

# LARC: Towards Human-level Constrained Retrosynthesis Planning through an Agentic Framework

Frazier N. Baker<sup>1</sup>, Daniel Adu-Ampratwum<sup>2</sup>, Reza Averly<sup>1</sup>, Botao Yu<sup>1</sup>, Huan Sun<sup>1</sup>, and Xia Ning<sup>1,2,3,4</sup>

<sup>1</sup> Department of Computer Science and Engineering

<sup>2</sup> Division of Medicinal Chemistry & Pharmacognosy

<sup>3</sup> Department of Biomedical Informatics

<sup>4</sup> Translational Data Analytics Institute

The Ohio State University, Columbus, OH, USA

{baker.3239,adu-ampratwum.1,averly.1,yu.3737,sun.397,ning.104}@osu.edu

## Abstract

Large language model (LLM) agent evaluators leverage specialized tools to ground the rational decision-making of LLMs, making them well-suited to aid in scientific discoveries, such as constrained retrosynthesis planning. Constrained retrosynthesis planning is an essential, yet challenging, process within chemistry for identifying synthetic routes from commercially available starting materials to desired target molecules, subject to practical constraints. Here, we present LARC, the first LLM-based Agentic framework for Retrosynthesis planning under Constraints. LARC incorporates agentic constraint evaluation, through an Agent-as-a-Judge, directly into the retrosynthesis planning process, using agentic feedback grounded in tool-based reasoning to guide and constrain route generation. We rigorously evaluate LARC on a carefully curated set of 48 constrained retrosynthesis planning tasks across 3 constraint types: avoiding carcinogens, pyrophoric substances, or a user-specified substance. LARC achieves a 72.9% success rate on these tasks, vastly outperforming LLM baselines and approaching human expert-level success in substantially less time. The LARC framework is extensible, and serves as a first step towards an effective agentic tool or a co-scientist to human experts for constrained retrosynthesis.

## 1 Introduction

Large language model (LLM) agents have recently shown great promise as evaluators [50]. An LLM agent evaluator, also called an Agent-as-a-Judge, provides grounded evaluations in complex settings, eliciting the strengths of general-purpose LLMs and domain-specific tools. These evaluators are poised to make substantial impacts in scientific discoveries, such as in chemistry, where complex evaluation settings abound [8, 4], accurate evaluations are paramount, and many domain-specific tools exist to ground evaluations [6, 2]. One such application is *constrained retrosynthesis planning* [12], an essential process in chemistry for identifying synthetic routes from commercially available starting materials to desired target molecules (products), subject to practical constraints (e.g., avoiding hazardous molecules).

LLM agents for constrained retrosynthesis planning are not explored. Current artificial intelligence (AI) methods are primarily focused on unconstrained retrosynthesis planning [10, 49, 27, 17, 11, 3, 14], aiming to generate synthetic routes that are feasible. Only a few AI methods have attempted to address constrained retrosynthesis planning [47, 20]. However, they only support a very simplistic constraint – including a user-specified molecule in the synthetic routes, and cannot be applied to more general and practical constraints, such as avoiding broad classes of hazardous molecules in the synthetic routes. These constraints are substantially more challenging to evaluate and enforce, requiring specialized knowledge of hazardous materials. Rather than guiding synthesis planning towards a single, clearly-specified goal, these constraints require guiding it away from many diverse hazards. LLM agents are well-suited for such a challenge, as they can leverage specialized chemistry tools to ground their evaluations and make rational decisions to guide the planning. Furthermore, LLM agents could support a variety of constraints, choosing the appropriate tools for each type of constraint. They have the potential to mimic the typical behaviors of human chemists, such as using reference materials [41, 24] to assess safety constraints during retrosynthesis planning [8], and thus, automate, accelerate, and optimize the reliability of outcomes in the constrained retrosynthesis planning process.

Here, we present LARC, an **LLM-based Agentic framework for Retrosynthesis planning under Constraints**. Figure 1 presents an overview of LARC. LARC uses an Agent-as-a-Judge, equipped with chemistry tools, to evaluate constraints during retrosynthesis planning. This agentic feedback is incorporated back into the retrosynthesis planning process, dynamically guiding and constraining route generation. It addresses key safety constraints in retrosynthesis planning, such as avoiding carcinogens, pyrophoric substances, or a user-specified substance. LARC is extensible by design, allowing it to improve or expand as future capabilities emerge. To the best of our knowledge, LARC is the first agentic framework for constrained retrosynthesis planning, representing an innovative paradigm for this complex scientific problem.

We rigorously evaluate LARC on a carefully curated set of 48 constrained retrosynthesis planning tasks across 3 constraint types. LARC achieves an impressive 72.9% success rate on these tasks, indicating that LARC is very effective at constrained retrosynthesis planning. We compare LARC against general-purpose LLMs and a human expert. The experiments show that LARC vastly outperforms the LLMs and approaches expert-level success in substantially less time. Case studies indicate that LARC can mimic human expert’s retrosynthesis planning logic and even produce better synthetic routes on some tasks. An ablation study reveals the key impact of agentic tooling in LARC, enabling high success rate and efficiency through deliberate and grounded evaluations. With further extension to cover comprehensive practical constraints, the LARC framework can serve as an effective agentic tool or a co-scientist to human experts for constrained retrosynthesis. The code and data for LARC are publicly available at <https://github.com/ninglab/LARC>.

## 2 Related work.

**Constrained retrosynthesis planning.** Recently, AI methods have emerged for constrained retrosynthesis planning. TANGO\* [20] and DESP [47] both constrain retrosynthesis planning to include a user-specified molecule in the synthetic routes. TANGO\* adapts an unconstrained retrosynthesis planner to constrained retrosynthesis planning, incorporating molecule similarity to the user-specified molecule as constraint guidance. DESP performs a double-ended search, expanding the synthetic route from the target molecule and the user-specified molecule until the route is connected and complete. While TANGO\* and DESP show that AI can perform constrained retrosynthesis planning, they do not support more practical constraints (e.g., avoiding hazardous molecules). Recently, Bran et al. [7] introduced LLM-based re-ranking of synthetic routes generated by an unconstrained planner to identify those satisfying the constraints. This

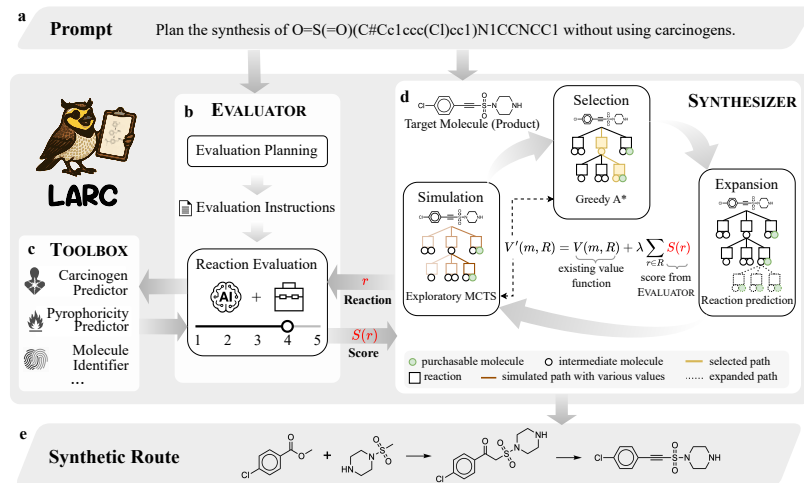


Figure 1: Overview of LARC. **a**, The user prompt specifies the target molecule (product) and the constraint for synthesis. **b**, EVALUATOR acts as a judge and evaluates each individual reaction involved in the retrosynthesis planning with respect to the constraint. **c**, The toolbox contains external tools to ground EVALUATOR’s decision-making **d**, SYNTHESIZER explores and constructs synthetic routes, incorporating feedback from EVALUATOR. **e**, LARC outputs a synthetic route that satisfies the constraint.

approach is computationally expensive and assumes that some generated synthetic routes will satisfy the constraint, which may not always occur. Moreover, it relies on an LLM’s intrinsic knowledge alone to evaluate routes, which may be insufficient for some constraints. Thus, there remains a need for a framework that directly incorporates agentic constraint evaluation into the planning process, enabling synthetic route generation under practical constraints. Additional related work on LLMs for chemistry is discussed in Appendix C.

### 3 LARC: An Agentic Framework

LARC is an agentic framework for constrained retrosynthesis planning. For a given target molecule (i.e., product), it effectively plans its synthetic routes that satisfy constraints specified by user prompts. LARC incorporates constraint evaluation directly into its planning process, leveraging agentic feedback grounded in tool-based reasoning to dynamically guide and constrain route generation. LARC features two key components: **(1)** EVALUATOR, which acts as a judge and evaluates each individual reaction involved in the retrosynthesis planning with respect to the constraint, and **(2)** SYNTHESIZER, which explores and constructs synthetic routes, incorporating feedback from EVALUATOR. Through these components, LARC integrates and elicits the strengths of LLMs, cheminformatics tools, and search and planning algorithms, representing an innovative, extensible, agentic paradigm, with Agent-as-a-Judge in the loop, for computer-aided constrained retrosynthesis planning.

#### 3.1 Evaluator

EVALUATOR, an LLM-based Agent-as-a-Judge [50], assesses whether reactions satisfy the constraints specified in the user prompts. EVALUATOR approaches this task in two phases: **(1)** evaluation planning, where it generates evaluation instructions specific to user constraints, and **(2)** reaction evaluation, where, following the evaluation instructions, it evaluates reactions and provides feedback to SYNTHESIZER.

**Evaluation planning.** Evaluation planning occurs once for each user prompt at the beginning, during which EVALUATOR details its plan for reaction evaluation, including tools it plans to use

and its scoring strategy, into concrete and structured evaluation instructions. These instructions will be consistently followed by EVALUATOR throughout the entire retrosynthesis planning process. The system instructions for evaluation planning are provided in Appendix B. They were designed to encourage strategic, step-by-step evaluation while providing key information on tool syntax and semantics, including input/output formats.

