# A Scientific Reasoning Model for Organic Synthesis Procedure Generation

Guoqing Liu[1*], Junren Li[1,2*†], Zihan Zhao[1,3*†], Eray Inanc[1], Krzysztof Maziarz[1], Jose Garrido Torres[1], Victor Garcia Satorras[1], Shoko Ueda[1], Christopher M. Bishop[1], Marwin Segler[1]

[1]Microsoft Research AI for Science; [2]Peking University; [3]Shanghai Jiao Tong University;
*Equal contribution; †During an internship at Microsoft Research AI for Science
Correspondence to: {guoqingliu, marwinsegler}@microsoft.com

## Abstract

Solving computer-aided synthesis planning is essential for enabling fully automated, robot-assisted synthesis workflows and improving the efficiency of drug discovery. A key challenge, however, is bridging the gap between computational route design and practical laboratory execution, particularly the accurate prediction of viable experimental procedures for each synthesis step. In this work, we present QFANG[1], a scientific reasoning language model capable of generating precise, structured experimental procedures directly from reaction equations, with explicit chain-of-thought reasoning. To develop QFANG, we curated a high-quality dataset comprising 905,990 chemical reactions paired with structured action sequences, extracted and processed from patent literature using large language models. We introduce a Chemistry-Guided Reasoning (CGR) framework that produces chain-of-thought data grounded in chemical knowledge at scale. The model subsequently undergoes supervised fine-tuning to elicit complex chemistry reasoning. Finally, we apply Reinforcement Learning from Verifiable Rewards (RLVR) to further enhance procedural accuracy. Experimental results demonstrate that QFANG outperforms advanced general-purpose reasoning models and nearest-neighbor retrieval baselines, measured by traditional NLP similarity metrics and a chemically-aware evaluator using an LLM-as-a-judge. Moreover, QFANG generalizes to certain out-of-domain reaction classes and adapts to variations in laboratory conditions and user-specific constraints. We believe that QFANG's ability to generate high-quality synthesis procedures represents an important step toward bridging the gap between computational synthesis planning and fully automated laboratory synthesis.

---

[1]QFANG is a moniker derived from the Chinese phrase *Qianfang* (meaning "thousands of recipes"). One suggested pronunciation is *Chien-fahng*.

## Introduction

Organic synthesis is the foundational engine of molecular innovation, enabling the creation of a wide range of life-saving pharmaceuticals and other advanced functional molecules[1,2]. While modern algorithms can design millions of novel molecules *in silico*, the practical synthesis of these molecules remains a major bottleneck[2–4]. This process is typically resource-intensive and depends heavily on the tacit knowledge and expert intuition accumulated by chemists over years of practice[5].

To accelerate this process, the field of Computer-Aided Synthesis Planning (CASP) has made significant strides in automating reaction design[6–11]. Advanced algorithms can now predict single-step reactions[12–22], plan complex multi-step retrosynthetic routes[23–30], and even suggest suitable reaction conditions[31–35]. Yet, a critical gap remains between these high-level plans and their practical execution in the laboratory[36–38]. Converting strategic plans into precise, step-by-step experimental procedures still requires substantial human effort to specify essential operational details, such as the reagent addition sequence, reaction durations, temperature gradients, work-up, and purification methods. This challenge is further amplified in the context of robotic systems intended to automate and scale chemical reaction execution[39,40]. Early work to address this gap framed procedure generation as a sequence-to-sequence task, training Transformer- and BART-based models[41,42] from scratch to translate chemical equations into ordered action steps[38,43,44]. While pioneering, these methods were constrained by the representational capacity of early models, often struggling to produce lengthy, coherent procedures or to capture the underlying chemical principles of reactions.

The recent advent of Large Language Models (LLMs) provides a promising avenue for this domain[45–62]. Trained on extensive corpora compris-

## Reaction Input



```
[H][C@]1([C@H](C2=CC=CC=C2)O)O[C@@]3(CC(O[C
@@]3([C@H]1O)[H])=O)[H]>>[H][C@]4(C(C5=CC=CC=
C5)=O)O[C@@]6(CC(O[C@@]6([C@H]4O)[H])=O)[H]
```

## Scientific Challenge

- What's the experimental procedure for this reaction?
- Which **solvent** should I use?
- How to achieve **chemoselectivity**?
- How to properly **isolate** the product from the reaction mixture?
- ......

## QFANG Training



**Base Model**
Qwen3-series (8B/32B)

**Procedure Dataset**
Over 900k reaction-procedure pairs extracted by LLMs

**Chemistry-Guided Reasoning & SFT**
Generate chain-of-thoughts at scale that emulate chemists' reasoning processes and perform supervised fine-tuning

**Reinforcement Learning from Verifiable Reward**
Leverage verifiable rewards to Improve procedure accuracy

## Reasoning Trace <think></think>

*......The reaction involves the oxidation of a benzylic alcohol to a ketone, specifically targeting the hydroxyl group adjacent to the phenyl ring while preserving other functional groups such as esters, ethers, and aliphatic hydroxyls. This selectivity is critical, as indiscriminate oxidation could disrupt the ester or hydroxyl moieties elsewhere in the molecule. Manganese dioxide is chosen as the reagent because it is a well-established mild oxidant for benzylic alcohols, avoiding over-oxidation to carboxylic acids or breaking of carbon-carbon double bonds......*

## Predicted Procedure

```
Make a solution by dissolving
[H][C@]1([C@H](C2=CC=CC=C2)O)O[C@@]3(CC(O[C@@]3([C@H]1O)[H
])=O)[H] (9.06 mmol) in methylene chloride (50 mL) to get
Mixture 1.
Add manganese dioxide (3.95 g, 45.3 mmol) to Mixture 1 to
get Mixture 2.
Wait for 24.00 hours. Stirring.
Filter Mixture 2 using celite to get the filtrate Mixture
3 and the residue Mixture 4.
Chromatograph Mixture 3 to get Mixture 5.
Obtain
[H][C@]4(C(C5=CC=CC=C5)=O)O[C@@]6(CC(O[C@@]6([C@H]4O)[H])=
O)[H] from Mixture 5.
```
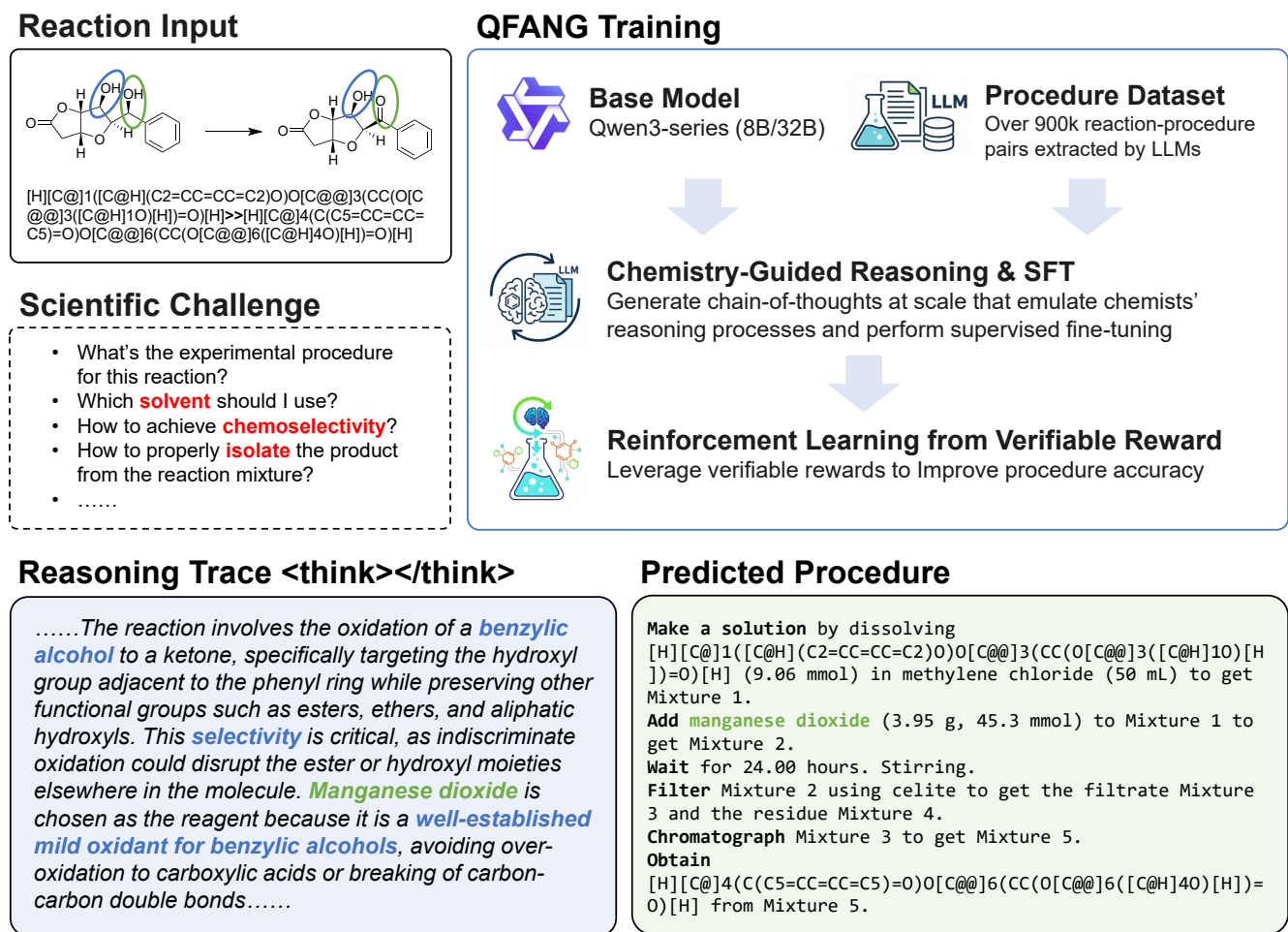
Figure 1: Overview of the inputs, training methodology, and outputs of QFANG. The training process comprises three stages: (i) A large-scale experimental procedure dataset with over 900k reaction-procedure pairs is is assembled, (ii) A chemistry-guided reasoning (CGR) framework is applied to generate chemically grounded chain-of-thoughts at scale, followed by supervised fine-tuning, and (iii) A reinforcement learning stage that leverages verifiable rewards to further improve model predictions. During inference, QFANG takes a reaction equation as input, generates a high-level reasoning trace capturing key decisions, and outputs a detailed, machine-readable experimental procedure.

ing chemical literature, reaction databases, and experimental protocols, LLMs internalize the statistical patterns that underlie established principles, precedent reactions, and methodological variations accumulated over decades of research. Combined with their ability to interpret complex contexts and perform multi-step reasoning[63–65], LLMs are well-suited for tasks like procedure generation. Unlike conventional reaction condition prediction, which focuses on the core reaction stage, full procedure generation needs to specify the work-up and purification operations. These operations require context-dependent chemical reasoning that current condition prediction models cannot provide[66, 67].

Initial efforts to harness LLMs have used few-shot, in-context learning with general-purpose models, guided by curated, action-level datasets[68]. However, in-context learning depends on analogical inference from limited examples, and thus struggles to develop a mechanistic understanding of chemistry.

To build robust and capable systems, models should move beyond analogical imitation and undergo explicit training in chemical reasoning, enabling generalization beyond seen examples[57, 59, 61].

To address this gap, we developed QFANG, a scientific reasoning model post-trained on a large corpus of reaction–procedure pairs collected through an LLM-based automated annotation pipeline. To elicit complex chemical reasoning, we propose Chemistry-Guided Reasoning (CGR), a two-stage framework that first programmatically constructs a *factual skeleton* capturing the core logic of each chemical transformation, and subsequently expands it into an expert-style chain-of-thought using an LLM. The model is supervised fine-tuned on the resulting CGR dataset to acquire chemical reasoning patterns. To further improve model accuracy, we incorporate Reinforcement Learning from Verifiable Rewards (RLVR), which applies rule-based, step-wise reward functions to enforce chemical robustness. A high-level overview of

QFANG is shown in Figure 1. We instantiate the first version of QFANG on a dataset comprising 905,990 reaction-procedure pairs, while noting its straightforward scalability to larger datasets in the future.

To evaluate QFANG's capabilities, we benchmarked it against advanced general-purpose reasoning models and nearest-neighbor retrieval baselines, including GPT-5 (High)[2], using established text-similarity metrics consistent with prior work. Recognizing the limitations of traditional metrics, we further developed a chemically-aware evaluation framework that leverages GPT-5 as an expert judge. On conventional metrics, QFANG achieves a BLEU-4 score of 61.3, surpassing the 54.4 score of a strong retrieval-augmented 3-shot GPT-5 baseline. Under expert-judge evaluation, this margin increases further. Additional analyses show that QFANG generalizes well to out-of-domain reactions, adapts procedural plans to chemist-specified constraints, and even corrects flawed procedures originating from its training data. These results demonstrate that QFANG exhibits a deeper chemical understanding that extends beyond simple pattern matching.

# Large-Scale Procedure Dataset Construction via LLM Annotation

To develop a reasoning model capable of generating viable experimental procedures, it is essential to first define a fine-grained and comprehensive action system[3], and collect large-scale, high-quality, structured procedure datasets.

## Action System Design

For an experimental procedure to be precise, structured, and machine-readable, it should be grounded in a well-defined action system that comprehensively covers all operations involved in chemical experimentation, with detailed specifications for each operation[36,37,69–76]. Rather than directly adopting the earlier action system proposed by Vaucher et al.[37] in 2020, which has been employed in several subsequent works such as ActionIE[70] and OpenExp[71], we employ an enhanced action system that is both more expressive and more comprehensive.

Our system defines 24 distinct actions, extending the set with operations such as `Change pressure`, `Change atmosphere`, `Sample`, `Irradiate`, `Chromatograph`, and `Distill`. Each operation is supported by more detailed arguments and outputs, enabling finer granularity and richer semantics. In ad-

| O3-high Scoring | OpenExp | Ours |
|---|---|---|
| Substance Accuracy ↑ | 7.80 | 9.46 |
| Action Coverage ↑ | 5.72 | 8.62 |
| Order Correctness ↑ | 6.35 | 8.92 |
| Overall Score | 6.14 | 8.76 |

**Table 1:** Quality comparison between experimental synthesis procedures expressed using our action system and those of OpenExp.

dition, the system supports parallel operations across multiple mixtures, whereas OpenExp focuses on procedures for a single target mixture. Detailed definitions of all 24 actions are provided in the Appendix. To assess the validity of our action system, we compared procedures for identical reactions represented using our system and OpenExp's, evaluated with the OpenAI O3-high scoring metric (defined in the Appendix). As shown in Table 1, our action system achieves superior performance in substance accuracy, action coverage, and procedural order.

# Annotation of Structured Experimental Procedures from Text Paragraphs

Constructing a high-quality, large-scale dataset of reaction–procedure pairs is essential for training reasoning models. Yet such data are typically expensive and labor-intensive to obtain[77,78]. Most existing chemical databases provide noisy, unstructured experimental descriptions mined from patents, rather than clean, step-by-step procedures. To overcome this limitation and build a high-quality, large-scale dataset, we leverage general-purpose LLMs. While these models may lack deep chemical expertise, they are highly effective at extracting structured information and producing outputs in customized formats. Specifically, we employ GPT-4o[79] to transform free-text experimental descriptions from Pistachio[80], one of the largest chemical reaction databases, into structured, step-by-step procedures following our defined action schema. The overall annotation pipeline consists of three main steps as illustrated in Figure 2.

**(i) Coreference resolution.** In a single textual description of an experiment, the same substance may be referred to by different names or pronouns. To facilitate the subsequent translation into structured actions, we first instruct GPT-4o to perform coreference resolution on each paragraph. Specifically, we provide GPT-4o with the existing component annotations from the Pistachio database and ask it to replace all chemical entity mentions with their corresponding code names. In addition, GPT-4o is allowed

---

[2]https://openai.com/index/introducing-gpt-5

[3]A standard set of laboratory operations, e.g., Make solution, Add, Change temperature, Quench.
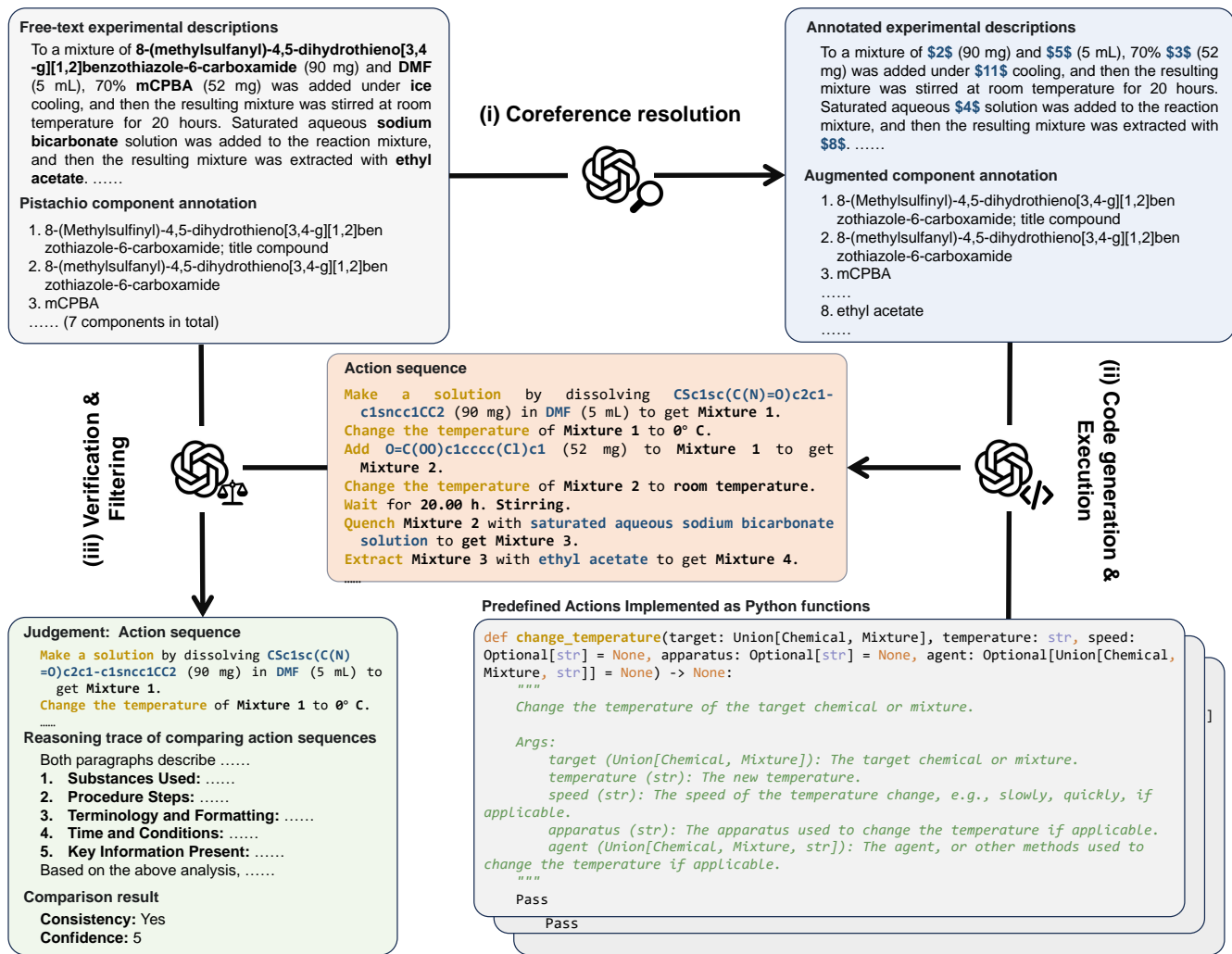
Figure 2: Overview of the procedure annotation pipeline, which consists of three main steps: (i) Coreference resolution: GPT-4o annotates free-text experimental descriptions by replacing chemical entity mentions with their code names, assisted by component annotations from Pistachio. (ii) Code generation: GPT-4o translates the annotated experimental descriptions into structured, step-by-step procedures in Python. (iii) GPT-4o compares the generated procedures with the corresponding descriptions to validate annotation accuracy.

to augment the component list with newly identified molecules or materials when necessary. To ensure the accuracy of the coreference resolution process, we reconstruct the paragraphs and compute the edit distance between the original and reconstructed versions. Entries with edit distances exceeding a predefined threshold are discarded (Figure 2(i)).

**(ii) Code generation and execution.** In this step, we implement each action in our action system as a python function, accompanied by supporting classes such as `Chemical`, `Mixture`, and `TimePeriod`. GPT-4o is then instructed to translate each textual description into the corresponding Python code that invokes these actions. Using code as an intermediate representation provides two key advantages: 1) it allows us to make use of advanced code generation abilities of LLMs, and 2) it enforces structural and type constraints, thereby ensuring that the resulting actions adhere to our pre-defined action specifications. After LLM-based annotation, the generated python code is executed and converted into the final action sequences, while samples with syntax errors or type-check failures are removed (Figure 2(ii)).

**(iii) Verification and filtering.** In this step, we instruct GPT-4o to verify that the generated action sequences accurately express the same experimental procedure as their corresponding source paragraphs. GPT-4o is instructed to produce a reasoning trace before providing a judgment score for consistency (Yes, No, or Uncertain). It then assigns a confidence score for the given judgment ranging from 0 to 5 (where 5 indicates the highest confidence, and 0 indicates the lowest). We retain all entries that pass the consistency verification with a confidence score above 3 and discard the rest (Figure 2(iii)).

Detailed prompts used for all three steps are provided in the Appendix. Through the above action system design and automatic structured exper-

imental procedure annotation pipeline, we obtain a high-quality reaction-procedure dataset from raw Pistachio, consisting of 905,990 chemical reactions paired with their corresponding structured action sequences. This dataset forms the foundation for the subsequent chemistry-guided chain-of-thought data generation, supervised fine-tuning, and reinforcement learning with verifiable rewards.

# Chemistry-Guided Reasoning and Supervised Fine-Tuning

While general-purpose LLMs have showed advanced reasoning in complex tasks such as mathematics and coding, their performance in specialized scientific domains, such as organic chemistry, remains a significant challenge[55, 58, 62]. When tasked with generating a precise, in-depth chain-of-thought (CoT) to derive suitable synthesis procedures, LLMs often generate reasoning that is verbose, factually inaccurate (hallucinated), or driven by superficial textual patterns rather than fundamental chemical principles[58, 81, 82]. In the Appendix, we provide raw reasoning traces from general-purpose LLMs (Qwen3-Max[83] and Phi-4-reasoning[84])[4] prompted with the example reaction shown in Figure 1, illustrating the model's difficulty in parsing the chemical structure.

Training a model to reason like a chemist requires data with two key criteria. First, the CoT reasoning must be grounded in chemical facts, accurately identifying strategic challenges such as chemoselectivity. Second, to support large-scale model training, such high-quality reasoning must be generated at scale, covering hundreds of thousands of chemical reactions like those present in the Pistachio dataset. Manual curation by domain experts is both time-consuming and prone to stylistic inconsistency, making it impractical for large datasets.

To address this, we developed **Chemistry-Guided Reasoning (CGR)**, a hybrid two-stage framework for generating large-scale, high-quality CoT datasets. In the first stage, a deterministic, cheminformatics-based programmatic analysis is applied to chemical reactions to systematically identify key bond changes, functional group conversions, and strategic considerations. This ensures that the underlying chemical logic is both factually accurate and structurally consistent across the entire dataset. In the second stage, these curated chemical facts serve as structured guidance for an LLM to generate coherent, expert-style reasoning narratives. By combining the precision of programmatic analysis with the expressive power of LLMs, CGR enables scalable

construction of chemically rigorous and stylistically consistent CoT datasets for training QFANG.

**Programmatic Generation of Factual Skeletons.** The first stage of CGR programmatically generates a structured, cheminformatics-based reasoning trace for each reaction–procedure pair. It starts with a detailed analysis of the chemical transformation. Using an atom-mapping algorithm (e.g., LocalMapper[85]), we compare the atomic connectivity and bond types in the reactants and products to identify which atoms and bonds are directly involved in the reaction. This structural analysis gives a clear picture of the transformation at the molecular level.

Building on this, our algorithm iterates through a predefined library of 243 functional groups, cataloging their presence in both reactants and products. It categorizes each group as being transformed, newly formed, or unchanged. This automated annotation allows the system to programmatically flag key strategic challenges that a chemist would need to address. For instance, in a selective oxidation of a benzylic alcohol in the presence of an aliphatic one, the system generated a key insight: *"The functional group 'alcohol' is converted to the functional group 'aldehyde', while the other functional group 'alcohol' remains unchanged. It should be noted that the 'alcohol' functional groups are selectively transformed in this reaction. The procedure should take care of the selectivity issue."*

Concurrently, the associated experimental procedure is deconstructed. The algorithm programmatically separates the *reaction* phase from the *workup* phase by identifying keywords such as `Quench`, `Extract`, or `Filter`. Within each phase, it identifies all chemical entities and assigns their roles such as *reactant*, *catalyst*, *reagent*, or *solvent* following the raw annotations[77]. The script also captures critical procedural details, including the order of addition of these components and environmental conditions like the use of a protective nitrogen atmosphere. High-level context, such as the formal reaction name (e.g., "Horner-Wadsworth-Emmons reaction"), is also appended when identifiable. This stage produces a structured collection of verified chemical facts.

**LLM-based Enhancement for Expert-like Narratives.** In the second stage, the factual skeleton generated in the first stage serves as grounding context for an LLM (e.g., Qwen3-235B-Thinking-2507[86]). The input prompt provides the chemical reaction query, the ground-truth experimental procedure, and the programmatically generated factual skeleton. The model is explicitly instructed to synthesize these discrete facts into a coherent, explana-

---

[4]We chose these open-weight models because they allow access to their raw reasoning process.

tory narrative, adopting the perspective of an expert chemist designing the experiment. The prompt directs the model to emphasize the causal links between the identified challenges and the procedural decisions made.

Revisiting the selective oxidation example, where the proper oxidant here is manganese dioxide($MnO_2$), the LLM enhances the rule-based skeleton to produce a final, expert-like reasoning process: *"The reaction involves the oxidation of a benzylic alcohol to a ketone, specifically targeting the hydroxyl group adjacent to the phenyl ring while preserving other functional groups such as esters, ethers, and aliphatic hydroxyls. This selectivity is critical, as indiscriminate oxidation could disrupt the ester or hydroxyl moieties elsewhere in the molecule. Manganese dioxide is chosen as the reagent because it is a well-established mild oxidant for benzylic alcohols, avoiding over-oxidation to carboxylic acids or breaking of carbon-carbon double bonds."*

This two-stage methodology ensures that the final CoT is not only factually accurate but also mirrors the causal, step-by-step logic of an expert chemist.

**Supervised Fine-Tuning.** Following the CGR procedure, all the 905,990 reactions with structured action sequences were paired with their corresponding chemical reasoning traces. Using this dataset, we performed a time-based split, reserving the most recent 10% of the entries as a held-out test set. Specifically, the training set contains reactions published in source patents from 1971 to July 2023, while the held-out test set covers those published from July 2023 to June 2024. This time-based split was chosen over a random split because it more accurately reflects real-world deployment, where a model trained on historical data is expected to predict future, previously unseen chemical procedures.

The remaining 90% of the entries constituted the initial training corpus. A substantial portion of this first version dataset is derived from patents which, despite their scale, are often noisy, exhibiting issues such as stoichiometric inconsistencies, incomplete workup descriptions, and procedures optimized for intellectual property protection rather than experimental reproducibility. To curate a high-fidelity training dataset, we implemented an LLM-based filtering protocol. Specifically, we employed the Qwen3-235B-Thinking-2507 model as an expert chemical evaluator to score each annotated procedure (on a scale of 0–10) across four critical axes: (i) *Reactant Completeness*, (ii) *Workup and Purification Completeness*, (iii) *Condition Completeness*, and (iv) *Reaction Safety*. Only procedures achieving an average score of 5.0 or higher were prioritized, with the additional strict constraint

that no single axis score could fall below 3.0. The detailed prompt used for scoring is provided in the Appendix.

With this curated dataset, we initiated model training via Supervised Fine-Tuning (SFT). We selected the Qwen-3 family[86] as base models, fine-tuning both the 8-billion and 32-billion parameter versions to examine the effect of model scale on performance. During the SFT stage, the models were trained to produce coherent reasoning chains followed by structured experimental procedures. These resulting SFT models, denoted as QFANG-8B (SFT) and QFANG-32B (SFT), not only demonstrated strong initial performance but also served as the starting policy for the subsequent reinforcement learning stage.

# Reinforcement Learning with Verifiable Rewards

After SFT, we leverage Reinforcement Learning from Verifiable Rewards (RLVR) to further enhance the prediction accuracy of QFANG. However, predicting experimental procedures poses unique challenges for reward design. In particular, exact matches between predicted and ground-truth procedures are difficult to achieve, and conventional text-similarity metrics like BLEU fail to capture chemical plausibility. As a result, commonly used outcome-based reward functions in RLVR are not very suitable for this task.

To address this, we design a verifiable, step-wise reward function to guide the RL training. For simplicity of verification, we assume that the predicted action sequence must strictly follow the ground-truth sequence. Thus, synonymous but equivalent actions are treated as incorrect. This assumption enables step-by-step evaluation and the assignment of dense rewards. Specifically, the accuracy reward for each action at step $t$ consists of three parts.

**Format reward** $R_{\text{format}}^t$. It checks whether the predicted action follows the correct format of any defined operation. If the action fails the format check, a negative reward of –1 is assigned, and no further rewards are computed for that step. If the format is correct, the reward for this component is zero.

**Type reward** $R_{\text{type}}^t$. It checks whether the type of the current action matches the ground truth action type. If the types do not match, the reward is set to zero and subsequent rewards for this step are not computed. If the types match, a reward of 1 is assigned for this component.

**Necessary/Optional parameter reward** $R_{\text{nec}}^t/$ $R_{\text{opt}}^t$. We categorize the parameters for all

operations into two groups: *necessary* and *optional*. Necessary parameters define the core functionality of an operation and allow minimal variation, whereas optional parameters are less critical, and differences may not affect the operation's functionality. We evaluate the necessary and optional groups independently. For each group, we compute a score in the range [0, 1] as the average matching quality of the parameters within that group.

Overall, the accuracy reward for the $t$-th action $R_{\text{acc}}^t$, is computed as:

$$R_{\text{acc}}^t = R_{\text{format}}^t + R_{\text{type}}^t + R_{\text{necc}}^t + R_{\text{opt}}^t.$$

In addition to this accuracy reward, we introduce two auxiliary rewards to further stabilize the RL training process and mitigate potential reward hacking during the RLVR phase.

The first auxiliary reward is designed to penalize actions that exceed the ground-truth sequence. In our current formulation, such over-predicted actions receive a fixed negative reward. However, setting this penalty too small or too large can bias the model toward producing overly long or overly short action sequences. To address this, we propose an adaptive penalty scheme for these exceeding actions. Specifically, we aggregate the rewards from all non-exceeding actions at the same time step across the other samples in the batch. We then assign penalties to the exceeding actions such that the average reward for predictions at the current step is zero across the entire batch. Formally, the exceeding reward for sample $i$ at step $t$ is calculated as:

$$R_{i,\text{exc}}^t = \begin{cases} -\frac{1}{N_{\text{exc}}^t} \sum_{j \in \mathcal{B}^t} R_{j,\text{acc}}^t, & i \notin \mathcal{B}^t, \\ 0, & i \in \mathcal{B}^t, \end{cases}$$

where $N_{exc}^t$ is the number of samples in the current batch where the model predicts an exceeding action at step $t$, while $\mathcal{B}^t$ denotes the set of the samples in the current batch whose ground truth action sequence has at least $t$ steps.

The second auxiliary reward is the action type distribution modifier. Due to the varying difficulty in achieving rewards across different action types, the model may develop an undesirable bias towards predicting types that are easier to reward. To address this issue, we introduce a reward modifier based on the distributional differences of the action types. Specifically, we calculate the distribution of action types in the ground-truth annotations and in the model predicted action sequences within the current batch, respectively. The distribution modifier for the

action at step $t$ is then calculated as:

$$R_{\text{dist}}^t = \begin{cases} \frac{p_{\text{gt}}^t - p_{\text{pred}}^t}{\max\left(p_{\text{gt}}^t, p_{\text{pred}}^t\right)}, & \frac{p_{\text{gt}}^t - p_{\text{pred}}^t}{\max\left(p_{\text{gt}}^t, p_{\text{pred}}^t\right)} > \theta, \\ 0, & \text{otherwise}, \end{cases}$$

where $p_{\text{gt}}^t$ denotes the proportion of ground-truth actions in the current batch that share the same type as the predicted action at step $t$. Similarly, $p_{\text{pred}}^t$ denotes the proportion of predicted actions in the current batch that belong to this action type, and $\theta$ is a predefined threshold.

Finally, the total reward for the action at step $t$ used during training is given by:

$$R^t = R_{\text{acc}}^t + R_{\text{exc}}^t + R_{\text{dist}}^t.$$

A detailed algorithm describing the computation process of the final reward is presented in the Appendix. Future versions will consider incorporating additional factors, such as enforcing consistency between the mass and the number of moles used.

**PPO Training.** Building on the designed reward function, we fine-tune the QFANG model using the Proximal Policy Optimization (PPO) algorithm[87]. PPO is an empirically stable choice for language model alignment and complex reasoning tasks[63,88,89], particularly well-suited for our dense reward setting. The PPO objective function is defined as:

$$\mathcal{L}_{\text{PPO}}(\theta, \psi) = \mathcal{L}_{\text{clip}}(\theta) - w_1 \text{KL}(\theta) - w_2 \mathcal{L}(\psi),$$

where $\theta$ represents the language model parameters, and $\psi$ represents the critic model parameters. The clipped objective, $\mathcal{L}_{\text{clip}}(\theta)$, maximizes the expected reward while restricting large policy updates. The KL penalty $\text{KL}(\theta)$ limits divergence from the reference policy (e.g., the SFT model). Lastly, the critic loss, $\mathcal{L}(\psi)$, minimizes value estimation errors, improving reward prediction and overall training stability.

**GRPO Training.** Alternatively, the dense reward structure can be reformulated as an outcome-based setting, allowing the use of Group Relative Policy Optimization (GRPO)[90] to fine-tune QFANG. The GRPO objective function is defined as:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathcal{L}_{\text{clip}}(\theta) - w_1 \text{KL}(\theta).$$

A key innovation of GRPO over PPO is the elimination of a separate critic model. Instead, GRPO estimates the advantage by comparing the reward of each response to the average reward obtained from a group of responses generated by the same policy. Given a question $x$ from the dataset, GRPO samples $G$ responses $y_1, \ldots, y_G$, and assigns rewards $r_1, \ldots, r_G$ to

| Model | BLEU-2 | BLEU-4 | LEV avg | LEV 90% | LEV 75% | LEV 50% | Rouge-1 | Rouge-2 | Rouge-L | METEOR | Seq-O |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Nearest Neighbor | 54.5 | 39.1 | 46.6 | **1.0** | 2.4 | 30.2 | 60.0 | 38.4 | 47.3 | 50.0 | 65.1 |
| 1-shot GPT-4o | 39.6 | 29.9 | 31.3 | 0.0 | 0.0 | 1.4 | 50.4 | 26.2 | 36.5 | 55.1 | 6.1 |
| 3-shot GPT-4o | 41.5 | 31.8 | 32.7 | 0.0 | 0.0 | 2.4 | 53.1 | 30.1 | 39.4 | 57.4 | 8.4 |
| 1-shot o4-mini (high) | 55.2 | 43.8 | 42.7 | 0.0 | 0.1 | 17.6 | 60.4 | 36.9 | 46.6 | 58.9 | 26.6 |
| 3-shot o4-mini (high) | 56.8 | 45.5 | 44.7 | 0.0 | 0.3 | 24.0 | 62.9 | 39.9 | 49.2 | 61.5 | 26.2 |
| 1-shot GPT-5 (high) | 67.1 | 56.4 | 51.5 | 0.8 | 2.9 | 54.1 | 67.6 | 49.4 | 55.3 | 64.7 | 59.9 |
| 3-shot GPT-5 (high) | 65.0 | 54.4 | 51.7 | 0.6 | 2.6 | 55.4 | 68.5 | 50.1 | 55.9 | 66.9 | 61.1 |
| QFANG-8B (SFT) | 69.9 | 59.4 | 56.8 | 0.1 | 2.8 | 77.9 | 71.2 | 52.9 | 59.8 | 68.6 | 69.4 |
| QFANG-32B (SFT) | 70.2 | 59.7 | 57.2 | 0.1 | 3.2 | 78.8 | 71.5 | 53.3 | 60.2 | 68.8 | 70.1 |
| QFANG-8B (RL) | **72.0** | **61.3** | **57.4** | 0.3 | **4.3** | 78.8 | **72.1** | **54.5** | **61.1** | **69.7** | **70.9** |

**Table 2:** Performance comparison of QFANG against baselines on traditional NLP metrics. Metrics are in the interval [0, 100].



Figure 3: LLM-as-a-judge evaluator and performance comparison. (a) Overview of the evaluation process, where the GPT-5 judge receives the reference procedure, the generated procedure, and a detailed scoring rubric. (b) Bar chart comparing QFANG's final scores with baseline models; horizontal dashed lines mark the Oracle baseline (upper) and the lowest-performing negative baseline (lower).

each. For the $i$-th response in the group, the corresponding group-relative advantage is computed as:

$$\hat{A}_i = \frac{r_i - \text{mean}(r_1, \ldots, r_G)}{\text{std}(r_1, \ldots, r_G)}.$$

Additional details on the loss components are provided in the Appendix.

# Results

**Quantitative Evaluation: Benchmarking against Baselines.** To assess the performance of QFANG, we conducted a comprehensive quantitative evaluation against several strong baselines using standard text-similarity metrics, along with a new LLM-as-a-judge assessment.

To contextualize the performance of QFANG, we established a set of competitive baselines, comprising both retrieval-based and advanced generative approaches. The simplest baseline is a non-generative Nearest Neighbor (NN) method. For a given target reaction, this approach identifies the single most similar reaction from the training set, based on Tanimoto similarity of their reaction fingerprints (DRFP[91]) and directly outputs its associated procedure as the prediction. We further tested against a suite of state-of-the-art LLMs, including GPT-4o, o4-mini (high) and GPT-5 (high), leveraging a retrieval-augmented in-context learning strategy. For each target reaction, we retrieved the top-$k$ most similar reactions and their ground-truth procedures from the training set using the same DRFP-based similarity metric. These retrieved pairs were then integrated into the prompt as in-context examples to guide the model's generation. We evaluated this strategy in both 1-shot ($k$=1) and 3-shot ($k$=3) settings. During inference, QFANG operates without in-context learning.

**Superior Performance on Traditional NLP Metrics.** Similar to previous studies, we evaluated the lexical similarity between the generated and reference procedures using a suite of established metrics[38,68], with results presented in Table 2. These included BLEU, a metric based on n-gram precision originally developed for machine translation[92]; ROUGE, which measures n-gram recall[93]; and ME-

TEOR, which incorporates synonyms and stemming for more robust alignment[94]. We also report metrics based on normalized Levenshtein (LEV) similarity, which quantifies character-level edit distance[95]; specifically, LEV X% denotes the fraction of predictions achieving a normalized similarity score of X% or greater when compared to the reference[38]. Finally, we include Seq-O, a metric designed to measure the similarity of the core action verb sequences, providing insight into the model's ability to capture the procedural workflow[96].

Across all evaluation metrics, QFANG consistently outperforms all baseline models. The RL trained variant, QFANG-8B (RL), achieves a BLEU-4 score of 61.3 and a ROUGE-L score of 61.1. This substantially surpasses the strongest generative baseline, 3-shot GPT-5 (high), which scored 54.4 and 55.9, respectively. This trend holds for our SFT models as well; the QFANG-32B (SFT) variant reached a METEOR score of 68.8 and a Seq-O score of 70.1, indicating a strong capability to generate text that is both syntactically correct and lexically aligned with the ground-truth data. Notably, the performance of retrieval-augmented LLMs is highly dependent on the number of in-context examples, increasing the number of examples does not necessarily lead to better scores. We also find that the simple Nearest Neighbor baseline performs well on some metrics like BLEU-2 and Seq-O, suggesting that many reactions in the test set may have close analogues in the training data.

To further investigate whether the superior average performance of QFANG stems from true generalization or merely an improved ability to recall similar training examples, we conducted a more rigorous analysis. We stratified the test set based on its procedural and chemical similarity to the training data. The results, detailed in Appendix, demonstrate that QFANG maintains its performance advantage even on samples that are highly dissimilar to the training set, whereas the baseline models exhibit a much sharper decline. This provides strong evidence that QFANG has learned to reason from underlying chemical principles, validating the effectiveness of our CGR training methodology.

Despite the superior performance of QFANG on these metrics, we posit that such surface-level comparisons are insufficient for rigorously evaluating chemical procedure generation. Lexical metrics are inherently insensitive to chemical logic and safety. For example, a simple string comparison would penalize the substitution of a hazardous reagent like sodium hydride (NaH) with a milder base like sodium hydroxide (NaOH) with only a minor score reduction, yet this represents a fundamental and potentially dangerous misunderstanding of the required reaction conditions. Similarly, a chemically consistent substitution, such as writing $NEt_3$ instead of triethylamine, would be incorrectly penalized as a mismatch. While the strong performance of QFANG on these metrics is a prerequisite, a more meaningful assessment of chemical soundness and procedural viability is essential for a comprehensive evaluation of QFANG's capabilities. This motivates our development of a chemically-aware, LLM-based evaluation framework.

**LLM-as-a-Judge: A Chemically-Aware Evaluator.** To address the semantic gap of traditional NLP metrics, we designed and implemented a more rigorous evaluation framework that leverages an expert-level LLM, GPT-5 (high), as a judge. This framework moves beyond simple lexical comparison to assess generated procedures on their chemical viability, procedural completeness, and strict adherence to a machine-readable format, reflecting the practical requirements for synthesis planning and eventual laboratory automation.

The evaluation of each generated procedure is performed by providing the GPT-5 judge with a composite prompt that includes the ground-truth reference procedure (as the gold standard), the model-generated procedure, and a detailed scoring rubric. The rubric consists of four categories: *Reaction Score* (40 points), assessing the core transformation and stoichiometry; *Workup and Purification Score* (30 points), evaluating the separation process; *Conditions Score (20 points)*, examining solvents and reagents; and *Safety and Modern Practice Score* (10 points), penalizing usage of outdated and toxic components. The prompt also penalized deviations from a required machine-readable syntax.

To validate our LLM-as-a-judge framework and contextualize the scores, we established several control baselines designed to probe the judge's sensitivity to specific types of procedural errors. We created three *negative baselines* by systematically corrupting the ground-truth procedures: (1) *Reagent*, where a key reagent was replaced with a chemically nonsensical alternative; (2) *Swap Actions*, where the order of two critical steps was inverted; and (3) *Both*, which combined both errors. Conversely, we established an *Oracle baseline* by making only chemically benign modifications to the ground-truth procedures, such as replacing reagents with known synonyms (e.g., substituting methanol with MeOH). As shown in Figure 3, these baselines provide a clear scale for interpretation. The *Oracle baseline* achieves a near-perfect score of 90.5, confirming the judge's ability to tolerate reasonable chemical variations. In contrast, the *negative baselines* score significantly lower, dropping to 39.1, 39.7, and 26.8, respectively, demonstrating the
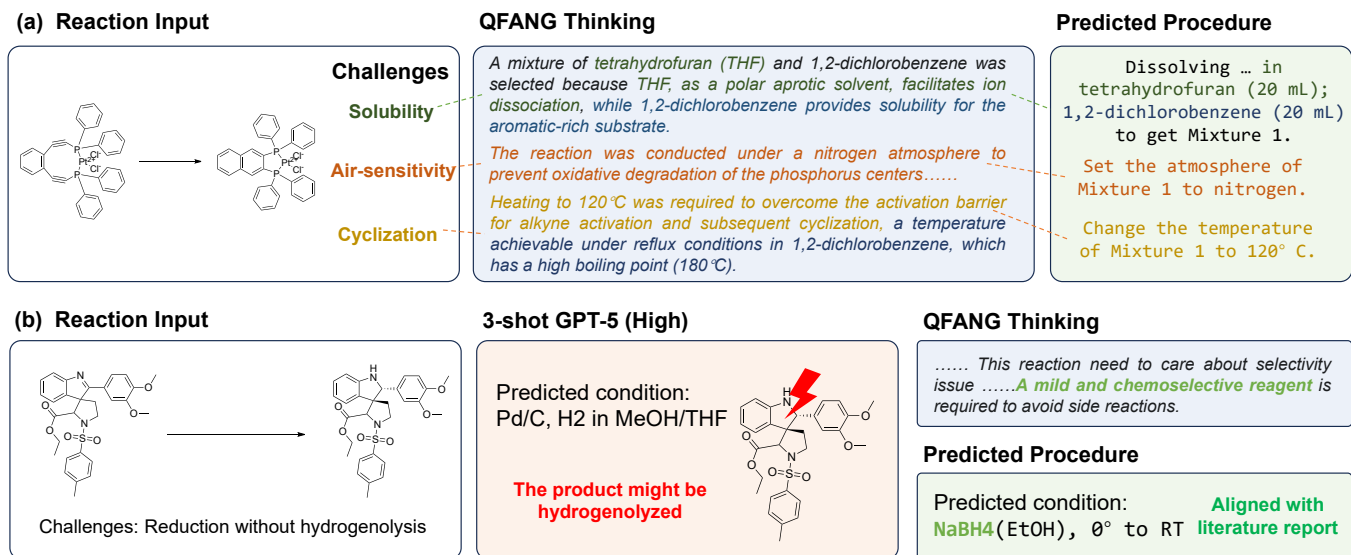
Figure 4: Demonstration of QFANG's ability to generalize across diverse chemical contexts. (a) For an out-of-distribution organometallic reaction, QFANG infers a viable synthetic procedure from first principles. It correctly identifies structural challenges like poor solubility and air-sensitivity, and proposes expert-level solutions (e.g., a binary solvent system and an inert atmosphere) in its generated procedure. (b) Achieving chemoselectivity on natural product intermediates. When tasked with a sensitive imine reduction, QFANG exhibits superior chemical judgment. In contrast to a 3-shot GPT-5 baseline, which suggests harsh and potentially destructive conditions (Pd/C, H₂), QFANG selects a mild and chemoselective reagent (NaBH₄), thereby preventing side reactions and aligning with the literature-reported gold standard.

framework's capability to effectively penalize chemically and structurally flawed procedures.

The results of this expert-level LLM evaluation, presented in Figure 3, further solidify the superior capabilities of QFANG. Our best-performing model, QFANG-8B (RL), achieves a final score of 78.2, surpassing the strongest baseline (3-shot GPT-5-high at 67.8) by a substantial margin and closely approaching the Oracle score. This high score indicates that the procedures generated by QFANG are not only chemically sound and logical but also adhere precisely to the stringent formatting required for robust, automated laboratory systems. The performance of the SFT models, QFANG-8B (SFT) and QFANG-32B (SFT) at 72.3 and 74.2, respectively, also shows a clear advantage over the retrieval-augmented LLMs. This outcome validates our training methodology, demonstrating that our model has learned the underlying principles of chemical procedure design rather than merely mimicking surface-level text patterns.

## Case Studies

**Generalizability: Navigating Out-of-Domain Chemical Challenges.** We evaluated the model's generalization ability through two distinct and highly challenging cases: one requiring reasoning in a completely novel chemical domain, and another demanding nuanced selectivity within the familiar context of natural product synthesis.

Our first challenge tested the model's ability to reason from first principles on a complex, out-of-distribution reaction from the domain of organophosphorus and organometallic chemistry, a field far removed from the drug synthesis patents that constitute its training data. The task was to devise a procedure for an intramolecular cycloaromatization within a platinum-diphosphine heterocyclic framework (Figure 4(a))[97]. Success in this task hinges not on recalling a known reaction template, but on a fundamental analysis of the substrate's unique structural features: its poor solubility, the presence of air-sensitive phosphorus centers, and the high activation energy typical of such bond reorganization reactions.

Remarkably, QFANG generated a highly plausible and expert-level procedure. It correctly inferred the need for a binary solvent system, pairing a polar aprotic solvent (THF) to solvate the ionic phosphonium portions with a high-boiling-point, nonpolar aromatic solvent (1,2-dichlorobenzene) to dissolve the hydrophobic backbone and provide the necessary thermal energy. Furthermore, it identified the air-sensitive nature of the phosphorus centers and correctly prescribed an inert nitrogen atmosphere. The proposed conditions of moderate-high temperature (120°C) and prolonged reaction time (overnight) are also fully consistent with the high activation barrier expected for such a Bergman cycloaromatization. This case demonstrates a profound capability to de-
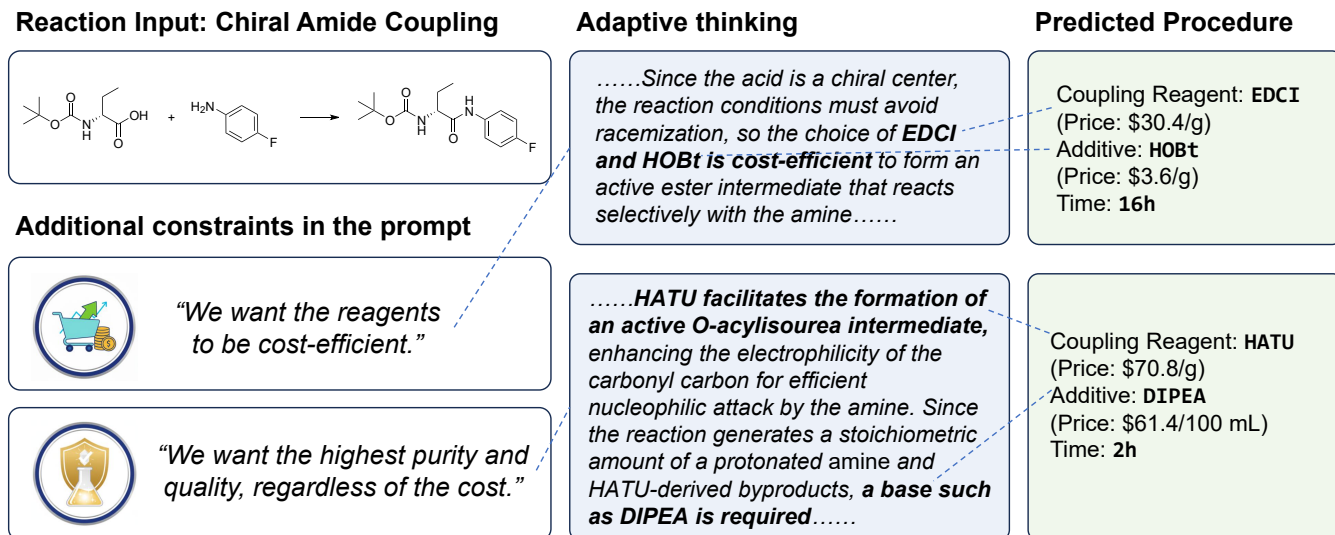
**Reaction Input: Chiral Amide Coupling**

**Additional constraints in the prompt**

*"We want the reagents to be cost-efficient."*

*"We want the highest purity and quality, regardless of the cost."*

**Adaptive thinking**

*……Since the acid is a chiral center, the reaction conditions must avoid racemization, so the choice of EDCI and HOBt is cost-efficient to form an active ester intermediate that reacts selectively with the amine……*

*……HATU facilitates the formation of an active O-acylisourea intermediate, enhancing the electrophilicity of the carbonyl carbon for efficient nucleophilic attack by the amine. Since the reaction generates a stoichiometric amount of a protonated amine and HATU-derived byproducts, a base such as DIPEA is required……*

**Predicted Procedure**

Coupling Reagent: **EDCI**
(Price: $30.4/g)
Additive: **HOBt**
(Price: $3.6/g)
Time: **16h**

Coupling Reagent: **HATU**
(Price: $70.8/g)
Additive: **DIPEA**
(Price: $61.4/100 mL)
Time: **2h**

Figure 5: Demonstration of QFANG's ability to adapt based on user-specific constraints. For the same chiral amide coupling reaction, the model proposes two distinct, chemically valid procedures, each tailored to different and potentially competing user objectives: (i) cost efficiency, and (ii) maximum purity and quality.

construct a novel molecular structure and devise a viable experimental plan based on core chemical principles of solubility, stability, and reactivity.

Having established the model's ability to navigate a new chemical domain, we next probed its capacity to handle fine-grained selectivity within a more familiar, yet notoriously complex area: natural product synthesis. The second challenge was a critical reduction step from R.B. Woodward's landmark total synthesis of strychnine[98]. The task involves the reduction of an imine within the dense, polyfunctionalized core of the natural product (Figure 4(b)). The primary difficulty lies in achieving chemoselectivity, as a non-selective or overly harsh reducing agent could easily lead to undesired side reactions, such as the hydrogenolysis of the newly formed C-N bond.

As shown in Figure 4, QFANG again successfully navigated this challenge by proposing the use of sodium borohydride, a mild and selective reducing agent. This prediction aligns perfectly with the established literature procedure, achieving the desired clean conversion. In stark contrast, a powerful baseline like GPT-5-high, when prompted, proposed catalytic hydrogenation with Palladium on carbon. While effective for simple imines, these powerful conditions are ill-suited for this delicate substrate and would likely cause product decomposition. This pair of case studies effectively illustrates the breadth and depth of QFANG's chemical understanding, showcasing its ability to reason from context-specific rules in complex domains, which is essential for reliable real-world synthesis planning. Beyond chemical novelty, we also evaluated operational generalization. In the Appendix, we demonstrate QFANG's ability to transition from discovery-scale to process-scale chem-
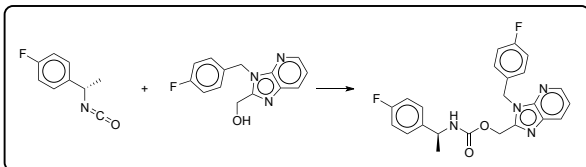
istry, successfully designing a chromatography-free, 50-kg manufacturing procedure under industrial constraints.

**Goal-Oriented Planning with User Constraints.** A key feature of a practical synthesis planning system is its ability to adapt and re-optimize plans based on user-specified high-level constraints. We illustrate this through a case study of chiral amide coupling, a common transformation where maintaining stereochemical integrity is crucial. This example highlights that there is no single "best" procedure; rather, the optimal conditions are determined by a nuanced trade-off between cost, time, and the desired stereochemical purity[99].

As illustrated in Figure 5, we prompted QFANG with the same chemical transformation but under two distinct, competing constraints. Initially, when prompted with the constraint *"We want the reagents to be cost-efficient."*, QFANG generated a procedure centered on the use of EDCI and HOBt. The reasoning trace identified the substrate's chirality and the inherent risk of racemization associated with simple carbodiimide activators. Consequently, it correctly included HOBt as a necessary additive to suppress this side reaction, proposing a 16 hour procedure that reflects a standard, cost-effective method suitable for a basic laboratory setup.

In a subsequent interaction, the model was given a new, quality-focused directive: *"We want the highest purity and quality, regardless of the cost."*. In response, QFANG dynamically altered its plan, replacing the previous system with the more expensive but superior coupling reagent, HATU, and the base DIPEA. The reasoning trace shifted to priori-

**Reaction Input: Carbamate Formation**

**Patent Procedure (excerpted)**

```
Make a solution by dissolving
OCc1nc2cccnc2n1Cc1ccc(F)cc1 (70 mg, 0.27 mmol in THF (1
mL) in round bottom flask fitted with magnetic stirrer
and reflux condenser to get Mixture 1.
Add C[C@H](N=C=O)c1ccc(F)cc1 (225 mg, 1.36 mmol) to
Mixture 1 to get Mixture 2.
Change the temperature of Mixture 2 to 0° C with
Mixture 1.
Wait for 1.50 h. Stirring.
Change the temperature of Mixture 2 to 40°C.
Wait until overnight.
Add EtQAç (15 mL) to Mixture 2 to get Mixture 3.
Wash Mixture 3 with water (10 mL) to get Mixture 4.
……
Obtain product from Mixture 9 with a percentage yield
of 4.00%(4 mg).
```

**QFANG Thinking**

*......Because the reaction requires an excess of isocyanate to ensure complete reaction, **1.2 equivalents** of the isocyanate should be added. To promote the reaction, DMAP can also be added to the reaction system as a catalyst. **DMAP** is a highly efficient nucleophilic catalyst that significantly increases the rate of reaction between isocyanates and alcohols, allowing the process to occur under mild conditions.*

**Predicted Procedure**

```
Make a solution by dissolving
OCc1nc2cccnc2n1Cc1ccc(F)cc1 (0.100 g, 0.39
mmol); DMAP (0.005 g, 0.04 mmol) in DCM (2 mL)
to get Mixture 1.
Add C[C@H](N=C=0)c1ccc(F)cc1 (0.070 g, 0.43
mmol) to Mixture 1 to get Mixture 2.
Wait for 1.00 h. Stirring.
……
```

Figure 6: Demonstration of QFANG's ability to correct flawed procedures. When presented with a stoichiometrically imbalanced reaction from the source dataset, the model produces a chemically valid and optimized alternative, incorporating the appropriate catalyst to ensure reaction feasibility.

tize efficiency and stereochemical fidelity, explaining that modern uronium salt-based reagents like HATU are renowned for their rapid kinetics and their exceptional ability to prevent racemization. This was reflected in the final procedure, which reduced the reaction time from 16 hours to just 2 hours while ensuring the highest possible enantiomeric purity—a standard expected in an advanced R&D setting where speed and quality are paramount.

This case study illustrates that the understanding of QFANG goes beyond surface-level reagent selection. It grasps the deep, practical trade-offs in synthetic chemistry and can select the appropriate methodology based on high-level user goals, demonstrating a crucial step towards creating a collaborative and intelligent tool for synthesis planning.

**Surpassing and Correcting Flawed Training Data.** An interesting finding from our study is that QFANG can identify and correct chemically flawed procedures present in its training data. This capability indicates that the model has internalized fundamental chemical principles, allowing it to generate procedures that are superior to some of the examples it was trained on. Its function as a corrective filter is critical for building reliable systems from vast but imperfect data sources, such as chemical patents. We illustrate this with a representative case from the training set, where the original patented procedure is stoichiometrically questionable, leading to a very low reported yield.

The case involves an acylation reaction where the original procedure calls for a four-fold excess of a highly reactive isocyanate (Figure 6). From a chemical standpoint, such a large excess is not only wasteful but also detrimental to the process. The unreacted electrophile will inevitably undergo hydrolysis during the aqueous workup, leading to significant byproduct formation (e.g., urea) that complicates purification and ultimately contributes to the low reported yield of just 4%[100].

Instead of merely replicating this flawed method, QFANG proposes a refined and more efficient procedure. Its CoT analysis correctly identifies the core transformation but implicitly rejects the erroneous stoichiometry. The model generates a new protocol that reduces the isocyanate to a more standard 1.2 equivalents, a quantity sufficient to drive the reaction to completion without creating an excessive purification burden. Furthermore, it introduces a catalytic amount of DMAP (4-Dimethylaminopyridine) to ensure an efficient reaction, a common best practice for acylations that was absent in the original text. Finally, it implements a more robust, multi-step basic and neutral wash protocol designed to effectively remove both acidic and neutral impurities before the final purification.

This case study shows that QFANG does not blindly reproduce the flawed methods present in its training corpus. Instead, it leverages its learned understanding of stoichiometry and reaction optimization to generate a chemically sound and superior alternative. This ability to critically assess and rectify suboptimal procedures underscores its potential as a

tool not only for novel synthesis planning but also for the crucial validation and refinement of existing experimental protocols.

## Conclusion

In this work, we addressed the critical gap between synthesis planning and machine-readable laboratory actions by introducing QFANG, a scientific reasoning model capable of generating high-fidelity chemical procedures. Our approach was built on three key pillars: (i) construction of a large-scale, structured procedure dataset through a novel LLM-based automatic annotation pipeline; (ii) development of a chemistry-guided reasoning framework to elicit chemical-principle-based reasoning in the model; and (iii) integration of a reinforcement learning with verifiable rewards stage to further enhance predictive accuracy.

Comprehensive evaluations demonstrate that QFANG achieves competitive performance on established NLP related metrics, surpassing strong baselines such as retrieval-augmented GPT-5, while maintaining high chemical validity under the LLM-as-a-judge assessment framework. Qualitative analyses further reveal the emergence of chemical reasoning capabilities, as reflected in the model's ability to generalize to out-of-domain reactions, adapt experimental plans to high-level user constraints, and even identify and amend flawed procedures within its own training data. By generating robust, well-reasoned, and machine-readable experimental protocols, QFANG represents a promising step toward bridging the loop between in-silico design and laboratory execution, contributing to the progress of next-generation autonomous platforms for scientific discovery.

## Acknowledgements

## References

[1] Li, J. *et al.* Synthesis of many different types of organic small molecules using one automated process. *Science* **347**, 1221–1226 (2015).

[2] Blakemore, D. C. *et al.* Organic synthesis provides opportunities to transform drug discovery. *Nature chemistry* **10**, 383–394 (2018).

[3] Campos, K. R. *et al.* The importance of synthetic chemistry in the pharmaceutical industry. *Science* **363**, eaat0805 (2019).

[4] Stanley, M. & Segler, M. Fake it until you make it? generative de novo design and virtual screening of synthesizable molecules. *Current Opinion in Structural Biology* **82**, 102658 (2023).

[5] Ley, S. V., Fitzpatrick, D. E., Ingham, R. J. & Myers, R. M. Organic synthesis: march of the machines. *Angewandte Chemie International Edition* **54**, 3449–3464 (2015).

[6] Vleduts, G. Concerning one system of classification and codification of organic reactions. *Information Storage and Retrieval* **1**, 117–146 (1963).

[7] Corey, E. J. & Wipke, W. T. Computer-assisted design of complex organic syntheses: Pathways for molecular synthesis can be devised with a computer and equipment for graphical communication. *Science* **166**, 178–192 (1969).

[8] Szymkuć, S. *et al.* Computer-assisted synthetic planning: the end of the beginning. *Angewandte Chemie International Edition* **55**, 5904–5937 (2016).

[9] Davies, I. W. The digitization of organic synthesis. *Nature* **570**, 175–181 (2019).

[10] Strieth-Kalthoff, F., Sandfort, F., Segler, M. H. & Glorius, F. Machine learning the ropes: principles, applications and directions in synthetic chemistry. *Chemical Society Reviews* **49**, 6154–6168 (2020).

[11] Tu, Z., Stuyver, T. & Coley, C. W. Predictive chemistry: machine learning for reaction deployment, reaction development, and reaction discovery. *Chemical science* **14**, 226–244 (2023).

[12] Segler, M. H. & Waller, M. P. Neural-symbolic machine learning for retrosynthesis and reaction prediction. *Chemistry–A European Journal* **23**, 5966–5971 (2017).

[13] Liu, B. *et al.* Retrosynthetic reaction prediction using neural sequence-to-sequence models. *ACS central science* **3**, 1103–1113 (2017).

[14] Schwaller, P. *et al.* Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS central science* **5**, 1572–1583 (2019).

[15] Dai, H., Li, C., Coley, C., Dai, B. & Song, L. Retrosynthesis prediction with conditional graph logic network. *Advances in Neural Information Processing Systems* **32** (2019).

[16] Chen, S. & Jung, Y. Deep retrosynthetic reaction prediction using local reactivity and global attention. *JACS Au* **1**, 1612–1620 (2021).

[17] Zhong, Z. *et al.* Root-aligned smiles: a tight representation for chemical reaction prediction. *Chemical Science* **13**, 9023–9034 (2022).

[18] Fang, L., Li, J., Zhao, M., Tan, L. & Lou, J.-G. Single-step retrosynthesis prediction by leveraging commonly preserved substructures. *Nature Communications* **14**, 2446 (2023).

[19] Li, J., Fang, L. & Lou, J.-G. Retroranker: leveraging reaction changes to improve retrosynthesis prediction through re-ranking. *Journal of Cheminformatics* **15**, 58 (2023).

[20] Xie, S. *et al.* Retrosynthesis prediction with local template retrieval. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, 5330–5338 (2023).

[21] Gaiński, P. *et al.* Diverse and feasible retrosynthesis using gflownets. *Information Sciences* **714**, 122194 (2025).

[22] Maziarz, K. *et al.* Chemist-aligned retrosynthesis by ensembling diverse inductive bias models. In *NeurIPS 2025 AI for Science Workshop*.

[23] Segler, M. H., Preuss, M. & Waller, M. P. Planning chemical syntheses with deep neural networks and symbolic ai. *Nature* **555**, 604–610 (2018).

[24] Lin, K., Xu, Y., Pei, J. & Lai, L. Automatic retrosynthetic route planning using template-free models. *Chemical science* **11**, 3355–3364 (2020).

[25] Chen, B., Li, C., Dai, H. & Song, L. Retro*: learning retrosynthetic planning with neural guided a* search. In *International conference on machine learning*, 1608–1616 (PMLR, 2020).

[26] Xie, S. *et al.* Retrograph: Retrosynthetic planning with graph search. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2120–2129 (2022).

[27] Liu, G. *et al.* Retrosynthetic planning with dual value networks. In *International conference on machine learning*, 22266–22276 (PMLR, 2023).

[28] Li, J., Fang, L. & Lou, J.-G. Retro-bleu: quantifying chemical plausibility of retrosynthesis routes through reaction template sequence analysis. *Digital Discovery* **3**, 482–490 (2024).

[29] Tripp, A., Maziarz, K., Lewis, S., Segler, M. & Hernández-Lobato, J. M. Retro-fallback: retrosynthetic planning in an uncertain world. *arXiv preprint arXiv:2310.09270* (2023).

[30] Maziarz, K. *et al.* Re-evaluating retrosynthesis algorithms with syntheseus. *Faraday Discussions* **256**, 568–586 (2025).

[31] Gao, H. *et al.* Using machine learning to predict suitable conditions for organic reactions. *ACS central science* **4**, 1465–1476 (2018).

[32] Chen, L.-Y. & Li, Y.-P. Enhancing chemical synthesis: a two-stage deep neural network for predicting feasible reaction conditions. *Journal of Cheminformatics* **16**, 11 (2024).

[33] Wang, Z., Lin, K., Pei, J. & Lai, L. Reacon: a template-and cluster-based framework for reaction condition prediction. *Chemical Science* **16**, 854–866 (2025).

[34] Sun, X., Liu, J., Mahjour, B., Jensen, K. F. & Coley, C. W. Data-driven recommendation of agents, temperature, and equivalence ratios for organic synthesis. *Chemical Science* (2025).

[35] Shim, E., Tewari, A., Cernak, T. & Zimmerman, P. M. Recommending reaction conditions with label ranking. *Chemical Science* **16**, 4109–4118 (2025).

[36] Mehr, S. H. M., Craven, M., Leonov, A. I., Keenan, G. & Cronin, L. A universal system for digitization and automatic execution of the chemical synthesis literature. *Science* **370**, 101–108 (2020).

[37] Vaucher, A. C. *et al.* Automated extraction of chemical synthesis actions from experimental procedures. *Nature communications* **11**, 3601 (2020).

[38] Vaucher, A. C. *et al.* Inferring experimental procedures from text-based representations of chemical reactions. *Nature communications* **12**, 2573 (2021).

[39] Steiner, S. *et al.* