consisting of 409 244 problems. IndustryOR, with 100 problems covers linear, integer, mixed-integer, and non-linear 410 programming across three difficulty levels, offers a more complex challenge. This complexity is cru-411 cial for testing the robustness of our method in navigating larger, ambiguous solution spaces. Since 412 neither benchmark provides enough information to fully evaluate optimization models-lacking 413 ground truth for decision variables, objectives, and constraints—we adopt an approximately correct framework, measuring accuracy using a ground truth value found in (Tang et al., 2024), following 414 the same methodology. All experiments were performed using GPT4o-mini, except for NL4OPT, 415 where GPT4-1106 was used for fair comparison. 416

417

419

418 5.1 BENCHMARK COMPARISON

For a comprehensive evaluation, we cover all relevant previous 420 approaches, including: Reflexion (Shinn et al., 2023), Chain-421 of-Experts (Xiao et al., 2023), and OptiMUS (AhmadiTeshnizi 422 et al., 2024) in Table 1. These approaches employ agents to it-423 eratively enhance both the mathematical formulation and the 424 optimization outcomes, delivering strong results. For consis-425 tency, we report the performance of these methods based on 426 both GPT-4 outputs. We also include methods based on fine-427 tuning through large synthetic data such as ORLM (Tang et al., 428 2024). Note that our approach is not directly comparable to

Table 1: Benchmark comparison							
Method	NL4OPT	IndustryOR					
Finetuned methods							
ORLM-LLama-3-8B	85.7%	38.0%					
Methods	based on GI	PT-4					
Standard	47.3%	28.0%					
Reflexion	53.0%	-					
Chain-of-Experts	64.2%	-					
OptiMUS	78.8%	-					
MCTS (iter-1)	85.24%	35.0%					
MCTS (iter-3)	92.21%	42.0%					
MCTS (All)	92.62%	48.0%					

429 previous approaches since it captures all possible interpretations of problems. We report the accuracy of the first solution found by MCTS, as well as the accuracy of the first three solutions and 430 after 16 rollouts. At each step, we limit the number of generated solutions to 10, with a maximum 431 of three children per node. A higher computational budget could potentially yield more solutions.

432 5.2 EVALUATING CRITIC CAPABILITIES OF LLM 433

434 We evaluate the critic capabilities of the LLM, which in our method is used to search the for-435 mulation space. We study backpropagation (global score) and prior score (local score) inde-436 pendently. (1) Getting a global score. The score in Eq. (4) compares complete formula-437 tions. Unlike prior scoring, we compare full formulations without sharing partial structures. To isolate this component, we extract all solutions from IndustryOR that contain correct an-438 swers and compare them to wrong answers. Our goal is to show that correct solutions are 439 generally preferred over incorrect ones when comparing formulations directly. Using this ap-440 proach, we obtain a point biserial correlation coefficient of 0.48 with a significant p-value of 441 2.0681×10^{-3} . We also compare this ranking method with a scoring method from the litera-442 ture Zhang et al. (2024b), where formulations are scored from 1 to 100, yielding a correlation of 443 0.2325 with a p-value of 1.1185×10^{-1} , which is not significant. (2) Getting a Local Scores. 444 To measure prior scores, we use a simple greedy evaluation. To 445 isolate the effect of the MCTS framework, we first construct a tree 446 using Depth-First Search (DFS) with up to three children per node, 447 and assign prior rewards via the LLM. For evaluation, we only 448 consider trees containing the ground truth and measure the success 449 rate using a greedy approach. Figure 3 shows the comparison with





EFFICIENCY OF DETECTING FUNCTIONALLY EQUIVALENT FUNCTIONS 5.3

random prior scoring, illustrating that the chances of reaching the

ground truth decrease as the number of tree leaves increases.

455 We evaluate the efficiency of our framework by analyzing the number of generated for-456 mulations at each step and the remaining formulations after two key filtering stages: (1) 457 grouping equivalent formulations using SMT to eliminate redundancy, and (2) selecting the top three solutions based on rankings provided by the LLM. The top of Fig. 458 shows the initial number of generated formulations as a percentage (ten per experiment in 459 practice) and tracks the results of each filtering step across all problems in IndustryOR. 460

The equivalent grouping step is the most impactful, re-461 ducing the solutions by a factor of five accross all steps. 462 The final filtering step, which selects the top-ranked for-463 mulations, discards very few solutions-a positive out-464 come as we aim to retain viable formulations. Based 465 on these results, the bottom of the figure contrasts effi-466 ciency with a non-hierarchical method, where errors in 467 earlier stages compound and affect later steps. The re-468 sults demonstrate that by the time the formulations reach 469 the inequality constraint step, efficiency increases a thousandfold. In other words, a non-hierarchical approach 470 without filtering would have required a thousand more 471 simulations to arrive at the same number of solutions. 472



Figure 4: (Top) Number of nodes filtered. (Bottom) Estimated efficiency.

5.4 FAILURE MODES 474

450

451 452 453

454

473

475 We examine the failure modes of our method by (1) analyzing 476 its search capabilities for detecting the ground truth and (2) 477 identifying the problems where the model itself fails. For the 478 first point, Fig. 5 illustrates the search evolution of our MCTS 479 relative to the number of rollouts, highlighting our method's 480 ability to benefit from additional exploration (more iterations) 481 to discover more correct solutions. Despite the excellent capa-482 bilities of our MCTS in finding novel solution, we notice that 483 as our MCTS is sensible to recommendation of best candidates based on greedy approaches as shown in the Figure 8 of Ap-484



pendix. This indicates the better critic methods are necessary to improve recommendation of best 485 candidates.

We analyze failure modes based on the difficulty and problem character-istics in IndustryOR. Accuracy is measured by whether the mode finds the correct solution. In both grouping methods, the model with the high-est accuracy also produces less entropic (less dense) trees. Similarly, when grouped by difficulty, the medium category exhibits lower accuracy and higher tree entropy. Our method does not show a significant weakness against any specific type of problems. For example, ORLM performs well in IPs but is weaker in MIPs (see App. D).

Category	Accuracy	Entropy
Grou	iped by Diffic	ulty
Easy	0.6750	1.9638
Medium	0.2895	3.0408
Hard	0.5000	2.7337
Gi	rouped by Typ	ре
IP	0.5484	1.6536
LP	0.4167	2.1716
MIP	0.5161	3.3201

6 DISCUSSIONS

In this paper, we formally define the problem of autoformulation for
 mathematical optimization models, establishing objectives, evaluation metrics, and a categoriza-

tion of problem types based on their challenges to autoformulators. We introduce a novel approach that frames autoformulation as a search problem and effectively leveraging the hierarchical struc-ture of optimization modeling. Our method integrates LLMs as hypothesis generators and evalua-tion functions of formulation correctness within an MCTS framework, systematically exploring the vast hypothesis space of possible formulations. The introduction of hypothesis pruning using SMT solvers further enhances efficiency by eliminating redundant formulations. Empirical evaluations across linear and mixed-integer programming benchmarks demonstrate our method's superior per-formance in formulating correct models, with notable efficiency gains from pruning and LLM-based correctness evaluation.

Future works. Looking ahead, we envision autoformulation as an exciting domain where LLMs can significantly augment human experts. Future research directions include developing collaborative frameworks to synergize with human expertise, exploring advanced LLM-based methods such as retrieval-augmented generation (Lewis et al., 2020) and test-time compute (Lightman et al., 2024). A particularly promising direction lies in fine-tuning models to enhance autoformulator capabilities, with special emphasis on process-supervised learning for hierarchical modeling steps (Lightman et al., 2023; Wan et al., 2024), which aligns naturally with the structured decomposition inherent in optimization formulation. An important consideration for future work is the potential correlation bias from using the same LLM for both generation and evaluation; while our composite evaluation strategy helps mitigate this through solver feedback and comparative ranking, developing specialized evaluation models or ensemble-based approaches could provide more robust assessment. To support these advancements, the development of large-scale, diverse benchmarks encompassing var-ious problem types and complexities, particularly those requiring intricate reformulations, will be crucial.

540 REFERENCES

549

550

- Ali AhmadiTeshnizi, Wenzhi Gao, and Madeleine Udell. Optimus: Scalable optimization modeling
 with (mi) lp solvers and large language models. *arXiv preprint arXiv:2402.10172*, 2024.
- Farid Alizadeh and Donald Goldfarb. Second-order cone programming. *Mathematical programming*, 95(1):3–51, 2003.
- Anastasios G Bakirtzis and Pandelis N Biskas. A decentralized solution to the dc-opf of intercon nected power systems. *IEEE Transactions on Power Systems*, 18(3):1007–1013, 2003.
 - Clark Barrett and Cesare Tinelli. Satisfiability modulo theories. *Handbook of model checking*, pp. 305–343, 2018.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron
 Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference on machine learning*, pp. 531–540. PMLR, 2018.
- 555 Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- Julien Bramel and David Simchi-Levi. *The Logic of Logistics: Theory, Algorithms, and Applications for Logistics Management.* Springer, 1997.
- David Brandfonbrener, Simon Henniger, Sibi Raja, Tarun Prasad, Chloe Loughridge, Federico Cassano, Sabrina Ruixin Hu, Jianang Yang, William E. Byrd, Robert Zinkov, and Nada Amin. Vermets: Synthesizing multi-step programs using a verifier, a large language model, and tree search, 2024. URL https://arxiv.org/abs/2402.08147.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Alphamath almost zero: process supervision without process, 2024a. URL https://arxiv.org/abs/2405.03553.
- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Step-level value preference optimization
 for mathematical reasoning, 2024b. URL https://arxiv.org/abs/2406.10858.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- ⁵⁷⁸ Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pp. 72–83. Springer, 2006.
- IBM ILOG Cplex. V12. 1: User's manual for cplex. International Business Machines Corporation, 46(53):157, 2009.
- George B Dantzig. Origins of the simplex method. In *A history of scientific computing*, pp. 141–151.
 1990.
- Erwin J Delgado, Xavier Cabezas, Carlos Martin-Barreiro, Víctor Leiva, and Fernando Rojas. An
 equity-based optimization model to solve the location problem for healthcare centers applied to
 hospital beds and covid-19 vaccination. *Mathematics*, 10(11):1825, 2022.
- Matthew DeLorenzo, Animesh Basak Chowdhury, Vasudev Gohil, Shailja Thakur, Ramesh Karri, Siddharth Garg, and Jeyavijayan Rajendran. Make every move count: Llm-based high-quality rtl code generation using mcts. *arXiv preprint arXiv:2402.03289*, 2024.
- 593 Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.

594 Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, 595 and Jun Wang. Alphazero-like tree-search can guide large language model decoding and train-596 ing. arXiv preprint arXiv:2309.17179, Sep 2023a. Available at https://arxiv.org/abs/ 597 2309.17179. 598 Xidong Feng, Ziyu Wan, Muning Wen, Ying Wen, Weinan Zhang, and Jun Wang. Alphazerolike tree-search can guide large language model decoding and training. arXiv preprint 600 arXiv:2309.17179, 2023b. 601 602 Robert Fourer, David M Gay, and Brian W Kernighan. A modeling language for mathematical 603 programming. Management Science, 36(5):519-554, 1990. 604 605 Benjamin Fuchs. Application of convex relaxation to array synthesis problems. *IEEE Transactions* 606 on Antennas and Propagation, 62(2):634-640, 2013. 607 Michel X Goemans and David P Williamson. Improved approximation algorithms for maximum cut 608 and satisfiability problems using semidefinite programming. Journal of the ACM (JACM), 42(6): 609 1115-1145, 1995. 610 611 Andrea J Goldsmith and Pravin P Varaiya. Capacity, mutual information, and coding for finite-state 612 markov channels. IEEE transactions on Information Theory, 42(3):868-886, 1996. 613 614 Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. Connecting large language models with evolutionary algorithms yields powerful 615 prompt optimizers. arXiv preprint arXiv:2309.08532, 2023. 616 617 Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2024. URL https://www. 618 gurobi.com. 619 620 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting 621 Hu. Reasoning with language model is planning with world model. In The 2023 Conference on 622 Empirical Methods in Natural Language Processing, 2023a. 623 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 624 Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992, 625 2023b. 626 627 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob 628 Steinhardt. Measuring massive multitask language understanding. In International Conference 629 on Learning Representations, 2021. 630 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot 631 planners: Extracting actionable knowledge for embodied agents. In International conference on 632 machine learning, pp. 9118–9147. PMLR, 2022. 633 634 Narendra Karmarkar. A new polynomial-time algorithm for linear programming. In Proceedings of 635 the sixteenth annual ACM symposium on Theory of computing, pp. 302–311, 1984. 636 637 Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In European conference on machine learning, pp. 282–293. Springer, 2006. 638 639 Javad Lavaei and Steven H Low. Zero duality gap in optimal power flow problem. *IEEE Transactions* 640 on Power systems, 27(1):92–107, 2011. 641 642 H. Lebret and S. Boyd. Antenna array pattern synthesis via convex optimization. *IEEE Transactions* 643 on Signal Processing, 45(3):526–532, 1997. 644 645 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-646 tion for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems, 33: 647 9459-9474, 2020.

663

679

680

681

682 683

684

685 686

687

688

689

690

691

- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
- Qingyao Li, Wei Xia, Kounianhua Du, Xinyi Dai, Ruiming Tang, Yasheng Wang, Yong Yu, and
 Weinan Zhang. Rethinkmcts: Refining erroneous thoughts in monte carlo tree search for code generation. *arXiv preprint arXiv:2409.09584*, 2024.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=v8L0pN6EOi.
- Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew-Soon Ong. Large language models as
 evolutionary optimizers. In 2024 IEEE Congress on Evolutionary Computation (CEC), pp. 1–8.
 IEEE, 2024a.
- Tennison Liu, Nicolás Astorga, Nabeel Seedat, and Mihaela van der Schaar. Large language models
 to enhance bayesian optimization. In *The Twelfth International Conference on Learning Representations*, 2024b. URL https://openreview.net/forum?id=00xotBmGol.
- Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated process supervision. *arXiv preprint arXiv:2406.06592*, 2024.
- ⁶⁷⁴ Zhi-Quan Luo and Shuzhong Zhang. Dynamic spectrum management: Complexity and duality.
 ⁶⁷⁵ *IEEE journal of selected topics in signal processing*, 2(1):57–73, 2008.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Mazura Mokhtar, A Shuib, and D Mohamad. Mathematical programming models for portfolio optimization problem: A review. *International Journal of Mathematical and Computational Sciences*, 8(2):428–435, 2014.
 - David Monniaux. A survey of satisfiability modulo theory. In Computer Algebra in Scientific Computing: 18th International Workshop, CASC 2016, Bucharest, Romania, September 19-23, 2016, Proceedings 18, pp. 401–425. Springer, 2016.
 - Rindranirina Ramamonjison, Timothy Yu, Raymond Li, Haley Li, Giuseppe Carenini, Bissan Ghaddar, Shiqi He, Mahdi Mostajabdaveh, Amin Banitalebi-Dehkordi, Zirui Zhou, et al. Nl4opt competition: Formulating optimization problems based on their natural language descriptions. In *NeurIPS 2022 Competition Track*, pp. 189–203. PMLR, 2023.
 - Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023.
- ⁶⁹³ Zhengyang Tang, Chenyu Huang, Xin Zheng, Shixi Hu, Zizhuo Wang, Dongdong Ge, and Benyou
 ⁶⁹⁴ Wang. Orlm: Training large language models for optimization modeling. *arXiv preprint* ⁶⁹⁵ *arXiv:2405.17743*, 2024.
- Ziyu Wan, Xidong Feng, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training. In *Forty-first International Conference on Machine Learning*, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.

702 703	H Paul Williams. Model building in mathematical programming. John Wiley & Sons, 2013.
704	Ziyang Xiao, Dongxiang Zhang, Yangjun Wu, Lilin Xu, Yuan Jessica Wang, Xiongwei Han, Xiaojin
705	Fu, Tao Zhong, Jia Zeng, Mingli Song, et al. Chain-of-experts: When Ilms meet complex opera-
706	tions research problems. In The Twelfth International Conference on Learning Representations,
707	2023.
708	Yuxi Xie Anirudh Goval Wenvue Zheng Min-Ven Kan Timothy P Lillicran Kenii Kawaguchi
709	and Michael Shieh Monte carlo tree search boosts reasoning via iterative preference learning
710	arXiv preprint arXiv:2405.00451, 2024.
711	
712	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun
713	Chen. Large language models as optimizers. In <i>The Twelfth International Conference on Learning</i>
714	<i>Representations</i> , 2024. UKL https://openreview.net/forum?id=Bb4vGowELI.
715	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
716	Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Ad-
717	vances in Neural Information Processing Systems, 36, 2024.
718	Roy D Yates A framework for unlink power control in cellular radio systems <i>IFFF Journal on</i>
719	selected areas in communications, 13(7):1341–1347, 1995.
720	
721	Wei Yu and Raymond Lui. Dual methods for nonconvex spectrum optimization of multicarrier
722	systems. <i>IEEE Transactions on communications</i> , 54(7):1310–1322, 2006.
723	Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts*: Llm
724	self-training via process reward guided tree search. arXiv preprint arXiv:2406.03816, 2024a.
725	$\mathbf{D}'_{\mathbf{T}}$
726	Di Zhang, Xiaoshui Huang, Dongzhan Zhou, Yuqiang Li, and Wanli Ouyang. Accessing gpt-4 level
727	mathematical orymptad solutions via monte carlo tree sen-renne with nama-5 80, 20240.
720	Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge for
729	large-scale task planning. Advances in Neural Information Processing Systems, 36, 2024.
731	Ruizhe Zhong, Xingho Du, Shixiong Kai, Zhentao Tang, Siyuan Xu, Hui-Ling Zhen, Jianye Hao
732	Oiang Xu, Mingxuan Yuan, and Junchi Yan. Llm4eda: Emerging progress in large language
733	models for electronic design automation. arXiv preprint arXiv:2401.12224, 2023.
734	
735	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language
736	agent tree search unities reasoning acting and planning in language models, 2025.
737	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Lan-
738	guage agent tree search unifies reasoning, acting, and planning in language models. In Forty-first
739	International Conference on Machine Learning, 2024.
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
750	
752	
754	
755	

756 A Additional Details on Method

A.1 FORMULATION EQUIVALENCE CHECKS

760 SMT solvers offer a powerful approach for verifying equivalence between various components of 761 optimization models (Barrett & Tinelli, 2018). These tools can rigorously check if different for-762 mulations of objective functions, sets of equality constraints, or sets of inequality constraints are logically equivalent. By encoding the components as logical formulas within appropriate theories 763 764 (such as linear arithmetic), SMT solvers can determine if the formulations are satisfiable under the same conditions. For objective functions, the solver can check if the difference between two func-765 tions is always zero across the feasible region. This is formally described in Eq. (5). For constraint 766 sets, it can verify if they define identical feasible regions by checking that each constraint in one set 767 is implied by the other set and vice versa, formally described in Eqs. (6) and (7). This approach not 768 only ensures the correctness of model transformations or reformulations but also aids in identifying 769 redundant constraints and simplifying complex models. However, the effectiveness of SMT solvers 770 in this context depends on the nature of the optimization problem, as nonlinear or highly complex 771 formulations may pose challenges for current solvers.

1. For objective functions $f^{(i)}$ and $f^{(j)}$:

772

774

775 776

787 788

789

790 791

792 793

794

796 797 798

799

$$\mathsf{Equivalent}(f^{(i)}, f^{(j)}) \iff \forall \mathbf{x} \in \mathcal{X}, f^{(i)}(\mathbf{x}) = f^{(j)}(\mathbf{x})$$
(5)

2. For sets of equality constraints $g^{(i)} = \{g_k^{(i)}\}_k^K$ and $g^{(j)} = \{g_l^{(j)}\}_l^L$:

$$\operatorname{Equivalent}(g^{(i)}, g^{(j)}) \iff \forall \mathbf{x} \in \mathcal{X}, (\bigwedge_{k} g_{k}^{(i)}(\mathbf{x}) = 0) \iff (\bigwedge_{l} g_{l}^{(j)}(\mathbf{x}) = 0)$$
(6)

3. For sets of inequality constraints $h^{(i)} = \{h_k^{(i)}\}_k^K$ and $h^{(j)} = \{h_l^{(j)}\}_l^L$:

$$\operatorname{Equivalent}(h^{(i)}, h^{(j)}) \iff \forall \mathbf{x} \in \mathcal{X}, (\bigwedge_{k} h_{k}^{(i)}(\mathbf{x}) \le 0) \iff (\bigwedge_{l} h_{l}^{(j)}(\mathbf{x}) \le 0)$$
(7)

A.2 PROMPT DESIGN

Template instruction

```
I have a problem in operational research:
------
###PROBLEM DESCRIPTION###
------
I have the following formalization:
formalization_dict = {"parameters": {}, "decision_variables2: {},
"objective": {}, "equality_constraints": {}, "inequality_constraints":
{}}
```

Parameters template (root node generation)

```
800
       You are an optimization modeling expert. Complete
       formalization_dict based on the problem description, you should
801
       complete the "parameters" field, which consists of assigning
802
       constants to descriptive variable names.
803
       Only complete "parameters" and nothing else.
                                                      Follow these
804
       guidelines:
805
806
       1. Your primary responsibility is to define all the parameters
807
       from the problem description that will later be used to define
808
       decision variables, the objective, and constraints (both
809
       equality and inequality).
```

810 2. You may include additional parameters in a format suitable 811 for facilitating the subsequent tasks of defining decision 812 variables, the objective function, and constraints. 813 3. For parameters that involve multiple indices (e.g., x[i] or 814 x[i,j]), use the most appropriate data structure, such as lists, 815 dictionaries, or dictionaries with tuple keys, to represent 816 them. 4. For each parameter, include a clear, descriptive comment 817 818 explaining its meaning. 5. Ensure that the parameter names (keys) are descriptive and 819 intuitive. 820 821 Return only the python dictionary update (i.e., 822 formalization_dict["parameters"] = ...) following the 823 described requirements. 824 825 826 **Decision Variables Template (Depth == 1)** 827 828 You are an optimization modeling expert. Complete only the 829 "decision_variables" field within the "formalization_dict" based 830 on the provided problem description. 831 Ensure the decision_variables comprehensively cover all 832 essential elements to accurately model the optimization problem. 833 834 Each key-value pair in the dictionary must adhere to the following structure: 835 836 <key>: { 837 "description": <description>, 838 "type": <type>, 839 "iteration_space": <space> } 840 841 The structure should meet these requirements: 842 843 1. Each <key> represents a decision variable that will later 844 be used to implement the objective, equality, and inequality constraints in a Python program. 845 2. Replace <key> with a symbolic name representing the decision 846 variable. Ensure that each <key> represents a distinct decision 847 variable with a unique symbolic name. 848 3. Replace <description> with a detailed explanation of the 849 role of the decision variable in the optimization model. 850 4. Replace <type> with a string representing the Gurobi 851 variable type (e.g., GRB.INTEGER), as this will be used to 852 create the variable via Gurobi's addVar function. 853 If the decision variable is indexed, replace <space> with a 5. 854 string representing Python for-loop using list comprehension 855 syntax to represent the index space. For this, assume direct access to these parameter variables (i.e., avoid using 856 parameters[variables] syntax). 857 6. If the variable is not indexed, set <space> to None. 858 7. If the variable is indexed, do not write the index in the 859 symbol (do not put the index when writing <key>). 860 8. You are encouraged to create decision variables that are 861 general. If two decision variables represent the same concept 862 write them as one key, creating an appropriate iteration space. 863

```
864
       Return only the Python dictionary update (i.e.,
865
       formalization_dict["decision_variables"] = ...) following the
866
       described requirements.
867
868
       Objective Template (Depth == 2)
869
870
       You are an optimization modeling expert. Complete only the
871
       "objective" field within the "formalization_dict" based on the
872
       provided problem description.
873
       Do not complete any other fields. Follow these requirements:
874
       1. Write the objective function mathematically using decision
875
       variables.
876
       2. Preface the key-value pair with a Python comment explaining
877
       the rationale behind the objective. DO NOT make a commentary
878
       inside the mathematical description.
879
       3.
           Use parameter-defined variables instead of hard-coded
880
       values. Assume direct access to these parameter variables
881
       (i.e., avoid using parameters[variables] syntax).
882
       4. The dictionary key must be 'min' or 'max', reflecting the
883
       nature of the objective (minimization or maximization).
884
       5. The dictionary value must be a string representation of the
       objective function based on the problem description, written in
885
       valid Python syntax.
886
887
       Return only the Python dictionary update (i.e.,
888
       formalization_dict["objective"] = "max":
                                                 ... or
889
       formalization_dict["objective"] = "min": ...) following the
890
       described requirements.
891
892
       Equality Constraints Template (Depth == 3)
893
894
       You are an optimization modeling expert. Complete the
895
       formalization_dict by filling in the equality_constraints field
896
       based on the problem description and the decision variables
897
       provided.
       These constraints include border constraints, initialization,
898
       and equality constraints derived from the problem description.
899
       Do not complete the "inequality_constraints" field. Follow
900
       these requirements:
901
902
       1. Descriptive constraints: Each key in the dictionary should
903
       represent a unique, clearly named constraint, with the value
904
       being a string that describes the corresponding mathematical
905
       equality using "==".
906
       2.
          Parameter Variables: Use parameter-defined variables
907
       instead of hard-coded values. Assume direct access to these
908
       parameter variables (i.e., avoid using parameters[variables]
       syntax).
909
          Indexed Variables: For indexed decision variables, indicate
       3.
910
       the index within brackets (e.g., x[i]).
911
          Handling Multiple Constraints: For similar constraints that
       4.
912
       repeat across indices or variables, use Python for loops and
913
       list comprehensions for efficient representation.
914
       5. String mathematical description: Note, the value
915
       (mathematical description) should be a single string. DO NOT
916
       use .join() or anything else. Even if it represents multiple
917
       constraints using a for loop.
```

918 6. No Inequality Constraints: Only define equality 919 constraints. Inequality constraints will be handled separately 920 by a subsequent expert. 921 7. Comments: Include a Python comment before each key-value 922 pair, explaining the rationale behind the constraint. 923 924 Return only the Python dictionary update (i.e., 925 formalization_dict["equality_constraints"] = ...) following 926 these requirements. Important: If the problem contains only inequality 927 constraints and no equality constraints, return: 928 formalization_dict["equality_constraints"] = {None: None}. 929 This will signal the need to focus on inequality constraints 930 in subsequent modeling steps. 931 932 933 Inequality Constraints Template (Depth == 4) 934 You are an optimization modeling expert. Complete the 935 formalization_dict by adding the inequality_constraints field 936 based on the problem description. Follow these requirements: 937 938 1. Descriptive constraints: Each key in the dictionary should 939 represent a unique, clearly named constraint, with the value 940 being a string that describes the corresponding mathematical 941 inequality. 942 2. Parameter Variables: Use parameter-defined variables 943 instead of hard-coded values. Assume direct access to these 944 parameter variables (i.e., avoid using parameters[variables] syntax). 945 3. Indexed Variables: For indexed decision variables, indicate 946 the index within brackets (e.g., x[i]). 947 4. Handling Multiple Constraints: For similar constraints that 948 repeat across indices or variables, use Python for loops and 949 list comprehensions for efficient representation. 950 5. String mathematical description: Note, the value 951 (mathematical description) should be a single string without 952 using join or anything else. Even if it represents multiple 953 constraints using a for loop. 954 6. Inequality Constraints Only: Include only inequality 955 constraints. Exclude any constraints already covered under equality_constraints. 956 7. Comments: Include a Python comment before each key-value 957 pair, explaining the rationale behind the constraint. 958 959 Return only the Python dictionary update (i.e., 960 formalization_dict["inequality_constraints"] = ...) following 961 these requirements. 962 Important: Think carefully of inequality constraints that 963 are not explicit in the problem description that should be 964 considered. If after thinking you conclude the problem contains 965 only equality constraints and no inequality constraints, return: 966 formalization_dict["inequality_constraints"] = {None: None}. 967 968 **Group Decision Variables template** 969

- Objective:

970

972 As an expert in optimization modeling, your role is to evaluate 973 multiple sets of decision variables provided for an operations 974 research problem. You are responsible for determining if two or 975 more sets of decision variables should be grouped together based 976 on their equivalency from an optimization perspective. 977 - Task Breakdown: 978 979 Your grouping decision is critical for assisting a subsequent 980 optimization expert, who will define the objective function, equality constraints, and inequality constraints for each group. 981 To facilitate this process, follow these precise guidelines: 982 983 - Equivalency Criteria: 984 985 1. Same Objectives and Constraints: Two sets of decision 986 variables should be grouped together if they result in the 987 definition of the same objective function, equality constraints, 988 and inequality constraints, even if the variable names differ. 989 2. Conceptual Equivalency: Variable sets should be grouped 990 together if, despite having different variable names, they define the same underlying concepts that ultimately lead 991 to identical objectives and constraints (both equality and 992 inequality). 993 Non-Equivalency Conditions: Two sets of decision variables 3. 994 should not be grouped together if they lead to differences in 995 any of the following: Objective function, Equality constraints, 996 Inequality constraints. 997 4. Naming Convention Irrelevance: The names of the decision 998 variables are irrelevant for grouping purposes. Only the 999 functional impact of the variables on the objective function 1000 and constraints should be considered. If two sets of variables lead to the same results, group them together, even if the names 1001 differ. 1002 1003 By following these guidelines, you will help ensure that 1004 decision variable sets are clearly classified for the next 1005 expert in the process. 1006 1007 Please list your clusters as follows: 1008 ### 1009 groups = $\{$ 1010 1: group_1, 1011 . . . , 1012 n: group_n} 1013 ### 1014 Where group_i is a python list containing the names (string) of 1015 all the set of decision variables that are equivalent. One set 1016 of decision variables can only belong to one group. The list 1017 should consider at least one element. 1018 1019 Important: Think carefully STEP BY STEP about your grouping 1020 decision, then conclude your assessment using the structured 1021 format provided above. 1022 Here are the current solutions: 1023 1024 solutions = {} 1025

```
1026
       Ranking Template to obtain prior rewards
1027
1028
       You are an expert in optimization modeling. Using the
1029
       formalization_dict as your current progress, you are tasked with
       selecting the optimal #VARIABLE# from the provided options.
1030
1031
       Please follow these steps:
1032
1033
       1.
           Carefully evaluate each potential #VARIABLE#.
1034
1035
       2.
           Rank the variables from best to worst based on their
1036
       suitability.
1037
1038
       Present your rankings in the following format:
1039
       ###
1040
       rank = \{
1041
       1: solution_1,
1042
       ...,
1043
       n: solution_n}
1044
       ###
1045
       Where:
1046
       - solution_1 represents the best #VARIABLE#.
1047
       - solution_n represents the least suitable #VARIABLE#.
1048
1049
       Important: Think carefully STEP BY STEP about your ranking
1050
       decision. Then conclude by listing the solutions in string
       format as structured above.
1051
1052
       Here are the possible solutions:
1053
1054
       solutions = \{\}
1055
1056
      A.3 EVALUATION METRICS
1057
```

1058 EXPECTED CORRECTNESS

1063

1073 1074

1075 1076

Mathematical Definition: Expected Correctness (EC) quantifies the probability of reaching a correct leaf node when starting from a given node and following a probabilistic policy based on the node rankings *R*. It is computed recursively:

• For a leaf node n:

$$EC(n) = \begin{cases} 1, & \text{if } n \text{ is correct} \\ 0, & \text{if } n \text{ is incorrect} \end{cases}$$

• For an intermediate node n with children $\{c_i\}$:

$$\operatorname{EC}(n) = \sum_{i} p(c_i \mid n) \times \operatorname{EC}(c_i)$$

where the probability of choosing child c_i is given by:

$$p(c_i \mid n) = \frac{R(c_i)}{\sum_j R(c_j)}$$

1077 Intuitive Explanation: Expected Correctness measures the likelihood that, by following the ranking-based probabilities at each decision point, we will eventually arrive at a correct solution.
 1079 It reflects the overall effectiveness of the ranking system in guiding the search process toward correct outcomes.

1080 PRECISION AND RECALL

Mathematical Definition:

 • **Precision** at a node is the ratio of relevant children (those leading to correct leaves) to all retrieved children, averaged over all positions:

$$Precision = \frac{1}{n} \sum_{k=1}^{n} \frac{\text{Relevant}_k}{k}$$

where Relevant_k is the number of relevant children in the top k positions.

• **Recall** at a node is the ratio of relevant children retrieved to the total number of relevant children, averaged over all positions:

$$\operatorname{Recall} = \frac{1}{n} \sum_{k=1}^{n} \frac{\operatorname{Relevant}_{k}}{\operatorname{Total Relevant}}$$

Intuitive Explanation: Precision indicates how well the top-ranked children (based on *R*) correspond to those that lead to correct solutions. High precision means most top-ranked choices are relevant. Recall measures the ability of the ranking to capture all relevant children among its selections. High recall implies that the ranking method successfully identifies most of the correct paths.

1103 NORMALIZED DISCOUNTED CUMULATIVE GAIN (NDCG@K)

Mathematical Definition: NDCG@k evaluates the quality of the ranking up to position k, ac-1105 counting for the position of relevant items:

1. Compute DCG@k (Discounted Cumulative Gain):

$$\mathsf{DCG}@k = \sum_{i=1}^{k} \frac{2^{\mathsf{rel}_i} - 1}{\log_2(i+1)}$$

where rel_i is the relevance score at position i (1 if the child leads to a correct leaf, 0 otherwise).

2. Compute IDCG@k (Ideal DCG) by ordering the children perfectly:

IDCG@
$$k = \sum_{i=1}^{k} \frac{2^{\operatorname{rel}_{i}^{*}} - 1}{\log_{2}(i+1)}$$

 $NDCG@k = \frac{DCG@k}{IDCG@k}$

where $\operatorname{rel}_{i}^{*}$ is the ideal ordering of relevance scores.

3. NDCG@k is the ratio:

1124 Intuitive Explanation: NDCG@k assesses not just whether relevant children are present in the top k positions but also how highly they are ranked. It rewards rankings that place relevant children earlier, reflecting the practical importance of quickly finding correct solutions.

1128 TOP-RANK SUCCESS RATE

Mathematical Definition: The Top-Rank Success Rate is the probability that always choosing the top-ranked child (highest *R*) at each decision point leads to a correct leaf:

Top-Rank Success Rate =
$$\frac{\text{Number of times top-path leads to correct leaf}}{\text{Total number of root nodes}}$$

Intuitive Explanation: This metric evaluates the effectiveness of a greedy strategy based solely on the ranking *R*. A high success rate indicates that the top-ranked paths are reliable guides to correct solutions, simplifying decision-making in the tree traversal.

B COMPARISON TO RELATED WORKS: MCTS METHODS

1143

1144

11/5

1138

Table 3: Comparison of LLM+MCTS variants. This table presents a comparative analysis of LLM+MCTS methods, highlighting the methodological innovation of our method. Specifically, we compare (1) expansion mechanisms for generating child nodes; sources of evaluation for (2) newly expanded nodes; and (3) terminal nodes; whether (4) evaluations employ comparative evaluations; and finally the (5) search space. $V, \mathcal{E}, \mathcal{N}$ represent evaluation signals obtained from LLMs, external environments (e.g. test cases in coding tasks), and comparisons with other solutions, respectively. $\mathcal{T}, \mathcal{P}, \Sigma$ represent space of natural language thoughts, formal programs, and tokens, respectively.

Method	Expansion mechanism	Initial node evaluation	Reward used in backpropagation	Evaluation via comparisons?	Search space
RAP (Hao et al., 2023b)	-	ν	V V	X	T
LATS (Zhou et al., 2023)	-	ν	\mathcal{V}, \mathcal{E}	X	τ
MCTSr (Zhang et al., 2024b)	Self-refine	ν	v	Х	τ
VeriGenMCTS (DeLorenzo et al., 2024)	Filter by functionality	-	E	X	Σ
RethinkMCTS (Li et al., 2024)	Rethink	\mathcal{V}, \mathcal{E}	\mathcal{V}, \mathcal{E}	Х	τ
VerMCTS (Brandfonbrener et al., 2024)	Joint Expansion-Evaluation	ε	ε	X	$ \mathcal{P} $
AlphaZero like MCTS (Feng et al., 2023b; Zhang	-	ν	\mathcal{V}, \mathcal{E}	Х	Τ
DPO like MCTS (Xie et al., 2024; Chen et al., 2024b)	-	\mathcal{V}	\mathcal{V}, \mathcal{E}	X	τ
Ours	SMT Pruning	$\mathcal{V}, \mathcal{E}, \mathcal{N}$	$\mathcal{V}, \mathcal{E}, \mathcal{N}$	1	$M_{1:4}$

1153 1154

1163

1168

1169

1170

1171

1172

Recent advancements in MCTS have demonstrated that combining MCTS with LLMs significantly enhances reasoning capabilities by refining the thinking process (Zhang et al., 2024b). Complementary improvements have been achieved using Process Reward Models, which further optimize reasoning using training with step-wise reward information (Feng et al., 2023a). In parallel, the integration of MCTS and LLMs has been applied to code generation (Zhong et al., 2023; Brandfonbrener et al., 2024; Li et al., 2024). We present a detailed comparison against these related works in Tab. 3.

1162 Our MCTS framework introduces three key innovations specifically tailored to autoformulation:

- Structured hierarchical search: We leverage the inherent structure of optimization modeling to decompose the search space. Unlike conventional MCTS approaches, which assume fixed search spaces, our hierarchical organization of search spaces both reduces search complexity and increases formulation diversity.
 - 2. **SMT-based pruning:** Our analysis shows that 80% of generated formulations are trivially equivalent (see Fig. 4). By integrating SMT solvers to prune these redundant formulations, we achieve a 400x improvement in search efficiency, avoiding exponential growth in search complexity.

3. **Comparative formulation evaluation:** We introduce novel pairwise comparative evaluation for assessing formulation correctness, which is distinct from the standard approach where LLMs evaluate solutions in isolation. This comparative framework enables more reliable preference-based evaluation to improves search efficiency.

- 1173 1174
- 1175 1176

C REFORMULATION STRATEGIES

1177 1178

Here, we enumerate commonly used strategies to reformulate the problem. Detailed descriptions can be found in (Boyd & Vandenberghe, 2004, Chapter 4.1.3).

- Change of variables
 - Transformation of functions
- 11831183Slack variables
- Eliminating equality constraints
- Adding equality constraints
- Optimizing over some variables
 - Epigraph form

1188 D ADDITIONAL RESULTS

1190 1191

D.1 OBTAINING LOCAL SCORES (PRIOR RANKING).

1193 1194 1195

1192

1196 We evaluated our node ranking method on the IndustryOR dataset, aggregating metrics across all 1197 problems. The assessment used five key metrics: Average Expected Correctness, Average Precision, Average Recall, NDCG@3, and Top-Rank Success Rate. Average Expected Correctness measures 1198 the likelihood of reaching a correct leaf by following a probabilistic policy based on RR scores. Pre-1199 cision and Recall assess the accuracy and coverage of ranking correct subtrees. NDCG@3 evaluates 1200 ranking quality within the top three positions, and Top-Rank Success Rate gauges the success of a 1201 greedy strategy that always selects the top-ranked child. Table in Figure 6 presents detailed results, 1202 with the top section showing Precision and Recall by formulation step and the bottom displaying 1203 aggregated results across all nodes. 1204

Our results shows an Average Expected Correctness of 1205 0.6692 indicating a significant likelihood of reaching cor-1206 rect solutions using the probabilistic policy based on $\mathbb R$ 1207 scores. The high Average Precision (0.8911) and Recall 1208 (0.8643) show that our ranking method effectively priori-1209 tizes correct subtrees while identifying most of them. The 1210 NDCG@3 of 0.6776 reflects strong ranking performance 1211 within the top three positions, crucial when resources for 1212 exploration are limited. Finally, the Top-Rank Success 1213 Rate that 73% of the time selecting the top-ranked child leads to correct solutions. These results validate the ro-1214



bustness of our ranking approach in guiding the search process toward correct solutions.

- 1216 1217
- 1218
- 1219

1220

1221 1222

1223

D.2 ASSESSING THE CRITIC EVALUATION CAPABILITIES OF LLMS

1224 Prior rewards. To evaluate whether the prior re-1225 ward serves as a reliable indicator of a good node 1226 during node selection, we compute the *normalized* 1227 accumulated prior rewards. This is the sum of all prior rewards from the visited nodes, normalized by 1228 dividing the total by 4 (corresponding to 4 steps). 1229 Using this metric, we analyze the results across the 1230 IndustryOR dataset, grouping the outcomes based 1231 on whether the formalization produces a score equal 1232 to the ground truth or not. The results of this 1233 study is presented in Figure 7. The mean accumu-1234 lated reward for successful formalizations (predic-1235 tion equals ground truth) is 0.744 ($\sigma_{succ} = 0.133$), 1236 notably higher than the mean for failed formaliza-1237 tions, which is 0.651 ($\sigma_{fail} = 0.122$). An independent t-test confirms that this difference is statistically



Figure 7: **Top**. Distribution of *accumulated prior rewards* of correct prediction vs failed prediction. **Bottom** Statistical test comparing both groups.

significant (t = 4.86, p < 0.001), indicating that successful formalizations consistently yield higher rewards. Additionally, a point-biserial correlation analysis shows a moderate positive correlation (r = 0.344, p < 0.001) between accumulated rewards and success, suggesting that higher accumulated rewards are associated with an increased likelihood of arriving at the successful formalization.

D.3 ADDITIONAL COMPARISONS AGAINST ORLM

Table 4: Comparison between GPT4-omini. For the MCTS approach the accuracy if the ground truth is found by the MCTS search.

Method		Difficulty			Q	uestion Ty	pes	
	Easy	Medium	Hard	LP	NLP	IP	MIP	Others
ORLM-LLaMA-3-8B	57.5%	20.0%	35.0%	36.1%	0.0%	61.2%	19.3%	0.0%
MCTS (ours)	67.50%	28.95%	50.00%	41.67 %	0.0%	54.84%	51.61%	0.0%

D.4 GREEDY SEARCH AFTER 16 ROLLOUTS MCTS





1297					,
1298		.	Corrector de la décorr		
1299		Formalization Step	Correlation	p-value	
1300		Decision Variables	0.1985	0.1159	
1301		Objective	0.2143	0.2470	
1302		Equality Constraints	0.3711	0.0004	
1303		Inequality Constraints	0.3290	0.0033	
1304	Table 5. Spec	amon's Completion and n	values for diff.	mant famaalia	ation stans
1305	Table 5: Spea	arman's Correlation and p	-values for diffe	erent formaliz	ation steps.
1306					
1307	Experimental setup. 1	o analyze partial model e	valuation, we c	ompute corre	lations between inter-
1308	mediate node evaluation	ns in the MCIS and grou	nd-truth correct	iness scores.	For each intermediate
1309	noue, we define its grou	ind-truth correctness as th	e percentage of	correct rear r	iodes in its subtree.
1310	Observations. Analysis	s of Spearman's correlation	ons (Tab. 5) reve	eals two key p	patterns:
1311	1 Evolution volution	and the survey of the stand to	Completion et		
1312	1. Evaluation robusting $(n = 0.1085)$ t	o inequality constraints (m	-0.3200 with	th correspond	ing going in statistical
1313	significance. This al	igns with intuition. deepe	= 0.5250), wi	n more compl	ete formulation infor-
1314	mation, enabling mo	re accurate evaluation.	i nodes contan	i more compi	ete formulation milor
1315	2. Component-specific	c evaluation challenges.	Earlier compor	ents show we	eaker correlations due
1316	to interdependencies	s between modeling elem	ents. For insta	ance, evaluati	ng decision variables
1317	in isolation is challe	enging without understand	ding their role	in objectives	and constraints. This
1318	increased uncertaint	y at earlier stages reflects	the inherent dif	ficulty in asse	ssing partial formula-
1319	tions without full co	ntext.			
1320	These for these sectors			-1	that maintain dimension
1321	exploration paths partic	core the importance of m	erarchical sear	ionals are we	aker
1222	exploration pauls, partic	cutarry in carry stages with	cic evaluation s	ignais are we	akcı.
1323					
1324					
1326					
1327					
1328					
1329					
1330					
1331					
1332					
1333					
1334					
1335					
1336					
1337					
1338					
1339					
1340					
1341					
1342					
1343					
1344					
1345					
1346					
1347					
1348					
1349					

1296 D.5 Additional Analysis on Partial Formulation Evaluations

1351				
1352		Factor	Optimality gap	Computational efficiency
1252	Type I (origi-	Formulation	Minimal impact: any correct (equiv-	Medium impact: choice of variable representation
1555	nally convex)		alent) formulation can achieve global	(e.g. structure-preserving formulations) can improve
1354			optimality	solution time
1255		Solver	Minimal impact: most commercial	Medium impact: specialized solvers for specific
1555			solvers can achieve similar optimality	problem structures (LP, QP, SOCP) can be faster
1356			gaps	
1357	Type II (non-	Formulation	High impact: correct, convexified re-	High impact: reformulation complexity affects solu-
1007	convex, but con-		formulation enables achievable global	tion time, generally convex reformulations are solved
1358	vexifiable)		optimality	faster
1359		Solver	Medium impact: solver ability to han-	High impact: solver must efficiently handle the spe-
			dle reformulated structures affect solu-	cific structure of reformulation
1360			tion quality	
1361	Type III (non-	Formulation	High impact: quality of relaxation	Medium impact: relaxation complexity affects so-
1000	convex, requir-		directly affects optimality gap and	lution time, involving trade-offs between relaxation
1302	ing relaxation		bounds tightness	tightness and computational efficiency
1363		Solver	High impact: solver abilities on re-	Medium impact: solving relaxed problems may re-
1964			laxed problem affects solution quality	quire specialized solvers, although general purpose
1304				solvers are roughly comparable

Table 6: Impact of **formulation** and **solver** on optimality gap and computational efficiency.

E CATEGORIZATION OF AUTOFORMULATION CHALLENGES BY **OPTIMIZATION PROBLEM STRUCTURE**

1372 The exact challenges faced by an autoformulator depends on the nature of the problem d. Here, 1373 we provide a categorization of optimization problems and their characteristics. To help elucidate 1374 different types of problems, we introduce two concepts. First, we define the set of correct formu**lations** for a problem $\mathcal{M}(d) \subset \mathcal{M}$ as the set of all equivalent formulations that *correctly* model a 1375 problem d. Second, we introduce the set of **original forms** $\mathcal{M}_{o}(d) \subseteq \mathcal{M}(d)$ —the set containing 1376 the natural representations of the problem, typically the initial models an optimization expert would 1377 create. This is a set, as it can contain trivially equivalent formulations. Finally, we partition the set 1378 \mathcal{M} into the set of convex problems \mathcal{M}_{conv} and the set of non-convex problems \mathcal{M}_{nonc} . 1379

1380 1381

1382

1383

1384

1385

1386

1387

1389

1390

1391

1350

1365

1367

1369

1370 1371

> 1. **Type I problems.** These are problems where the original form is inherently convex, namely $\mathcal{M}_o(d) \subseteq \mathcal{M}_{conv}$. Examples include certain resource allocation problems that can be naturally formulated as linear programs. The challenge of solving Type I problems is to ensure that the problem is correctly represented (formulation correctness, i.e. $\mathcal{H}(d) \cap \mathcal{M}_{\alpha}(d) \neq \emptyset$), which would entail that it can be efficiently solved to global optimality.

- 2. Type II problems. These are problems where the original form is non-convex, but can be reformulated into an equivalent convex problem, namely $\mathcal{M}_o(d) \subseteq \mathcal{M}_{nonc}$ but $\mathcal{M}(d) \cap \mathcal{M}_{conv} \neq$ \emptyset . In addition to formulation correctness, another challenge of solving Type II problems 1388 is to ensure the autoformulator can identify and apply appropriate reformulation strategies (e.g. change of variables) to transform the non-convex into an *equivalent* convex form, namely $\mathcal{H}(d) \cap (\mathcal{M}(d) \cap \mathcal{M}_{conv}) \neq \emptyset$ (see App. C for discussion). For such problems, evaluation extends beyond correctness to include the ability to achieve global optimality through reformulation.
- 1392 3. Type III problems. These are problems where the original form is non-convex and cannot be 1393 reformulated into a convex problem, namely $\mathcal{M}(d) \subseteq \mathcal{M}_{nonc}$. In such cases, there are two 1394 general options: a) solve the non-convex problem using general-purpose algorithms (e.g. gradient descent), or b) relax into a convex problem that approximates, but is not equivalent to, the original 1395 problem (e.g. semidefinite relaxation of a Max-Cut problem (Goemans & Williamson, 1995)).

1398 A crucial nuance here is that mathematically equivalent models, even when both are convex, can 1399 exhibit vastly different computational complexities. An example of this is quadratic programming 1400 and second-order cone programming (SOCP) reformulations of the same problems (Alizadeh & Goldfarb, 2003). Although mathematically equivalent, SOCP formulations often allow for more 1401 efficient solution methods. Therefore, computational efficiency is an important evaluation metric 1402 across all three problem types, significantly impacting practical utility of model formulations. In 1403 App. F, we provide concrete examples to illustrate each type of optimization problems.

1404 F ILLUSTRATIVE EXAMPLES OF PROBLEM CATEGORIZATION

In this section, we provide examples of canonical problems in engineering and machine learing that
 belong to each of the identified problem types. Specifically:

- **Type I**: Problems that have a precise mathematical model, which is convex in its original form. Examples are provided in App. F.1.
- Type II: Problems that have a precise mathematical model, which is nonconvex in its original form but can be reformulated as a convex problem (sometimes additional assumptions are needed). Examples are provided in App. F.2.
- Type III: Problems that have a precise mathematical model, which is nonconvex in its original form but can be *relaxed* to a convex problem (sometimes additional assumptions are needed). The difference from Type II is that the convex relaxation is *not* equivalent to the original problem. Examples are provided in App. F.3.
- We also provide examples of problems where the problem description is inherently ambiguous in App. F.4
- 1421 F.1 EXAMPLES OF **TYPE-I** PROBLEMS

1423 [P1] Mutual Information Maximization (Mutual-Information)

• Reformulation strategies: None.

- **Difficulty in reformulation:** Easy.
- Difficulty in reformulation: Easy.
 Difficulty in solving the reformulated problem: Easy.

Mutual information is a quantity that measures the divergence between two random variables, with applications in wireless communications (Goldsmith & Varaiya, 1996) and in data science (Belghazi et al., 2018). Here we describe it in the context of maximizing Shannon capacity in wireless communications.

We consider a discrete memoryless channel with an input random variable $X \in \{1, ..., \ell\}$, an output random variable $Y \in \{1, ..., y\}$, and channel transition matrix $P \in \mathbb{R}^{y \times \ell}$ with the element on the *j*-th row and the *i*-th column being $p_{ji} = \text{prob}(Y = j \mid X = i)$.

Input
$$X \longrightarrow$$
 Transition Probability $P \longrightarrow$ Output Y

1437 1438 1439

1440

1448

1449 1450 1451

1457

1436

1420

1422

1424

1425

Our goal is to choose the optimal probability distribution of input X, denoted $x \in \mathbb{R}^{\ell}$ with $x_i = \text{prob}(X = i)$, in order to maximize the mutual information between input X and input Y

$$I(X;Y) = \sum_{i=1}^{\ell} \sum_{j=1}^{y} x_i p_{ji} \log_2 \frac{p_{ji}}{\sum_{k=1}^{\ell} x_k p_{jk}}$$

1445 The optimal value of the problem is called Shannon capacity.

1446 This problem is convex in its original form:

Maximize
$$\sum_{i=1}^{\ell} \left(\sum_{j=1}^{y} p_{ji} \log_2 p_{ji} \right) x_i - \sum_{j=1}^{y} \left(\sum_{i=1}^{\ell} p_{ji} x_i \right) \log_2 \left(\sum_{i=1}^{\ell} p_{ji} x_i \right)$$
subject to $x_i \ge 0, \quad i = 1, \dots, \ell,$
$$\sum_{i=1}^{\ell} x_i = 1.$$
(8)

1456 [P2] Power control for maximum throughput (PC-MaxRate-decoupled)

• Reformulation strategies: None.

1458 Difficulty in reformulation: Easy.

1459 • Difficulty in solving the reformulated problem: Easy. 1460

1461 We consider the problem of allocating a unit amount of total power over ℓ independent frequency or 1462 temporal wireless channels (Yu & Lui, 2006; Luo & Zhang, 2008). Each channel has a channel gain g_i and a noise power σ_i , resulting in a channel gain to noise ratio of $\alpha_i = g_i / \sigma_i$. 1463

1464 Our goal is to choose the optimal power allocation, denoted $x \in \mathbb{R}^{\ell}_{\perp}$ with x_i being the power 1465 allocated to channel *i*, in order to maximize the total throughput. 1466

Maximize $\sum_{i=1}^{\ell} \log_2 \left(1 + \alpha_i x_i \right)$

 $x_i > 0, \quad i = 1, \ldots, n.$

(9)

This problem is convex in its original form: 1467

1468 1469

1470 1471

1472

1473 1474

1475 1476

1477

1479

[P3] Power control to satisfy SINR requirements with minimum power (PC-MinPower)

subject to $\sum_{i=1}^{\ell} x_i \leq 1$,

 Reformulation strategies: Transformation of function. 1478

• **Difficulty in reformulation:** Easy (straightforward observation).

• Difficulty in solving the reformulated problem: Easy (the reformulated problem is LP). 1480

1481 We consider the problem of determine the transmit power of ℓ pairs of transceivers. They operate 1482 in the same frequency at the same time, hence causing interference to each other. The problem data 1483 is a channel gain matrix $\mathbf{G} \in \mathbb{R}^{\ell \times \ell}$, where g_{ij} is the channel gain from transmitter j to receiver i, 1484 the noise power vector $\sigma \in \mathbb{R}^{\ell}$ with σ_i as the noise power at receiver *i*, and the minimum SINR 1485 requirement vector $\gamma \in \mathbb{R}^{\ell}$ with γ_i as the minimum SINR required by transceiver *i*. 1486

Our goal is to choose the transmit power, denoted $x \in \mathbb{R}^{\ell}_{\perp}$ with x_i being the power of transmitter 1487 *i*, in order to minimize the total transmit power while satisfying the SINR requirements of each 1488 transceiver (Yates, 1995). 1489

1490 This problem is non-convex in its original form:

1494

1495 1496

1497

1498

 $\begin{array}{ll} \text{Minimize} & \sum_{i=1}^{\ell} x_i \\ \text{subject to} & \frac{g_{ii}x_i}{\sum_{j\neq i}g_{ij}x_j + \sigma_i} \geq \gamma_i, \quad i = 1, \dots, \ell. \end{array}$ (10)

But it is not hard to observe that the constraints of SINR requirements can be reformulated as linear constraints, resulting in a LP:

Minimize
$$\sum_{i=1}^{\ell} x_i$$

subject to $g_{ii}x_i \ge \gamma_i \left(\sum_{j \ne i} g_{ij}x_j + \sigma_i\right), \quad i = 1, \dots, \ell.$

$$(11)$$

1506

1507 [P4] Beamforming to minimize sidelobes (Beamform-MinSidelobe)

 Reformulation strategies: Transformation of functions, slack variables. 1509

- Difficulty in reformulation: Easy (standard techniques were used). 1510
- Difficulty in solving the reformulated problem: Medium (the reformulated convex problem is 1511 a second-order cone program (SOCP)).

We consider the problem of an antenna array with ℓ elements in a 2-D space, with the position of kth element denoted by the coordinate (x_k, y_k) . Given the complex-valued weight w_k (i.e., a voltage or current phasor) of each element, the output of the array at direction θ is

$$G(\mathbf{w};\theta) = \sum_{k=1}^{\ell} w_k e^{\mathbf{i}(x_k \cos \theta + y_k \sin \theta)}$$

1518

1515 1516 1517

1524

1525

1519 where $\mathbf{i} = \sqrt{-1}$ and $G(\mathbf{w}; \theta)$ is complex-valued.

¹⁵²⁰ Our goal is to achieve certain gain at the target direction θ_{tar} , while minimizing the energy radiated at directions $\theta_1, \ldots, \theta_m$ outside the target area (Lebret & Boyd, 1997).

This problem is convex:

 $\begin{array}{ll} \text{Minimize} & \max_{i=1,\ldots,m} |G(\mathbf{w};\theta_i)| \\ \text{subject to} & G(\mathbf{w};\theta^{\text{tar}}) = 1. \end{array}$ (12)

It is convex because the objective function is the point-wise maximum of convex functions (i.e., norms) and the constraint is linear.

However, the maximum operator makes the objective function non-differentiable. To make it amenable for standard solvers, we would reformulate it as the following SOCP:

1531

1532

1533

1534

1537

1548

1549

1550 1551

1552 1553 1554 Minimize tsubject to $|G(\mathbf{w}; \theta_i)| \le t$, i = 1, ..., m, (13) $G(\mathbf{w}; \theta^{\text{tar}}) = 1$.

[P5] Optimal power flow in direct current (DC-OPF)

• Reformulation strategies: None.

1538 • Difficulty in reformulation: Easy.

Difficulty in solving the reformulated problem: Medium (the number of variables is in the order of hundreds to ten thousands).

We consider the problem of dispatching ℓ electric power generators in a power system. Here we use the direct current (DC) problem formulation. The decision variables are the power generation of each generator P_i and the bus angle θ_i .

Our goal is to minimize the total generation cost, subject to the balance of supply and demand, the power flow equations, and the power generation limits (Bakirtzis & Biskas, 2003).

- ¹⁵⁴⁷ This problem is convex: (Bakirtzis & Biskas, 2003)
 - Minimize $\max_{i=1,\dots,\ell} C_i(P_i)$ subject to $\left| \frac{1}{x_{ij}} \left(\theta_i - \theta_j \right) \right| \le F_{ij}^{\max}, \quad \forall i, j = 1, \dots, \ell,$ $P_i^{\min} \le P_i \le P_i^{\max}, \quad i = 1, \dots, \ell,$ $\mathbf{B} \cdot \theta = \mathbf{P} - \mathbf{D},$ (14)

where $C_i(\cdot)$ is the generation cost function (usually quadratic), x_{ij} and F_{ij}^{max} are the admittance and the power flow limit of the transmission line connecting bus *i* and bus *j*, P_i^{\min} and P_i^{\max} are the lower and upper limits of power generator *i*, **B** is the admittance matrix, and **D** is the demand vector.

- 1559 F.2 EXAMPLES OF **Type-II** PROBLEMS
- **[P1] PC for maximum throughput with interference** (PC-MaxRate-interference)

1562• Reformulation strategies: Change of variables.

- **Difficulty in reformulation:** Hard (highly skilled techniques were used).
- **Difficulty in solving the reformulated problem:** Medium (the reformulated convex problem may not be recognized by standard solvers as convex).



1571 1572 1573

1585

1586 1587

1596

1598

1574 We consider the problem of determine the transmit power of ℓ pairs of transceivers. They operate in 1575 the same frequency at the same time, hence causing interference to each other. The problem data is 1576 a channel gain matrix $\mathbf{G} \in \mathbb{R}^{\ell \times \ell}$, where g_{ij} is the channel gain from transmitter j to receiver i, and 1577 the noise power vector $\sigma \in \mathbb{R}^{\ell}$ with σ_i as the noise power at receiver i.

Our goal is to choose the transmit power, denoted $x \in \mathbb{R}^{\ell}_+$ with x_i being the power of transmitter *i*, in order to maximize the total throughput.

1581 This problem is nonconvex and is proved to be NP-hard in its original form: (Luo & Zhang, 2008)

Maximize
$$\sum_{i=1}^{\ell} \log_2 \left(1 + \frac{g_{ii}x_i}{\sum_{j \neq i} g_{ij}x_j + \sigma_i} \right)$$

subject to
$$\sum_{i=1}^{\ell} x_i \le 1,$$
$$x_i \ge 0, \quad i = 1, \dots, \ell.$$
 (15)

Under a "weak interference" assumption, we can reformulate the problem as a convex problem. Specifically, we define $\mathbf{H} = \mathbf{G} + \sigma \cdot \mathbf{1}^T$, namely $h_{ij} = g_{ij} + \sigma_i$. We assume that \mathbf{H} is invertible with $\mathbf{H}^{-1} \triangleq \mathbf{I} - \mathbf{C}$ and that \mathbf{C} is a nonnegative matrix. This assumption holds true when $g_{ii} > \sum_{j \neq i} g_{ij}$ for all *i*, namely the interference is weak. Under this assumption, with the change of variables $\mathbf{y} = \mathbf{H}\mathbf{x}$, namely $y_i = \sum_{j=1}^n h_{ij}x_j$, the problem can be reformulated as the following convex problem:

Maximize
$$\sum_{i=1}^{n} \log_2 \left(\frac{y_i}{c_i^T y} \right)$$

subject to
$$\sum_{i=1}^{n} \left(y_i - c_i^T y \right) \le 1,$$
$$y_i - c_i^T y \ge 0, \ i = 1, \dots, n.$$
 (16)

We can prove that the objective function is concave by showing that the negative of its Hessian matrix is positive semidefinite. However, standard solvers such as CVXPY cannot recognize it as a concave function. Therefore, we need to define a custom objective function in standard solvers or write custom interior-point methods to solve it.

[P2] Optimal power flow in alternating current and radial networks (AC-OPF-radial)

Relaxation strategies: Semidefinite relaxation (SDR).

- Difficulty in reformulation: Medium (SDR techniques were used).
- **Difficulty in solving the reformulated problem:** Hard (the number of variables is in the order of hundreds to ten thousands).

Please see AC-OPF in App. F.3. It is shown that under certain conditions (e.g., the power network is a radial network, which is often the case for distribution networks), the convex relaxation of the problem is exact (Lavaei & Low, 2011). In other words, the non-convex original problem is equivalent to a convex problem.

1617 1618 1619

1610

F.3 EXAMPLES OF TYPE-III PROBLEMS

[P1] Beamforming with nonconvex problem definitions (Beamform-Nonconvex)

1620 • Relaxation strategies: Semidefinite relaxation (SDR). 1621

S11

• Difficulty in reformulation: Easy (standard SDR techniques were used).

• **Difficulty in solving the reformulated problem:** Medium (the relaxed convex problem is a SDP, and recovery methods are needed).

1625 We consider the same setting as Beamform-MinSidelobe. But here our goal is to maximize the gain at the target direction θ_{tar} , while limiting the ripple effect at directions $\theta_1, \ldots, \theta_m$ outside the 1626 target area. 1627

1628 This problem is non-convex: (Fuchs, 2013) 1629

1622

1623

1624

1630 1631

1656

1661

1663

1664

1665

Maximize
$$|G(\mathbf{w}; \theta^{\text{tar}})|$$

subject to $1/\delta \le |G(\mathbf{w}; \theta_i)| \le \delta, \ i = 1, \dots, m.$ (17)

It is non-convex because we maximize a convex function (i.e., norm) and have constraints on convex 1633 functions greater than or equal to a constant. 1634

In this case, the standard semidefinite relaxation technique can be used, which "lifts" the problem to 1635 higher dimensions. Specifically, we define a rank-1 semidefinite matrix $\mathbf{W} \triangleq \mathbf{w}\mathbf{w}^{H}$. Then the gain 1636 at direction θ satisfies 1637

1638
1639
1639
1640
1641
1642

$$|G(\mathbf{w};\theta)|^{2} = \operatorname{trace}\left(\mathbf{E}(\theta_{i}) \cdot \mathbf{W}\right),$$
1639
1640
1640
1641
1642

$$\mathbf{e}(\theta_{i}) = \mathbf{e}(\theta_{i})^{T} \in \mathbb{C}^{n \times n} \text{ with}$$
1640
1641
1642

$$\mathbf{e}(\theta_{i}) = \left[e^{\mathbf{i}(x_{1} \cos \theta + y_{1} \sin \theta)}, \dots, e^{\mathbf{i}(x_{n} \cos \theta + y_{n} \sin \theta)}\right]^{T}$$

1643 With the new matrix variable W, we have the following equivalent problem: 1644

1645Maximizetrace
$$(\mathbf{E}(\theta_{tar}) \cdot \mathbf{W})$$
1646subject to $(1/\delta)^2 \leq \text{trace } (\mathbf{E}(\theta_i) \cdot \mathbf{W}) \leq \delta^2, \ i = 1, \dots, m,$ (18)1647 $\mathbf{W} \succeq 0,$ $\operatorname{rank}(\mathbf{W}) = 1.$ (18)

1650 Here, the only nonconvexity comes from the rank constraint. By removing it, we get the following 1651 convex relaxation: 1652

1653Maximizetrace
$$(\mathbf{E}(\theta_{tar}) \cdot \mathbf{W})$$
1654subject to $(1/\delta)^2 \leq \text{trace } (\mathbf{E}(\theta_i) \cdot \mathbf{W}) \leq \delta^2, \ i = 1, \dots, m,$ (19)1655 $\mathbf{W} \succeq 0.$

1657 In general, we need to recover an approximate solution vector w from the solution matrix W. Under 1658 certain conditions (e.g., uniform linear arrays), we can guarantee to recover the exact solution vector.

1659 [P2] Optimal power flow in alternating current (AC-OPF) 1660

• Relaxation strategies: Semidefinite relaxation (SDR).

- 1662 • Difficulty in reformulation: Medium (SDR techniques were used).
 - Difficulty in solving the reformulated problem: Hard (the number of variables is in the order of hundreds to ten thousands).

This is the alternating current (AC) version of DC-OPF, and is more accurate in power system 1666 modelling. 1667

1668 A canonical version of AC-OPF is as follows: (Lavaei & Low, 2011) 1660

where $V \in \mathbb{C}^{\ell}$ is the vector of voltages at each bus, $C_i(\cdot)$ includes the generation cost and the power loss in transmission lines, the first constraints are voltage magnitude constraints to ensure stability of transmission lines, and the second constraints are power injection constraints derived from Ohm's law and Kirchhoff's laws.

1678 1679 This problem is non-convex because of the constraints. A common approach is to perform semidefinite relaxation and solve the convex relaxation. Please see (Lavaei & Low, 2011) for details.

1681 1682 F.4 Examples of Ambiguous Problems

1683 [P1] Power control for maximum throughput (PC-MaxRate) 1684

This prompt does not provide the context of whether these transceivers cause interference to each other. If they are using orthogonal frequency-division multiplexing (OFDM), the major multipleaccess protocol in the fourth-generation (4G) and fifth-generation (5G) communications, we can formulate the problem as the relatively easy **PC-MaxRate-decoupled**. If they are self-organizing as a mesh network, we can formulate the problem as the challenging **PC-MaxRate-interference**. In this case, the LLM should ask for the specific use case to determine the mathematical model.

1690 1691 [P2] Beamforming (Beamform)

1692 The design specifications are not clear. The user may want to minimize the sidelobe while keeping 1693 certain gain at the target direction, resulting in **Beamform-MinSidelobe**, or may want to maximize 1694 the gain at the target direction while limiting the sidelobe, resulting in **Beamform-Nonconvex**. In 1695 this case, the LLM should ask for a clear design specification form the user.

[696 [P3] Optimal power flow (OPF)

The problem statement is not clear about which version of the OPF problem to use. When the power system is not heavily loaded and when fast solutions are needed, we can use the DC approximation and formulate the problem as DC-OPF, which is a LP. If we need accurate solutions (e.g., solutions that inform us of the power loss on the transmission lines), we need to formulate it as AC-OPF, which is non-convex in general. However, if the underlying network is a distribution grid, which is a radial network, the resulting AC-OPF-radial has a convex relaxation, which has been proven to be exact (i.e., equivalent to the original non-convex problem).

1705 F.5 OVERVIEW

1 - 0	I D'	∩ .	· · · 1 · · ·	41	1100 14	• • •	1	1	C	· · · ·	41		1	
1/0/	In $F10$	y we v	isnanze	ine	annem v	1n	SOLVING	ana	reformina	ino	Ine	exam	nie :	proplems
1101	111 1 15.	/,	ISuulle	une	unneurry	111	SOLVING	unu	rerormula	ms	une	onum	pic	problems.

1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

1704



1782 SOLVER PERFORMANCE COMPARISON G 1783

1784 In this section, we present simulation results comparing the performance of various optimization 1785 solvers on the PC_MinPower problem, both in its original nonconvex form and in a reformulated 1786 convex form. We consider problem instances with $\ell = 10$ and $\ell = 100$ users (i.e., ℓ optimization 1787 variables). For each instance, we evaluate the solvers in terms of success rate, optimality gap, and 1788 average solve time over 100 random samples. 1789 1790 G.1 EXPERIMENTAL SETUP 1791 The PC_MinPower problem aims to minimize the total power consumption in a system while 1792 satisfying certain constraints. The original formulation of this problem is nonconvex, which can 1793 pose challenges for optimization algorithms. However, by reformulating the problem, it can be 1794 converted into an equivalent convex problem, which is generally easier to solve efficiently. 1795 1796 We evaluated the following solvers: 1797 • General-Purpose Solvers: 1799 - TRCA: Trust-Region Constrained Algorithm - **SLSOP**: Sequential Least Squares Programming. - COBYLA: Constrained Optimization BY Linear Approximations. 1801 - **COBYQA**: Constrained Optimization BY Quadratic Approximations. 1803 Convex Program Solvers: - CLARABEL: A conic optimization solver. 1805 - ECOS: Embedded Conic Solver. 1806 SCS: Splitting Conic Solver. 1807 - OSQP: Operator Splitting Quadratic Program Solver. 1808 1809 For each solver and problem instance, we recorded: 1810 • Success Rate: The percentage of runs where the solver successfully found a feasible solution. 1811 1812 • **Optimality Gap**: The difference between the objective value obtained by the solver and the known 1813 optimal value. • Average Solve Time: The average computation time (in seconds) required by the solver. 1814 1815 1816 G.2 RESULTS 1817 Tables 7 and 8 present the performance of the solvers for problem sizes n = 10 and n = 100, 1818 respectively. 1819 1820 Table 7: Solver Performance for $\ell = 10$ Users 1821 1822 Original Nonconvex Problem Reformulated Convex Problem Solver Type Success Success Optimality Gap Time (s) **Optimality Gap** Time (s) 1824 General-Purpose Solvers $7.31 imes 10^{-3}$ 1.25×10^{-3} TRCA General-Purpose 100% 0.0399 100% 0.0420 1825 6.47×10^{-7} 0.0019 6.48×10^{-7} SLSQP 100% 100% 0.0009 General-Purpose 1826 2.94×10^{-6} 0.0073 7.15×10^{-6} COBYLA General-Purpose 67% 100% 0.0039 1827

1	831
1	832
1	833
1	834

1835

1828

1829

1830

G.3 DISCUSSION

COBYQA

ECOS

OSQP

SCS

CLARABEL

Convex Program Solvers

The results demonstrate several key observations:

General-Purpose

Convex Solver

Convex Solver

Convex Solver

Convex Solver

0%

6%

100%

100%

100%

100%

9.80

 6.32×10^{-7}

 6.16×10^{-7}

 4.45×10^{-7}

 $6.48 imes 10^{-7}$

14.4067

0.0002

0.0001

0.0002

0.0003

	Туре	Original Nonconvex Problem			Reformulated Convex Problem		
Solver		Success	Optimality Gap	Time (s)	Success	Optimality Gap	Time (s)
General-Purpose Solvers							
TRCA	General-Purpose	100%	$7.75 imes10^{-2}$	0.6628	100%	$1.28 imes 10^{-2}$	0.6856
SLSQP	General-Purpose	100%	1.04×10^{-6}	0.0750	100%	1.04×10^{-6}	0.0298
COBYLA	General-Purpose	0%	_	6.0764	100%	2.86×10^{-5}	9.9629
COBYQA	General-Purpose	0%		_	0%	—	_
Convex Program Solvers							
CLARABEL	Convex Solver	_	_	_	100%	6.22×10^{-7}	0.0121
ECOS	Convex Solver	_		_	100%	$9.91 imes 10^{-7}$	0.0097
SCS	Convex Solver	_	_	_	100%	$1.05 imes 10^{-6}$	0.0055
OSQP	Convex Solver		—	—	100%	1.04×10^{-6}	0.0110

Table 8: Solver Performance for $\ell = 100$ Users

1848

1836

- Importance of Problem Reformulation: For general-purpose solvers, reformulating the original nonconvex problem into an equivalent convex problem significantly improves solution quality. This improvement is more pronounced for larger problem sizes (l = 100). For instance, COBYLA's success rate increased from 0% to 100% when the problem was reformulated.
- Solver Selection Matters: Different general-purpose solvers exhibit varying performance levels.
 SLSQP consistently achieves near-zero optimality gaps and high success rates with relatively low solve times across both problem formulations and sizes. In contrast, COBYQA fails to find feasible solutions in most cases, highlighting the necessity of careful solver selection.
- Performance of Convex Program Solvers: For the reformulated convex problem, convex program solvers (CLARABEL, ECOS, SCS, OSQP) show excellent and consistent performance. They all achieve 100% success rates, negligible optimality gaps, and minimal solve times. The differences among these solvers are minimal, suggesting that any of them would be suitable for solving the convex formulation efficiently.

These findings underscore the importance of problem reformulation and solver selection in optimization tasks. Reformulating a nonconvex problem into a convex one can significantly enhance the performance of general-purpose solvers. Additionally, selecting the appropriate solver is crucial, as it can greatly impact the success rate and computational efficiency.

1866 1867 1868

1869 1870

H BROADER PERSPECTIVE ON AUTOFORMULATION AND OPTIMIZATION MODELING

Our autoformulation framework addresses a specific, yet crucial challenge: translating problem
 descriptions into mathematical and computational models. However, it is important to acknowledge
 that this represents just one component of the broader modeling process, which typically involves:

- 18741. Information gathering from diverse stakeholders, often requiring multiple iterations and integration of implicit domain knowledge/conventions (i.e., 'tribal knowledge');
- 1876
 1877
 2. Handling complex problem characteristics such as stochasticity, time-varying dynamics, and large-scale variables;
- 1878
 3. Rigorous validation against real-world data, including sensitivity analysis of modeling assumptions and verification of edge cases;
- 4. Continuous communication between technical and business stakeholders to ensure practical utility.

Understanding these complexities helps identify which aspects of modeling can be effectively auto mated, thereby enabling OR practitioners to focus their expertise on more nuanced challenges such as stakeholder engagement and validation of modeling assumptions.

1886

1882

1887

1888