**Reaction evaluation.** Following the evaluation instructions, EVALUATOR evaluates reactions, assessing how well they satisfy the constraint. In doing so, EVALUATOR adaptively switches between two action modes: (1) leveraging its own intrinsic knowledge, acting as an AI expert itself, or (2) consulting external tools in the toolbox, inspired by the typical behaviors of synthetic chemists using reference materials [41, 16] during retrosynthesis planning. EVALUATOR calibrates the evaluation outputs into scores, which will guide SYNTHESIZER. These scores, denoted as  $S(r)$ , where  $r$  is the evaluated reaction, quantify the degree of constraint satisfaction on a discrete scale from 1 to 5, with 1 indicating complete violation and 5 indicating full satisfaction. EVALUATOR is extensible by design to incorporate new tools or adapt to tool updates, allowing it to improve and expand as future capabilities emerge.

### 3.2 Synthesizer

LARC uses SYNTHESIZER to generate synthetic routes for target molecules. SYNTHESIZER is built upon existing unconstrained retrosynthesis planners, adapting them to constrained retrosynthesis planning. Unconstrained planners typically leverage a search algorithm, such as A\* [21] or Monte Carlo Tree Search (MCTS) [13], to iterate backwards from the target molecule and search for intermediates. The search is expanded through single-step retrosynthesis planning on the current initial intermediates until a full synthetic route is found, starting from commercially available materials. The search expansion is guided by some value function,  $V(m, R)$ , which estimates the utility of expanding the search along route  $R$  upstream from its initial intermediate  $m$ .  $V(m, R)$  is independently pre-trained with respect to general objectives in unconstrained retrosynthesis planning, such as preference for short routes with feasible and chemically plausible reactions.

To constrain the search, SYNTHESIZER combines  $V(m, R)$  with EVALUATOR’s score  $S(r)$ , thus producing a new, *constraint-aware* value function  $V'(m, R)$  to guide the search:

$$V'(m, R) = V(m, R) + \lambda \underbrace{\sum_{r \in R} S(r)}_{\text{constraint evaluation by EVALUATOR over all reactions along the route } R}, \quad (1)$$

where  $m$  is the intermediate molecule,  $R$  is a synthetic route starting from  $m$ ,  $r$  is a reaction in  $R$ ,  $S(r)$  is the score from EVALUATOR, and  $\lambda > 0$  is a trade-off hyperparameter. Thus,  $V'$  incorporates the constraint evaluation by EVALUATOR on all the reactions along  $R$  to guide retrosynthesis planning, ensuring the entire route is maximally subject to the constraint. This represents an innovation, differentiating SYNTHESIZER from existing work on constrained retrosynthesis [47, 20]. Note that SYNTHESIZER can easily adapt any unconstrained retrosynthesis planner without retraining the original  $V(m, R)$ , allowing LARC to accommodate future advancements in retrosynthesis planning.

## 4 LARC Instantiation

The current implementation of LARC adapts MEEA\* [49], a state-of-the-art method for unconstrained retrosynthesis planning, as SYNTHESIZER. MEEA\* uses a two-step process to determine how to expand its search. The first step uses MCTS to simulate route planning over the expanded partial routes, selecting at most  $K$  candidate routes for further expansion. In this simulation, LARC uses  $V'_{\text{MCTS}}$  as its constraint-aware value function, that is,  $V'_{\text{MCTS}} =$

$V_{\text{MCTS}} + \lambda \sum_{r \in R} S(r)$  (Equation 1), where  $V_{\text{MCTS}}$  includes an upper confidence bound (UCB) term [28] to encourage exploration. In this case,  $S(r)$  is calculated as follows: For the reactions that have not been evaluated by EVALUATOR in previous expansions, an optimistic default score ( $S(r) = 5$ ) is used to further encourage exploration; for those that have been evaluated, their actual evaluation score  $S(r)$  is used. In the second step of MEEA\*, A\* search is used to select a single route for expansion from the  $K$  candidate routes. In this step, EVALUATOR first evaluates all the reactions in the  $K$  candidate routes. Then, LARC uses  $V'_{A^*}$ , that is,  $V'_{A^*} = V_{A^*} + \lambda \sum_{r \in R} S(r)$  (Equation 1), to select a single route for expansion. Implementation details are presented in Appendix D.

**Benchmark dataset for constrained retrosynthesis.** We carefully curated a benchmark set of constrained retrosynthesis tasks for a set of target molecules (products) from the USPTO-190 [10], each with a single constraint to satisfy. In this instantiation, we considered constraints of avoiding hazardous substances in the synthetic routes, which can pose serious safety risks to chemists, equipment, and the environment. Three types of hazardous substances are included: (1) carcinogens, which are capable of causing cancer based on the classification by the International Agency for Research on Cancer (IARC) [24]; (2) pyrophoric substances, which can ignite spontaneously upon exposure to air, according to the U.S. Navy report on air- and water-reactive materials [18]. and (3) a user-specified hazardous substance (e.g., phosgene). To ensure these tasks are non-trivial, we selected the tasks such that the state-of-the-art unconstrained retrosynthesis planning methods could generate a valid route but violate the constraint, and the target molecule was not used in  $V(m, R)$  pre-training. In the end, 48 tasks were constructed for the benchmark set. Figure 2b presents the distribution of the three types of constraints. Please note, these 48 tasks represent diverse constraint violations from the unconstrained planner, with 28 tasks covering 16 distinct carcinogens, 12 tasks covering 9 pyrophoric substances, and 8 tasks covering 8 user-specified molecules. The benchmark tasks are available in Appendix A.

**Tools for retrosynthesis constraints.** Three specific tools are supplied in the toolbox: (1) a carcinogen predictor, which predicts whether a given molecule is a carcinogen using the state-of-the-art ADMET-AI [39] model, (2) a pyrophoricity predictor, which predicts the pyrophoricity of molecules by comparing them with known pyrophoric substances [18], with higher molecule similarities indicating higher likelihood of pyrophoricity, and (3) a molecule identifier, which identifies specific hazardous molecules using their fingerprints [33, 36]. Outputs from the tools will be calibrated by EVALUATOR into scores ( $S(r)$  in Equation 1).

## 5 Experimental Settings

**Base LLMs and baselines.** We select MISTRAL NEMO [32] and CLAUDE 3.5 SONNET [1] as the base LLMs for EVALUATOR, resulting in LARC variations denoted as LARC<sub>Mistral</sub> and LARC<sub>Claude</sub>, respectively. MISTRAL NEMO is selected as a representative small, open-source LLM for its strong instruction-following capabilities. It allows for the evaluation of how LARC performs with a cost-effective, openly available model, thereby demonstrating its practicality and accessibility. CLAUDE 3.5 SONNET represents a state-of-the-art, closed-source LLM that has demonstrated strong performance in chemistry tasks and tool use [23, 46]. Note that LARC is not limited to these two base LLMs, as its design is model-agnostic and can incorporate other LLMs as the field advances.

We use general-purpose LLMs as the baselines for constrained retrosynthesis planning, including CLAUDE 3.5 SONNET [1], GPT-4O [34], DEEPSEEK R1 [15], and MISTRAL NEMO [32]. We also compare LARC against a human expert in retrosynthesis planning, denoted as EXPERT. EXPERT is an experienced synthetic chemist with a doctoral degree and over 17 years of experience in retrosynthesis planning. We detail our rigorous human retrosynthesis planning protocol

below and our LLM baselines in Appendix E. Please note, existing constrained retrosynthesis methods cannot be fairly compared as baselines due to the novel and challenging nature of this problem. DESP and TANGO\* do not support the constraints studied here, and the re-ranking approach [7] relies on multiple route generations.

**Constrained retrosynthesis planning by human experts.** We also compare LARC against a human expert in retrosynthesis planning, denoted as EXPERT. EXPERT is an experienced synthetic chemist with a doctoral degree and over 17 years of experience in retrosynthesis planning. Such human experts are a scarce resource, limiting our experiments to a single EXPERT. For each task, EXPERT was provided with clear instructions on the constraint and the desired chemical product. EXPERT had full access to necessary resources and reference materials, including lists of purchasable materials [10], carcinogenic chemicals, and pyrophoric chemicals, as well as external and online references such as Reaxys [16] and SciFinder [9]. However, EXPERT was not permitted to use computer-aided multi-step retrosynthesis planning tools, such as AIZynthFinder [17] or MEEA\* [49]. To ensure both the quality and efficiency of human planning, EXPERT was instructed to complete each task with careful attention to detail while minimizing planning time. To support sustained performance and maintain high-quality planning, EXPERT was encouraged to take a 15-minute break between tasks.

**Evaluation metrics.** The generated synthetic routes are evaluated according to the following criteria: **(1)** Route presence: the routes are not empty – they contain some molecules; **(2)** Route connectivity: the routes are fully connected – all intermediate molecules are synthesized from preceding precursors; **(3)** Target reachability: the routes lead to the target molecule; **(4)** Commercial availability: the starting materials of the routes are directly purchasable from eMolecules dataset [10]; **(5)** Molecule validity: the molecules involved in the routes are chemically correct; and **(6)** Constraint satisfaction: the routes meet the specific constraints (avoid certain substances). Most of these criteria can be evaluated using automatic, in silico methods; however, molecule validity may require manual verification to account for valid non-SMILES representations from LLM baselines are considered valid. Based on these criteria, the following metrics are used to evaluate synthetic routes, as illustrated in Figure 2a: **(1) Success rate:** the percentage (%) of routes that satisfy all the six criteria, that is, successful routes; **(2) Validity rate:** the percentage (%) of routes that satisfy route presence, route connectivity, target reachability, commercial availability, and molecule validity, that is, valid routes per se but they do not necessarily satisfy the constraint; **(3) Presence rate:** the percentage (%) of routes that satisfy route presence, that is, non-empty routes.

## 6 Results

We focus primarily on carcinogenicity-constrained tasks here, with results of the pyrophoricity-constrained and user-specified-constrained tasks in Appendix F.

### 6.1 Carcinogenicity-constrained retrosynthesis planning

Figure 2d shows that LARC<sub>Mistral</sub> achieves a 82.1% success rate and LARC<sub>Claude</sub> achieves 64.3% in planning synthetic routes that avoid carcinogenic substances, both vastly outperforming the best LLM baseline CLAUDE 3.5 SONNET (success rate only 25.0%). CLAUDE 3.5 SONNET, as a general-purpose LLM, lacks specialized training on retrosynthesis planning [1, 46], and fails in the vast majority of the cases. Unlike general-purpose LLMs, LARC<sub>Mistral</sub> and LARC<sub>Claude</sub> are agentic and intentionally designed for this constrained retrosynthesis, eliciting the strengths of LLM reasoning, domain-specific tools, and efficient search algorithms. The results demonstrate that LARC is very effective in generating synthetic routes avoiding carcinogenic substances.

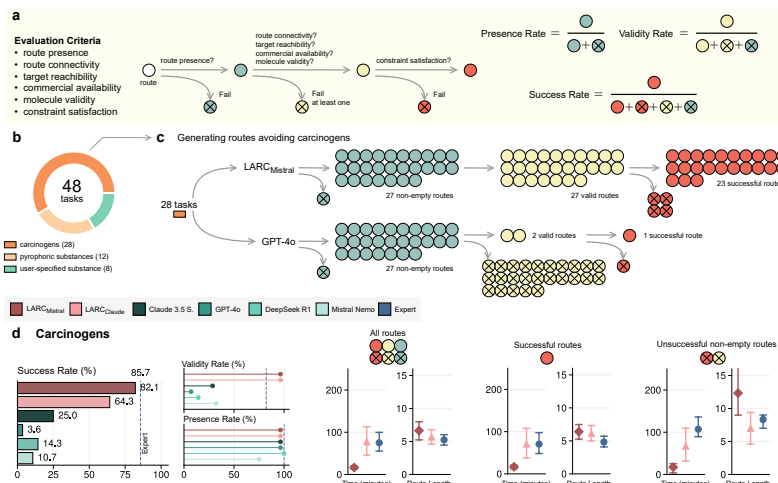


Figure 2: Evaluation. **a**, Evaluation criteria defining success, validity, and presence rates; **b**, Benchmark dataset of 48 curated constrained retrosynthesis tasks spanning 3 constraint types; **c**, Example evaluation showing calculation of presence, validity, and success rates; **d**, success, validity, and presence rates, planning time, and route length for benchmark tasks for avoiding carcinogens.

EXPERT generates 28 routes with a success rate of 85.7%, only slightly outperforming LARC<sub>Mistral</sub>. The successful routes by EXPERT have an average length of 4.83, which is shorter than that of LARC<sub>Mistral</sub> (6.39). This indicates that EXPERT may rely on domain-specific heuristics, intuitive shortcut strategies, or tacit knowledge not yet fully captured by LARC<sub>Mistral</sub>. However, EXPERT requires on average 70.42 minutes to generate each successful route, substantially slower than LARC<sub>Mistral</sub> (16.45 minutes) ( $p=2.8e-4$ , two-sided two-sample t-test). This highlights the potential for LARC<sub>Mistral</sub> to accelerate constrained retrosynthesis planning while maintaining near-human-level quality.

LARC and EXPERT exhibit different behaviors in their unsuccessful constrained retrosynthesis planning. LARC is diligent about planning valid routes; all of its non-empty routes are valid. However, LARC’s routes can sometimes violate the constraint. This is often the result of errors in the tools outputs and their interpretation, which is discussed further in Appendix F.4. In contrast, EXPERT is conscientious about constraint satisfaction, but this seems to distract from route validity. Specifically, EXPERT’s synthetic routes may not reach the target molecule, or it may require starting materials that are not commercially available. This highlights a key challenge of constrained retrosynthesis planning for both LARC and EXPERT: balancing constraint satisfaction with validity criteria.

Interestingly, LARC<sub>Mistral</sub> outperforms LARC<sub>Claude</sub> when planning routes that avoid carcinogenic substances. LARC<sub>Mistral</sub> uses the small, open-source MISTRAL NEMO with only 12 billion parameters as its base model [32], whereas LARC<sub>Claude</sub> uses the much larger, proprietary CLAUDE 3.5 SONNET with over 175 billion parameters [1]. LARC<sub>Claude</sub> tends to be slower – 77.56 minutes on average to generate each route, than LARC<sub>Mistral</sub> (16.45 minutes), and evaluate more reactions per route (59.54) on average than LARC<sub>Mistral</sub> (30.8) during the planning process. While it could be a general expectation that larger models perform better due to scaling laws [26], it is observed that LARC enables the use of smaller and thus cheaper models, such as MISTRAL NEMO, without sacrificing performance by grounding LLM reasoning with specialized chemistry tools. Therefore, LARC provides a fast, accurate, and inexpensive solution to constrained retrosynthesis planning with low costs and a high success rate.

Aside from CLAUDE 3.5 SONNET, the other LLM baselines also perform poorly, with the

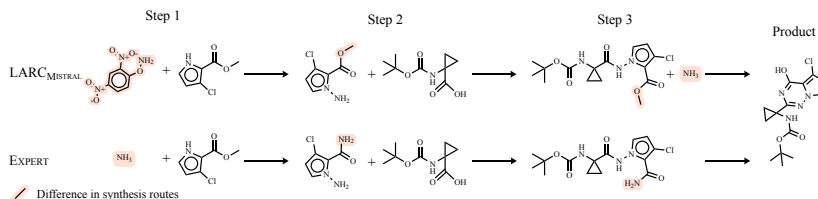


Figure 3: Synthesis route comparison of LARC<sub>Mistral</sub> against EXPERT for tert-butyl (1-(5-chloro-4-hydroxypyrrolo[2,1-f][1,2,4]-triazin-2-yl)cyclopropyl)carbamate.

best success rate of only 14.3%. For example, GPT-4O manages to generate only 1 successful route (success rate 3.6%) among only 2 valid routes. This further highlights the necessity of specialized tools for constrained retrosynthesis planning – a specific, very challenging yet important task, and LARC fills this gap.

## 6.2 Case Studies

Figure 3 shows two plans for the synthesis of the product tert-butyl(1-(5-chloro-4-hydroxypyrrolo[2,1-f][1,2,4]-triazin-2-yl)cyclopropyl)carbamate: one from LARC<sub>Mistral</sub> and one from EXPERT. Both LARC<sub>Mistral</sub> and EXPERT proposed very similar three-step routes to the final product. In both routes, Step 1 involves the amination of commercially available methyl 3-chloro-1H-pyrrole-2-carboxylate, followed by an amidation reaction in Step 2. Finally, Step 3 shows a condensation/cyclization sequence to provide the product.

The main difference between the two proposed routes lies in the source of the aminating agent in Step 1. LARC<sub>Mistral</sub> proposes O-(2,4-Dinitrophenyl)hydroxylamine, to prepare the hydrazine (methyl 1-amino-3-chloro-1H-pyrrole-2-carboxylate), which is known for being a very mild and selective aminating agent for introducing an NH<sub>2</sub> group into a molecule, making it a reagent of choice, especially when a metal-free process is desired. In contrast, EXPERT opted for ammonia (NH<sub>3</sub>) in Step 1. NH<sub>3</sub> is often a preferred reagent for large-scale amination reactions due to its availability and low cost [42]. However, using NH<sub>3</sub> can present some challenges in activation and selectivity, often necessitating sophisticated metal catalyst systems.

As shown in Figure 3, LARC<sub>Mistral</sub>’s use of the mild aminating agent in Step 1 will selectively introduce the NH<sub>2</sub> group on the pyrazole nitrogen without affecting the ester. Conversely, the use of non-selective NH<sub>3</sub> by EXPERT will convert the ester functional group to an amide. Both the ester and the amide can be respectively converted to the product in Step 3, completing the routes. Overall, both proposed routes have similar intermediates and transformations. Note that LARC<sub>Mistral</sub>’s use of O-(2,4-Dinitrophenyl)hydroxylamine is not an anomaly— 58 patented reactions in the USPTO reaction dataset [30] also use this reagent. Overall, LARC<sub>Mistral</sub> can mimic human retrosynthesis planning logic, even in this challenging constrained retrosynthesis setting. An additional case study, showing LARC can even generate better routes than human experts, is presented in Appendix F.3. Furthermore, an ablation study examining the impact of tooling on LARC is presented in Appendix F.4.

## 7 Discussions and Conclusions

LARC addresses only a very simplified version of practical constrained retrosynthesis planning, with the primary goal of demonstrating the potential of agentic AI in solving such complex scientific problems. Our experiments clearly show that, when equipped with appropriate tools for verifying constraint satisfaction, agentic AI can approach human expert-level performance while being more autonomous, effective, and scalable. Such traits are highly attractive in scientific workflows, as they reduce reliance on time-consuming and potentially inconsistent, and

error-prone manual efforts, and allow for the exploration beyond existing, potentially outdated knowledge of human experts.

Meanwhile, rigorous, systematic, and scalable evaluation and validation of results from agentic AI, including LARC, still fall far short. For retrosynthesis, there is not always a definitive and consensus “ground truth” for reaction feasibility. Compounded with the inherent biases of the search space in unconstrained retrosynthesis planners used in LARC, it is still likely that LARC generates synthetic routes that appear valid but are actually infeasible. Even worse, these routes may be difficult to detect and filter using either tools or human expertise. A more reliable option is to incorporate physics-based models, such as molecular dynamics simulations, to assess reaction feasibility in agentic AI models. Unfortunately, this approach requires a lot of customization (e.g., specific force fields) and is not readily scalable, undermining the advantages of agentic AI in being both autonomous and scalable. Ultimately, testing of the AI-generated synthetic routes in a laboratory will be needed to validate the results from agentic AI and truly translate its advantages into real impacts. This will require a selection of the most promising routes, which will eventually still rely on computational tools or human expertise, suffering from the same issues as *in silico* evaluation. Meanwhile, large-scale *in vitro* validation of *in silico*-generated reactions is still challenging. Thus, though very promising, this research on agentic AI for constrained retrosynthesis calls for in-depth and systemic investigation on its autonomous and scalable evaluation and validation strategies. Additional discussion is available in Appendix G.

**Conclusions.** LARC is the first agentic framework for constrained retrosynthesis planning, which uses an LLM-based Agent-as-a-Judge to evaluate reactions and guide the constrained planning. LARC achieves a 72.9% success rate over a carefully curated benchmark of 48 constrained retrosynthesis planning tasks spanning 3 constraint types. It vastly outperforms general-purpose LLM baselines and approaches human expert-level performance. LARC can mimic the retrosynthesis planning logic of human experts and can even produce synthetic routes better than human experts. These results establish LARC as a compelling proof-of-concept for leveraging agentic AI to advance next-generation scientific discovery in synthetic chemistry. More broadly, it illustrates the transformative potential of agentic AI to accelerate progress across the sciences.

## Acknowledgments

This project was made possible, in part, by support from the National Science Foundation grant nos. IIS-2133650 and IIS-2435819, and the National Library of Medicine grant no. 1R01LM014385. Any opinions, findings, conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the funding agency. We would like to thank our colleague, Xinyi Ling, for assisting with some of the figures.

## Ethics Statement

LARC is an agentic framework for constrained retrosynthesis planning, generating synthetic routes for target molecules under practical constraints. We acknowledge that not all target molecules or synthetic routes are safe, and LARC could generate harmful synthetic routes. Thus, we strongly recommend conscious oversight and intervention from human experts when using LARC. Synthetic chemists should confirm the safety of LARC’s generated synthetic routes before using them in laboratory experiments. We recommend using both automated and manual safety checks against external sources (e.g. local safety policies, GHS [41], etc.), and following standard procedures for laboratory safety. We advise all users of LARC to exercise their professional discretion and follow all applicable safety guidelines, laws, regulations, and ethical standards.

## References

- [1] Anthropic. Introducing Claude 3.5 Sonnet, 2024.
- [2] Reza Averly, Frazier N. Baker, and Xia Ning. LIDDIA: Language-based Intelligent Drug Discovery Agent. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2025.
- [3] Frazier N. Baker, Ziqi Chen, Daniel Adu-Ampratwum, and Xia Ning. RLSynC: Offline–Online Reinforcement Learning for Synthon Completion. *Journal of Chemical Information and Modeling*, 64(17):6723–6735, September 2024. Publisher: American Chemical Society.
- [4] Benjamin E. Blass. *Basic principles of drug discovery and development*. Academic Press, London, second edition edition, 2021.
- [5] Daniil A. Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, December 2023. Publisher: Nature Publishing Group.
- [6] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D. White, and Philippe Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, 6(5):525–535, May 2024. Publisher: Nature Publishing Group.
- [7] Andres M. Bran, Theo A. Neukomm, Daniel P. Armstrong, Zlatko Jončev, and Philippe Schwaller. Chemical reasoning in LLMs unlocks steerable synthesis planning and reaction mechanism elucidation, March 2025. arXiv:2503.08537 [cs].
- [8] Mike Butters. Route Design and Selection. In *Pharmaceutical Process Development - Current Chemical and Engineering Challenges*, pages 90–116. Royal Society of Chemistry (RSC), 2011.
- [9] CAS. Scifinder.
- [10] Binghong Chen, Chengtao Li, Hanjun Dai, and Le Song. Retro\*: Learning retrosynthetic planning with neural guided a\* search. In *The 37th International Conference on Machine Learning (ICML 2020)*, 2020.
- [11] Ziqi Chen, Oluwatosin R. Ayinde, James R. Fuchs, Huan Sun, and Xia Ning. G2Retro as a two-step graph generative models for retrosynthesis prediction. *Communications Chemistry*, 6(1):102, May 2023.
- [12] E.J. Corey and Xue-Min Cheng. *The Logic of Chemical Synthesis*. Wiley, 1989.
- [13] Rémi Coulom. Efficient selectivity and backup operators in Monte-Carlo tree search. In *Proceedings of the 5th international conference on Computers and games*, CG’06, pages 72–83, Berlin, Heidelberg, May 2006. Springer-Verlag.
- [14] Sean Current, Ziqi Chen, Daniel Adu-Ampratwum, Xia Ning, and Srinivasan Parthasarathy. DIFFER: categorical diffusion ensembles for single-step chemical retrosynthesis. *Journal of Cheminformatics*, 17(1):112, July 2025.
- [15] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, et al. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, January 2025. arXiv:2501.12948 [cs].
- [16] Elsevier Limited. Reaxys.
- [17] Samuel Genheden, Amol Thakkar, Veronika Chadimová, Jean-Louis Reymond, Ola Engkvist, and Esben Bjerrum. AiZynthFinder: a fast, robust and flexible open-source software for retrosynthetic planning. *Journal of Cheminformatics*, 12(1):70, November 2020.
- [18] Gibson, Jack R. and Weber, Joanne D. Handbook of Selected Properties of Air- and Water-Reactive Materials. Technical Report 144, Defense Technical Information Center, Crane, Indiana, March 1969.
- [19] Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, et al. Towards an AI co-scientist, February 2025. arXiv:2502.18864 [cs].
- [20] Jeff Guo and Philippe Schwaller. It Takes Two to Tango: Directly Optimizing for Constrained

- Synthesizability in Generative Molecular Design, October 2024. arXiv:2410.11527 [q-bio].
- [21] Peter E. Hart, Nils J. Nilsson, and Bertram Raphael. A Formal Basis for the Heuristic Determination of Minimum Cost Paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–107, July 1968.
- [22] Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani, and others. Biomni: A General-Purpose Biomedical AI Agent. *bioRxiv*, pages 2025–05, 2025. Publisher: Cold Spring Harbor Laboratory.
- [23] Yuqing Huang, Rongyang Zhang, Xuesong He, Xuyang Zhi, Hao Wang, Xin Li, et al. ChemEval: A Comprehensive Multi-Level Chemical Evaluation for Large Language Models, September 2024. arXiv:2409.13989 [cs].
- [24] IARC. IARC monographs on the identification of carcinogenic hazards to humans. Technical report, IARC, 2025.
- [25] Ruofan Jin, Zaixi Zhang, Mengdi Wang, and Le Cong. STELLA: Self-Evolving LLM Agent for Biomedical Research, July 2025. arXiv:2507.02004 [cs].
- [26] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling Laws for Neural Language Models, January 2020. arXiv:2001.08361 [cs].
- [27] Junsu Kim, Sungsoo Ahn, Hankook Lee, and Jinwoo Shin. Self-Improved Retrosynthetic Planning. In *Proceedings of the 38th International Conference on Machine Learning*, pages 5486–5495. PMLR, July 2021. ISSN: 2640-3498.
- [28] T. L. Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, March 1985.
- [29] Michael S. Lajiness, Gerald M. Maggiora, and Veerabahu Shanmugasundaram. Assessment of the Consistency of Medicinal Chemists in Reviewing Sets of Compounds. *Journal of Medicinal Chemistry*, 47(20):4891–4896, September 2004.
- [30] Daniel Lowe. Chemical reactions from US patents (1976-Sep2016), 2017. Artwork Size: 1494665893 Bytes Pages: 1494665893 Bytes.
- [31] Andrew D. McNaughton, Gautham Krishna Sankar Ramalaxmi, Agustin Kruel, Carter R. Knutson, Rohith A. Varikoti, and Neeraj Kumar. CACTUS: Chemistry Agent Connecting Tool Usage to Science. *ACS Omega*, 9(46):46563–46573, November 2024. Publisher: American Chemical Society.
- [32] Mistral AI. Mistral NeMo, July 2024.
- [33] H. L. Morgan. The Generation of a Unique Machine Description for Chemical Structures-A Technique Developed at Chemical Abstracts Service. *Journal of Chemical Documentation*, 5(2):107–113, May 1965.
- [34] OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, et al. GPT-4o System Card, October 2024. arXiv:2410.21276 [cs].
- [35] Robert J. Ouellette and J. David Rawn. Aldehydes and Ketones. In *Organic Chemistry*, pages 629–657. Elsevier, 2014.
- [36] RDKit: Open-source cheminformatics.
- [37] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A Graduate-Level Google-Proof Q&A Benchmark. In *Proceedings of the Conference on Language Modeling (COLM)*, August 2024.
- [38] Marwin H. S. Segler, Mike Preuss, and Mark P. Waller. Planning chemical syntheses with deep neural networks and symbolic AI. *Nature*, 555(7698):604–610, March 2018. Publisher: Nature Publishing Group.
- [39] Kyle Swanson, Parker Walther, Jeremy Leitz, Souhrid Mukherjee, Joseph C Wu, Rabindra V Shivnaraine, et al. ADMET-AI: a machine learning ADMET platform for evaluation of large-scale chemical libraries. *Bioinformatics*, 40(7):btac416, July 2024.
- [40] Yuji Takaoka, Yutaka Endo, Susumu Yamanobe, Hiroyuki Kakinuma, Taketoshi Okubo, Youichi

- Shimazaki, et al. Development of a Method for Evaluating Drug-Likeness and Ease of Synthesis Using a Data Set in Which Compounds Are Assigned Scores Based on Chemists’ Intuition. *Journal of Chemical Information and Computer Sciences*, 43(4):1269–1275, July 2003.
- [41] UNECE. Globally Harmonized System of Classification and Labelling of Chemicals (GHS Rev. 10, 2023), 2023.
- [42] Rahul Verma, Yaru Jing, Honghu Liu, Varun Aggarwal, Harish Kumar Goswami, Ekta Bala, et al. Employing Ammonia for Diverse Amination Reactions: Recent Developments of Abundantly Available and Challenging Nitrogen Sources. *European Journal of Organic Chemistry*, 2022(25), July 2022. Publisher: Wiley.
- [43] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, et al. MMLU-Pro: A More Robust and Challenging Multi-Task Language Understanding Benchmark. *Advances in Neural Information Processing Systems*, 37:95266–95290, December 2024.
- [44] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Systems Processing Systems*, 2022.
- [45] Botao Yu, Frazier N. Baker, Ziqi Chen, Xia Ning, and Huan Sun. LlaSMol: Advancing Large Language Models for Chemistry with a Large-Scale, Comprehensive, High-Quality Instruction Tuning Dataset. In *Proceedings of the Conference on Language Modeling (COLM)*, August 2024.
- [46] Botao Yu, Frazier N. Baker, Ziru Chen, Garrett Herb, Boyu Gou, Daniel Adu-Ampratwum, Xia Ning, and Huan Sun. ChemToolAgent: The Impact of Tools on Language Agents for Chemistry Problem Solving. In *Findings of the Association for Computational Linguistics: NAACL 2025*, 2025.
- [47] Kevin Yu, Jihye Roh, Ziang Li, Wenhao Gao, Runzhong Wang, and Connor W. Coley. Double-Ended Synthesis Planning with Goal-Constrained Bidirectional Search. *Advances in Neural Information Processing Systems*, 37:112919–112949, December 2024.
- [48] Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, et al. Optimizing generative AI by backpropagating language model feedback. *Nature*, 639(8055):609–616, March 2025. Publisher: Nature Publishing Group.
- [49] Dengwei Zhao, Shikui Tu, and Lei Xu. Efficient retrosynthetic planning with MCTS exploration enhanced A\* search. *Commun Chem*, 7(1):1–12, March 2024. Publisher: Nature Publishing Group.
- [50] Mingchen Zhuge, Changsheng Zhao, Dylan Ashley, Wenyi Wang, Dmitrii Khizbullin, Yunyang Xiong, et al. Agent-as-a-Judge: Evaluate Agents with Agents, October 2024. arXiv:2410.10934 [cs].

## A Dataset

Table A1 shows the benchmark dataset curated for our experiments. Tasks are provided as natural language prompts, organized by constraint type.

Table A1: Constrained Retrosynthesis Planning Dataset

ID	Constraint Type	Task
C1	Carcinogen	Find the shortest synthesis route for <chem>O=S(=O)(C#Cc1ccc(C1)cc1)N1CCNCC1</chem> that avoids all carcinogenic substances.

*Continued on next page*

ID	Constraint Type	Task
C2	Carcinogen	Identify a synthesis route for <chem>CC(C)(C)OC(=O)NC1(c2nc(O)c3c(Cl)ccn3n2)CC1</chem> that avoids known or suspected carcinogens.
C3	Carcinogen	Plan the synthesis of <chem>CC(=O)OCc1nc2cnc(Br)cc2n1C(C)(C)COC(C)=O</chem> , avoiding any carcinogenic substances.
C4	Carcinogen	Plan the synthesis of <chem>FC(F)(F)Cn1ncnc1-c1cc2n(n1)-c1cc(C3CCNCC3)ccc1OCC2</chem> , avoiding any known or suspected carcinogens.
C5	Carcinogen	Perform synthesis planning for <chem>C=C(C[C@@H](Cc1ccc(-c2ccccc2)cc1)NC(=O)OC(C)(C)C)C(=O)O</chem> without using any carcinogenic substances.
C6	Carcinogen	Provide a synthesis route for <chem>C0c1cc2c(Oc3cc(C)c(C)nc3-c3cccc(C)n3)ccnc2cc1OCCNCCO</chem> that avoids carcinogens.
C7	Carcinogen	Find the shortest synthesis route possible for <chem>C#CC1(O)C(C)=CC2(CC1(C)C(F)(F)F)OC(C)C(C)O2</chem> without using carcinogens.
C8	Carcinogen	Identify the shortest possible synthetic route for <chem>ClC1ccc2c(c1)Nc1ncnc1S2</chem> that avoids carcinogens.
C9	Carcinogen	Identify the best synthesis route for <chem>COC(=O)c1ccc2c(c1)C=CC(=C(Cl)Cl)CO2</chem> that avoids carcinogenic materials.
C10	Carcinogen	Design a synthesis plan for <chem>C[C@@H](O)c1nc2cnc3ccsc3c2n1[C@H]1CC[C@H](CO)CC1</chem> without using any carcinogens.
C11	Carcinogen	Provide the shortest synthesis route for <chem>Cn1oc(=O)nc1/C(=N\OCc1cccc(N)n1)c1cccc1</chem> that does not use any carcinogens.
C12	Carcinogen	Plan the synthesis of <chem>CC(C)(CO)n1c(CO)nc2cnc(Br)cc21</chem> without using carcinogens.
C13	Carcinogen	Plan the synthesis of <chem>CCCC[Sn](/C=C/C1(O)C(C)=CC2(CC1(C)C(F)(F)F)OC(C)C(C)O2)(CCCC)CCCC</chem> , but do not use any carcinogens in your synthesis route.
C14	Carcinogen	Find the shortest synthesis path for <chem>COC(=O)CCc1cc2cc(-c3noc(-c4ccc(O(C)C)c(Cl)c4)n3)ccc2n1C</chem> that does not use any carcinogenic substances.
C15	Carcinogen	Plan the synthesis of <chem>CC(C)(C)OC(=O)N1CC=C(c2ccc3c(c2)-n2nc(-c4ncnn4CC(F)(F)F)cc2CCO3)CC1</chem> without using any carcinogens.
C16	Carcinogen	Find the shortest synthesis route for <chem>C[Si](C)(C)CCOCn1cc(C2CCc3c(C(=O)O)nn(COCC[Si](C)(C)C)c3C2)cn1</chem> that doesn't use any carcinogenic substances.
C17	Carcinogen	Plan the synthesis of <chem>O=C(OCc1cccc1)N1CC[C@H]2CCCN(Cc3cccc3)C[C@H]21</chem> without using carcinogens.
C18	Carcinogen	Find a synthesis route for <chem>Cn1oc(=O)nc1/C(=N\OCc1cccc(NC(=O)OCCc2cccc2)n1)c1cccc1</chem> that avoids all known or suspected carcinogens.
C19	Carcinogen	Plan the shortest synthesis route for <chem>CC(C)(C)OC(=O)NC1(c2nc(NC3ccccn3)c3c(Cl)ccn3n2)CC1</chem> . Do not use any carcinogens in the route.

*Continued on next page*

ID	Constraint Type	Task
C20	Carcinogen	Generate a synthesis plan for the following compound without using any carcinogens: <chem>CC[C@@H](OC(=O)c1ccccc1)[C@H]1CCCN(C(=O)OC(C)(C)C)C1</chem>
C21	Carcinogen	Plan the synthesis of <chem>CC(C)(C)OC(=O)N[C@@H]1c2cccnc2[C@H](O)CC[C@H]1c1cccc(F)c1F</chem> . Avoid carcinogenic materials.
C22	Carcinogen	Find the shortest synthesis route that doesn't use carcinogens for the following product: <chem>CC(C)(C)OC(=O)N[C@@H]1c2cccnc2[C@H](N)CC[C@H]1c1cccc(F)c1F</chem> .
C23	Carcinogen	Design a synthesis path for <chem>CC(C)(C)OC(=O)N[C@@H]1c2cccnc2C(=O)CC[C@H]1c1cccc(F)c1F</chem> . Do not use any carcinogens.
C24	Carcinogen	Plan the shortest possible synthesis route for <chem>C0CCCc1cc(CN(C(=O)[C@H]2CNCC[C@@H]2c2ccc(OCCOc3c(C1)cc(C)cc3C1)cc2)C2CC2)cc(OCCOC)c1</chem> , but do not use any carcinogens in the route.
C25	Carcinogen	Find the shortest synthesis path for <chem>C0CCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(C1)cc(C)cc3C1)cc2)C2CC2)cc(OCCOC)c1</chem> . Do not use any carcinogens.
C26	Carcinogen	Generate a synthesis plan without carcinogens for <chem>CC#CCn1c(Br)nc(C=O)c1C(=O)OC</chem> .
C27	Carcinogen	Find the shortest synthesis path for <chem>COC(=O)c1ccc2c(c1)C=CC(=CC1)C02</chem> that does not use carcinogenic substances.
C28	Carcinogen	Find a synthesis path for <chem>O=C(Nc1cccc(C1)c1)N1CCc2[nH]nc(C(=O)N3CC(F)C03)c2C1</chem> that doesn't use carcinogens.
P1	Pyrophoric	Identify a synthesis route for <chem>C[C@H](O[Si](C)(C)C(C)(C)C)[C@@H]1CC(=O)CC(C)(C)N1</chem> that does not use pyrophoric substances.
P2	Pyrophoric	Find the shortest synthesis plan for <chem>C[C@H](c1ccccc1)N1C[C@]2(C(=O)OC(C)(C)C)C=CC[C@@H]2C1=S</chem> that avoids all pyrophoric and water-reactive substances.
P3	Pyrophoric	Plan the shortest synthesis route for <chem>CC(=O)c1ccc2c(c1)C=CC(O)(CO)C02</chem> without using any pyrophoric or water reactive substances.
P4	Pyrophoric	Synthesize <chem>CC[C@@H](OC(=O)c1ccccc1)[C@H]1CCCN(C(=O)OC(C)(C)C)C1</chem> without using any pyrophoric or water-reactive substances.
P5	Pyrophoric	Perform synthesis planning for <chem>CC1=NC2(N=C1N)c1cc(Br)ccc1CCC21CC1</chem> , avoiding pyrophoric materials (substances that ignite in moisture or air) in your synthesis route.
P6	Pyrophoric	Plan a synthesis route for <chem>C0c1cc2ncc3c(N)nc(-c4cncc(OCCN(Cc5ccc(F)cc5)C(=O)OC(C)(C)C)c4)cc3c2cc1OC</chem> that uses no pyrophoric or water-reactive substances.

*Continued on next page*

ID	Constraint Type	Task
P7	Pyrophoric	Perform synthesis planning for <chem>O[C@H]1C[C@H](c2cnn3c(N[C@H]4CCc5ccccc54)ncnc23)C=C1COCc1cccc1</chem> without using pyrophoric substances in your synthesis plan.
P8	Pyrophoric	Synthesize <chem>CC(=O)N1c2ccc(N3CCNCC3)cc2[C@H](Nc2ccccc2)[C@@H](C)[C@@H]1C</chem> without using any pyrophoric or water-reactive reagents.
P9	Pyrophoric	Find the shortest synthesis route for <chem>COc1cc2c(=O)[nH]c(=O)n([C@@H]3O[C@H](CO)[C@H]4OC(C)(C)O[C@H]43)c2cc1OC</chem> , avoiding pyrophoric substances.
P10	Pyrophoric	Identify a synthesis route for <chem>Oc1ccc2c3c(ccc2c1)Cc1cccc1OC3c1ccc(OCCN2CCCC2)cc1</chem> that does not use pyrophoric substances.
P11	Pyrophoric	Plan the synthesis of <chem>C[C@@H](O)C[C@H]1OC[C@@H](C2CCCC2)N(c2cc(C#CC(C)(C)C)sc2C(=O)O)C1=O</chem> . Do not use any pyrophoric substances.
P12	Pyrophoric	Find the shortest synthesis route for <chem>OC[C@H]1C[C@@H](c2cnn3c(N[C@H]4CCc5ccccc54)ncnc23)C[C@@H]1O</chem> that avoids using any pyrophoric substances.
S1	User-Specified	Find the shortest synthesis path for <chem>C[C@@H]1CCCN1CCc1nnc2cc(Br)ccc2c1O</chem> , but avoid using <chem>C=C[Sn](CCCC)(CCCC)CCCC</chem> in your synthesis route.
S2	User-Specified	Plan a synthesis route for <chem>CC(=O)NC[C@H]1CN(c2ccc3c(c2)CCCc2c(C(C)C)n[nH]c2-3)C(=O)O1</chem> . Avoid using phosgene ( <chem>O=C(Cl)Cl</chem> ) in your synthesis.
S3	User-Specified	Find the shortest synthesis route for <chem>C[C@@H]1CNC(=O)c2cc3cc(OCCCN4CCCC4)ccc3n21</chem> that does not use trimethyl borate ( <chem>COB(OC)OC</chem> ).
S4	User-Specified	Plan the synthesis of <chem>COCCc1cc(CN(C(=O)[C@H]2CNCC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCCOC)c1</chem> without using hexane ( <chem>CCCCCC</chem> ).
S5	User-Specified	Identify a synthesis plan for <chem>COc1cc2c(=O)[nH]c(=O)n([C@@H]3O[C@H](CO)[C@H]4OC(C)(C)O[C@H]43)c2cc1OC</chem> that does not use methanol ( <chem>CO</chem> ).
S6	User-Specified	Plan the synthesis of <chem>CC(C)c1ccc2c(c1)OC1(O)c3cccc3C(=O)C21NC(=O)c1cc(-c2ccccc2)n[nH]1</chem> . Do not use phenol ( <chem>Oc1ccccc1</chem> ).
S7	User-Specified	Synthesize <chem>CCOC(=O)/C(N)=N/Nc1cc(Cl)ccc1[N+](=O)[O-]</chem> without using nitric acid ( <chem>O=[N+]( [O-] )O</chem> ).
S8	User-Specified	Find the shortest synthesis route for <chem>C#CC1(O)C(C)=CC2(CC1(C)C(F)(F)F)OC(C)C(C)O2</chem> that does not use toluene ( <chem>Cc1ccccc1</chem> ).

## B LARC System Instructions

Here we present the system instructions for implementing the EVALUATOR in LARC<sub>Mistral</sub> and LARC<sub>Claude</sub>. Each set of system instructions was carefully engineered to elicit good instruction following from MISTRAL NEMO and CLAUDE 3.5 SONNET, respectively.

## B.1 LARC<sub>Mistral</sub> System Instructions

The following system instructions are used to implement the EVALUATOR in LARC<sub>Mistral</sub>. Values surrounded by curly braces {} are replaced with task-specific information. Please note, reaction evaluation is broken into two separate instructions: one to call tools and another to generate a score based on the results.

**Evaluation planning.** This set of instructions instructs MISTRAL NEMO to perform evaluation planning, generating the evaluation instructions.

You are an expert chemist evaluating a reaction. Your task is to use tools to determine if the user’s provided reaction SMILES (reactants>>product) satisfies only the following constraint:

To help you evaluate the constraint, you may use any of the following tools:

- AIExpert(question): Asks a self-contained question of an AI chemistry expert to get information not available from the other tools or infer reaction conditions (e.g. catalysts, temperature, etc.).
- Carcinogenicity(SMILES1, SMILES2, ...): returns a probabilistic score for each molecule indicating whether the provided SMILES represents a carcinogenic molecule. Can be called with one or more molecules; batching improves efficiency.
- Pyrophoricity(SMILES1, SMILES2, ...): returns a score for each molecule indicating the Tanimoto similarity between the provided SMILES and its nearest neighbor on an authoritative list of pyrophoric substances (materials that spontaneously ignite in air or water). A score of 1.000 indicates the substance is on the list. Can be called with one or more molecules. The term "pyrophoric" includes water-reactive substances, as these may ignite in air due to the presence of water vapor.
- Similarity(SMILES1, SMILES2): returns a similarity between exactly 2 molecules, represented by SMILES.

Create a plan of the actions you would take to evaluate this specific constraint for any user-provided reaction. Use as few actions as possible to evaluate the constraint, it is acceptable to use a single action. Explicitly identify actions taken using the following template:

```

...
ActionName(comma-separated action inputs)
ActionName(comma-separated action inputs)
...

```

Action names must be one of: AIExpert, Carcinogenicity, Pyrophoricity, Similarity.

In place of actual inputs, place descriptions of the values you would use. Be efficient and use only the tools necessary. Then, create a rubric for assigning a score of 1-5 where 5 is assigned if the constraint is completely satisfied and 1 is assigned if the constraint is completely violated.

Constraint: {constraint}

**Reaction evaluation 1.** This set of instructions instructs MISTRAL NEMO to follow the evaluation instructions to select tools to evaluate a specific reaction. It uses the messages from evaluation planning (including the generated evaluation instructions) as context.

```

{evaluation planning context}
  User-provided reaction:
  Reactants: {reactants}
  Product: `{product}`

```

When multiple reactants are present, they are separated by dots.  
 Fill in your plan with the specific inputs for this reaction. Format your actions using the following format, one per line, with inputs wrapped in quotes:  
 ...  
 ActionName(ActionInput1, ActionInput2, ...)  
 ActionName(ActionInput1, ActionInput2, ...)  
 ...  
 Do not include any other information. Await the results.

**Reaction evaluation 2.** This set of instructions instructs MISTRAL NEMO to generate a score between 1 and 5. It uses the evaluation planning messages, the messages from Reaction Evaluation 1 and the tool outputs as context.

{evaluation planning context}  
 {reaction evaluation context}  
 {tool outputs}  
 Analyze these results and provide your final score. Use the following format:  
 Reasoning: (your analysis)  
 Answer: N  
 where N is between 1 and 5.

## B.2 LARCClaude System Instructions

The following system instructions are used to implement the EVALUATOR in LARCClaude. Values surrounded by curly braces {} are replaced with task-specific information. Please note, reaction evaluation is broken into two separate instructions: one to call tools and another to generate a score based on the results.

**Evaluation planning.** This set of instructions instructs CLAUDE 3.5 SONNET to perform evaluation planning, generating the evaluation instructions.

You are an expert chemist assisting in retrosynthesis planning. You aim to determine if the user’s provided reaction (reactants>>product) satisfies the following constraint:  
 {constraint}  
 Note that there may be reagents omitted from the user’s provided reaction.  
 To help you answer the user question, you have access to these tools:

- Carcinogenicity(SMILES1, SMILES2, ...): returns a probabilistic score for each molecule indicating whether the provided SMILES represents a carcinogenic molecule. Can be called with one or more molecules; batching improves efficiency.
- Pyrophoricity(SMILES1, SMILES2, ...): returns a score for each molecule indicating the Tanimoto similarity between the provided SMILES and its nearest neighbor on an authoritative list of pyrophoric substances (materials that spontaneously ignite in air or water). A score of 1.000 indicates the substance is on the list. Can be called with one or more molecules. The term "pyrophoric" includes water-reactive substances, as these may ignite in air due to the presence of water vapor.
- Similarity(SMILES1, SMILES2): returns a similarity between exactly 2 molecules, represented by SMILES.
- AIExpert(question): Asks a self-contained question of an AI chemistry expert, for instance,

to infer reaction conditions or reagents not listed in the reaction, or to get information not available from the other tools. Include all relevant details, including reaction information and SMILES strings.

- Answer(value): Answer with a score of 1-5 based on only the constraint above is satisfied, where 5 is completely satisfied and 1 is not at all satisfied. Keep in mind these scores will be used to prioritize further planning.

Create a plan of the actions you would take to evaluate this specific constraint for any user-provided reaction. Explicitly identify actions taken using the following template:

Action: ActionName(ActionInput)

In place of actual inputs, place descriptions of the values you would use. Use only the tools necessary. Think step by step.

**Reaction evaluation 1.** This set of instructions instructs CLAUDE 3.5 SONNET to follow the evaluation instructions to select tools to evaluate a specific reaction.

You are an expert chemist assisting in retrosynthesis planning. You aim to determine if the user's provided reaction (reactants>>product) satisfies the following constraint:

{constraint}

To help you answer the user question, you have access to these tools:

- Carcinogenicity(SMILES1, SMILES2, ...): returns a probabilistic score for each molecule indicating whether the provided SMILES represents a carcinogenic molecule. Can be called with one or more molecules; batching improves efficiency.
- Pyrophoricity(SMILES1, SMILES2, ...): returns a score for each molecule indicating the Tanimoto similarity between the provided SMILES and its nearest neighbor on an authoritative list of pyrophoric substances (materials that spontaneously ignite in air or water). A score of 1.000 indicates the substance is on the list. Can be called with one or more molecules. The term "pyrophoric" includes water-reactive substances, as these may ignite in air due to the presence of water vapor.
- Similarity(SMILES1, SMILES2): returns a similarity between exactly 2 molecules, represented by SMILES.
- Answer(value): Answer with a score of 1-5 based on only the constraint above is satisfied, where 5 is completely satisfied and 1 is not at all satisfied. Keep in mind these scores will be used to prioritize further planning.

In one response, specify all of the actions you would take to gather the necessary information to evaluate the user's reaction. Follow this plan:

{evaluation instructions}

Follow the plan step by step. Replace all references to the AIExpert tool with your own expert chemistry knowledge on the user's reaction. Keep in mind there may be multiple correct answers to each question (e.g. multiple ways to catalyze a reaction), and you should discuss every possibility in detail. Reactants and products provided by the user should be considered required for the reaction, changing these would be considered a different reaction. Provide accurate and valid SMILES representations for molecules in your answer.

For the tools, construct a single, unified code block wrapped in triple backticks at the end of your response, including specific inputs based on the user's provided reaction. Ignore the Answer tool until you have results from your code block actions. Wrap all SMILES in backticks.

**Reaction evaluation 2.** This set of instructions instructs CLAUDE 3.5 SONNET to generate a score between 1 and 5. It uses the messages from Reaction Evaluation 1 and the tool outputs as context.

```

{reaction evaluation context}
{tool outputs}
Given this information, please provide a final score for the reaction using ``Answer(X)``,
where X is your score between 1-5 on only this constraint:
{constraint}
Do not worry about any other constraints, as we may assess these separately.

```

## C Additional Related Work

**LLM agents for chemistry.** LLM agents have recently shown great promise in chemistry applications. CoSCIENTIST [5] and CHEMCROW [6], are notable examples, combining LLM reasoning with cheminformatics tools, web search, and robotic laboratory equipment to perform chemistry tasks, including retrosynthesis planning. However, neither of these methods are well-equipped for constrained retrosynthesis planning. CoSCIENTIST relies on general-purpose LLMs for retrosynthesis planning, which lack accuracy on chemical reaction tasks compared to models with specialized training [45]. CHEMCROW generates synthetic routes using an external tool for unconstrained retrosynthesis planning, making it incapable of incorporating constraints directly into the planning process. CACTUS [31] and CHEMTOOLAGENT [46] are similar to CHEMCROW, but with different focuses. CACTUS focuses only on molecule property prediction with in silico tools, and CHEMTOOLAGENT focuses on analyzing the impact of tool use across various chemistry tasks. [19] introduces an LLM agent co-scientist similar to CoSCIENTIST, but designed for broad applicability across multiple scientific disciplines, not specifically for chemistry or retrosynthesis planning. LIDDIA [2] is an LLM agent for in silico drug discovery that automates generation, screening, and optimization of molecules, but it does not consider retrosynthesis planning. While these efforts demonstrate the promise of LLM agents in chemistry applications, none of the existing LLM agents can perform constrained retrosynthesis planning.

Additional LLM-based efforts, such as TEXTGRAD [48], STELLA [25], and BIOMNI [22], have also been used in biomedical and chemistry applications; however, none of these methods are equipped for retrosynthesis planning tasks.

## D Implementation Details for Reproducibility

For the experiments in this paper, we used the following hyperparameters and settings. To control cost and runtime, we enforced a limit on the number of expansions and evaluations. The search terminates after 500 expansions, following the established convention in multi-step retrosynthesis literature [10, 49, 27]. Additionally, after 300 evaluations, all subsequent evaluations were replaced with the optimistic default score (i.e. 5 on the 1-5 scale). To unify scaling, the raw values from MEEA\*'s  $V_A^*$  are min-max normalized to a  $[0, 1]$  scale at each selection step. The MEEA\* simulation step already included min-max scaling, so no change was required for  $V_{MCTS}$ . Additionally, the evaluator scores  $S(r)$  are each normalized to a  $[0, 1]$  scale, where 0 indicates complete constraint violation and 1 indicates complete constraint satisfaction. For hyperparameters, we used  $\lambda = 2$  to ensure the constraint was followed. While  $\lambda$  can be tuned as needed in practice, we observed during development that LARC<sub>Mistral</sub> had relatively high success rates for of  $\lambda \in \{2, 3\}$ , moderate success rates for  $\lambda \in \{4, 5\}$ , and low success rate with  $\lambda = 1$ . Additionally, planning time generally increased with larger  $\lambda$ , thus our choice of  $\lambda = 2$ . We set  $K = 5$  to control the number of candidate routes subject to evaluation. Following the MEEA\* paper, we used a UCB scaling term of 4.

**Algorithm 1** LARC

---

**Require:** Target molecule  $p$ ; constraints  $c$ ; commercially available molecule set  $B$ ;  
 unconstrained MCTS value function  $V_{MCTS}$ ; unconstrained  $A^*$  value function  $V_{A^*}$ ;  
 expansion function EXPAND;  
 expansion limit  $N_{exp}$ ; evaluation limit  $N_{eval}$ ;  
 number of simulations  $K$ ; constraint weight  $\lambda$ ; default score  $S_{def}$

**Ensure:** Synthetic route  $R_*$

```

1:  $n_{exp} \leftarrow 0$  ▷ Initialize expansion count
2:  $T \leftarrow \text{Tree}(\text{nodes}=\{p\}, \text{edges}=\{\})$  ▷ Initialize search tree
3:  $S \leftarrow \{\}$  ▷ Initialize reaction evaluations

4:  $P \leftarrow \text{EVALUATOR-plan}(c)$  ▷ Evaluation planning generates evaluation instructions  $P$ 

5: while  $R_*$  is empty  $\wedge n_{exp} < N_{exp}$  do
6:    $\mathcal{C} \leftarrow \text{SYNTHESIZER-MCTS}(T, V_{MCTS}, S, S_{def}, \lambda, K)$  ▷ Simulate to get candidate set  $\mathcal{C}$ 
7:
8:   for all  $(m, R) \in \mathcal{C}$  do
9:     for all  $r \in R$  do
10:      if  $|S| < N_{eval} \wedge r \notin S$  then
11:         $S(r) \leftarrow \text{EVALUATOR-evaluate}(r, P)$  ▷ Reaction evaluation
12:      else if  $r \notin S$  then
13:         $S(r) \leftarrow S_{def}$  ▷ Use default after evaluation limit
14:      end if
15:    end for
16:  end for

17:   $(m_{exp}, R_{exp}) \leftarrow \text{SYNTHESIZER-A}^*(\mathcal{C}, V_{A^*}, S, \lambda)$  ▷ Selection of  $(m_{exp}, R_{exp})$  from  $\mathcal{C}$ 

18:   $\mathcal{R} \leftarrow \text{EXPAND}(T, R_{exp}, m_{exp})$  ▷ Expansion updates  $T$  with set of reactions  $\mathcal{R}$ 
19:   $n_{exp} \leftarrow n_{exp} + 1$ 

20:  for all  $r \in \mathcal{R}$  do
21:    if  $R_{exp} \oplus r$  is a complete route from  $B$  to  $p$  then ▷ Check for complete route
22:       $R_* \leftarrow R_{exp} \oplus r$ 
23:    end if
24:  end for
25: end while

26: return  $R_*$ 

```

---

The algorithm for LARC is presented in Algorithm 1, illustrating how LARC uses both the EVALUATOR and SYNTHESIZER to perform constrained retrosynthesis planning. LARC’s EVALUATOR has two steps: evaluation planning (denoted EVALUATOR-plan in Algorithm 1) and reaction evaluation (denoted EVALUATOR-evaluate). Both of these steps are implemented using the system instructions. LARC’s SYNTHESIZER adapts MEEA\*’s two step process:

MCTS simulation (denoted SYNTHESIZER-MCTS) to generate  $K$  candidate routes, and A\* search (denoted SYNTHESIZER-A\*) to select the best of the  $K$  candidates. LARC’s constrained adaptations of these steps are given in Algorithms 2 and 3, respectively. LARC’s SYNTHESIZER uses the same expansion function (denoted EXPAND in Algorithm 1) as MEEA\* [49]. Please note, these algorithms assume  $V_{A^*}$  and  $S(r)$  are already normalized, per the implementation details above. A full implementation of LARC can be found at <https://github.com/ninglab/LARC>.

---

**Algorithm 2** SYNTHESIZER-MCTS
 

---

**Require:** search tree  $T$ ; unconstrained value function (includes UCB)  $V_{MCTS}$ ;  
EVALUATOR scores  $S$ ; default score  $S_{\text{def}}$  constraint weight  $\lambda$ ; number of simulations  $K$ ;

**Ensure:** expansion candidates  $\mathcal{C} = \{(m, R), \dots\}$  of intermediate-route pairs

```

1:  $\mathcal{C} \leftarrow \{\}$ 
2:  $V'_{MCTS}(m, R) := V_{MCTS}(m, R) + \lambda \sum_{r \in R} (S(r) \text{ if } r \in S \text{ else } S_{\text{def}})$  ▷ Equation 1

3: for  $k = 1$  to  $K$  do
4:    $R \leftarrow \langle \rangle$ ;  $m \leftarrow \text{root}(T)$  ▷ start from target molecule and empty route
5:    $A \leftarrow$  reactions that produce  $m$  in  $T$ 

6:   while  $|A| > 0$  do ▷ Simulate route from reactions in  $T$ 
7:      $r^*, m^* \leftarrow \arg \max_{r \in A, m \in \text{reactants}(r)} V'_{MCTS}(m, R \oplus r)$  ▷ select reaction
8:      $R \leftarrow R \oplus r^*$ ;  $m \leftarrow m^*$  ▷ update route
9:      $A \leftarrow$  reactions that produce intermediate  $m$  in  $T$ 
10:  end while

11:   $\mathcal{C} \leftarrow \mathcal{C} \cup \{(m, R)\}$  ▷ add candidate
12: end for

13: return  $\mathcal{C}$ 

```

---



---

**Algorithm 3** SYNTHESIZER-A\*
 

---

**Require:** candidate routes  $\mathcal{C}$ ; unconstrained value function  $V_{A^*}$ ;  
EVALUATOR scores  $S$ ; constraint weight  $\lambda$ ;

**Ensure:** intermediate molecule  $m_*$ , route  $R_*$

```

1:  $\mathcal{C} \leftarrow \{\}$  ▷ Initialize empty candidate set
2:  $V'_{A^*}(m, R) := V_{A^*}(m, R) + \lambda \sum_{r \in R} S(r)$  ▷ Equation 1

3:  $m_*, R_* \leftarrow \arg \max_{(m, R) \in \mathcal{C}} V'_{A^*}(m, R)$  ▷ Select best candidate based on constrained A* policy
4: return  $m_*, R_*$ 

```

---

## E Additional Experimental Settings

**Constrained retrosynthesis planning by LLMs.** We use general-purpose LLMs, including CLAUDE 3.5 SONNET [1], GPT-4O [34], DEEPSEEK R1 [15], and MISTRAL NEMO [32] as the baselines for constrained retrosynthesis planning. CLAUDE 3.5 SONNET and GPT-4O are representative state-of-the-art closed-source LLMs, which have both shown promising results on chemistry knowledge benchmarks [43, 37, 46]. DEEPSEEK R1 serves as a representative of open-source reasoning LLMs, leveraging an internal chain of thought [44] process to improve its responses. MISTRAL NEMO serves as a representative of small, open-source LLMs, selected due to its instruction-following capabilities [32]. Note that baselines do not include models specifically designed for unconstrained retrosynthesis planning, such as MEEA. This is because the benchmark is constructed by selecting the tasks for which MEEA can generate valid routes but violate the constraints. As MEEA is the state of the art for unconstrained retrosynthesis planning, we believe other unconstrained retrosynthesis planning methods [10, 27, 17, 38] will not succeed on the benchmark data.

## F Additional Results

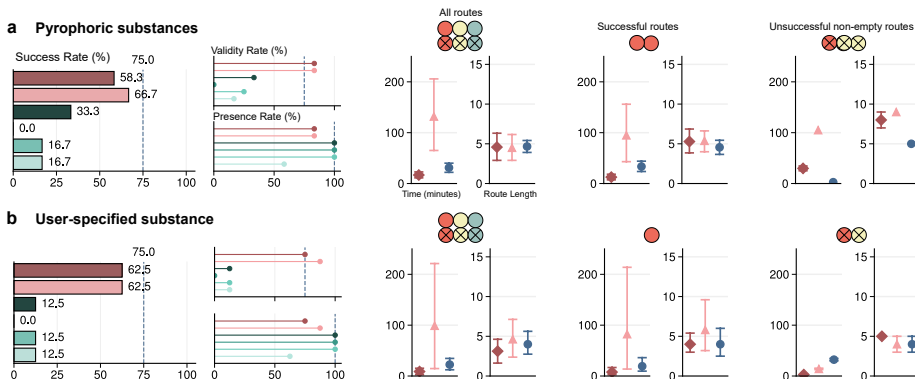


Figure F4: Evaluation **a,b** success rate, validity rate, presence rate, planning time, and route length for the benchmark tasks, organized by constraints for avoiding pyrophoric substances, and a user-specified substance, respectively. The symbols and colors in this figure follow the same conventions as Figure 2.

### F.1 Pyrophoricity-constrained retrosynthesis planning

Figure F4a presents the results for retrosynthesis planning tasks that are constrained to avoid pyrophoric substances.  $LARC_{Mistral}$  and  $LARC_{Claude}$  continue to achieve high success rates of 58.3% and 66.7%, respectively.  $LARC_{Claude}$  moderately underperforms EXPERT, who secures a success rate of 75.0%. Meanwhile,  $LARC_{Mistral}$  is significantly faster, with 12.58 minutes on average to generate one successful route, compared to EXPERT (33.48 minutes) ( $p=0.020$ , two-sided two-sample t-test) and  $LARC_{Claude}$  (94.81 minutes) ( $p=0.026$ ). This suggests that, for retrosynthesis planning avoiding pyrophoricity, LARC is still behind human experts, but could be a decent option with satisfactory/acceptable success rates, considering its efficiency and automation. LARC can also serve as an alternative to human experts, and provide additional, different solutions, which human experts can further select, utilize, or optimize based on their domain knowledge.

The LLM baselines fall short on these tasks, with the best LLM baseline, CLAUDE 3.5 SONNET, achieving only 33.3% success rate. Pyrophoric substances can vary greatly in their composition and the underlying mechanism of pyrophoricity [18], demanding a nuanced understanding of chemical reactivity. These results demonstrate LARC’s ability and efficiency to avoid a broad set of dangerous materials in its retrosynthesis planning compared to general-purpose LLMs.

## F.2 User-specified constrained retrosynthesis planning

As Figure F4b shows, LARC<sub>Claude</sub> and LARC<sub>Mistral</sub> achieve the same success rate of 62.5% when planning 8 routes that must each avoid a single, user-specified substance. LLMs still struggle, with a very low success rate of 12.5%, due to the lack of ability to generate valid synthetic routes. On the other hand, LARC generalizes well to the highly diverse constraint types, fulfilling a key need in the dynamic context of retrosynthesis planning. EXPERT achieves a higher success rate of 75.0% with an average 22.47 minutes on generating one route. LARC<sub>Mistral</sub> is still faster on average than EXPERT (8.82 minutes), though the difference is not very significant ( $p=0.089$ , two-sided two-sample t-test).

## F.3 Case Studies

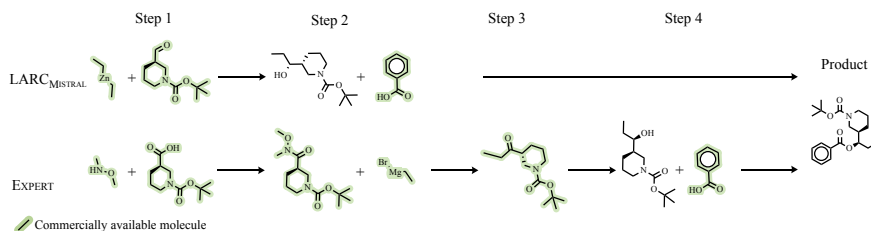


Figure F5: Synthesis route comparison of LARC<sub>Mistral</sub> against EXPERT tert-butyl (R)-3-((R)-1-(benzoyloxy)propyl)-piperidine-1-carboxylate.

### F.3.1 LARC<sub>Mistral</sub> can generate better routes than human experts

Figure F5 compares routes from LARC<sub>Mistral</sub> and EXPERT for the synthesis of tert-butyl (R)-3-((R)-1-(benzoyloxy)propyl)-piperidine-1-carboxylate. LARC<sub>Mistral</sub> proposed a very concise two-step reaction to the product using commercially available reagents: diethyl zinc, tert-butyl (S)-3-formylpiperidine-1-carboxylate, and benzoic acid. Its Step 1 involves a stereoselective, nucleophilic addition of an ethyl group to tert-butyl (S)-3-formylpiperidine-1-carboxylate using the diethylzinc reagent. This results in an alcohol that can be converted to the product (an ester) using benzoic acid in Step 2.

In contrast, EXPERT proposed a four-step reaction to the same product. Step 1 involves conversion of commercially available (S)-1-(tert-butoxycarbonyl)piperidine-3-carboxylic acid to the Weinreb amide intermediate, tert-butyl (S)-3-(methoxy(methyl)carbamoyl)piperidine-1-carboxylate. This step was unnecessary, as its product is also commercially available. Step 2 is the Grignard addition of an ethyl group to obtain tert-butyl (S)-3-propionylpiperidine-1-carboxylate. Stereoselective reduction of the ketone in Step 3 results in the alcohol intermediate obtained in LARC<sub>Mistral</sub>’s Step 1. A final ester formation, similar to that proposed by LARC<sub>Mistral</sub>, yields the product.

One potential reason for EXPERT’s longer route is the EXPERT’s inherent bias towards the stability of some functional groups and their commercial availability. For example, it is a generally accepted notion that aldehydes are not very stable [35], so chemists tend to synthesize aldehydes and use them as needed. In most cases, an aldehyde equivalent such as the Weinreb amide, which was proposed in Step 1 by EXPERT, is used instead, followed by a reduction step when the product is an alcohol. Thus, as the Weinreb amide is stable and easily synthesizable, EXPERT may tend to propose a synthesis of this intermediate in the synthetic route. This may be another reason why EXPERT proposed this route, relying on their familiarity with this synthesis strategy rather than referring to the commercial availability of proposed intermediates. In contrast, LARC<sub>Mistral</sub> benefits from SYNTHESIZER’s thorough checks for commercial availability while simultaneously relying on feedback from the EVALUATOR to ensure the constraint is satisfied.

## F.4 Study on tooling in LARC

For LARC<sub>Mistral</sub>, tools are critical: across all 48 benchmark tasks, LARC<sub>Mistral</sub> achieves a success rate of 72.9% using tools in addition to EVALUATOR’s internal knowledge, and 45.8% without tooling. Specifically, LARC<sub>Mistral</sub> with tooling generates 43 valid routes (validity rate of 89.6%), of which 35 are successful, and without tooling, 41 valid routes (validity rate of 85.4%), of which 22 are successful. This indicates that tooling primarily impacts LARC through deliberately assessing and enforcing constraint satisfaction to ensure successful routes. Moreover, with tooling, LARC<sub>Mistral</sub> generates successful routes faster (14.53 minutes on average) than without tooling (22.08 minutes), whereas in the latter case, EVALUATOR has to act as the AI expert and conduct reasoning via MISTRAL NEMO.

EVALUATOR’s interpretation of results from different tools also significantly impacts LARC performance. For example, in pyrophoricity-constrained retrosynthesis planning, LARC<sub>Claude</sub> performs better than LARC<sub>Mistral</sub> (success rate of 66.7% vs 58.3%). One reason for this lies in how CLAUDE 3.5 SONNET- and MISTRAL NEMO-based EVALUATOR interprets and leverages evaluation results from the pyrophoricity predictor. Although with similar instructions, CLAUDE 3.5 SONNET considers a high value (e.g., close to 1.0) from the predictor as an indicator of high pyrophoricity, and thus, a low  $S(r)$ . This allows SYNTHESIZER to explore more extensively to avoid pyrophoricity, and thus, a high success rate in the generated routes. However, LARC<sub>Mistral</sub> could interpret a relatively low value (e.g., 0.333) for high pyrophoricity, thus discouraging the exploration of routes that satisfy the constraint, resulting in more failed routes.

## G Additional Discussion

Constrained retrosynthesis planning is highly challenging in the practice of synthetic chemistry, as it requires finding synthetic routes that not only lead to the target molecule but also satisfy additional, often complex, user-specified requirements. These constraints, such as limiting the number of steps, avoiding specific reagents or reaction types, or adhering to cost, safety, or environmental guidelines, can drastically narrow and fragment the feasible search space. Moreover, verifying constraint satisfaction is non-trivial: it may involve detailed reaction feasibility assessments, availability checks for intermediates, or compliance with regulatory and safety standards. Computational tools for such verifications can be resource-intensive to run, less reliable, or even unavailable, while manual verification can be very slow, biased, and inconsistent – it is not uncommon that chemists disagree with each other in retrosynthesis planning [40, 29].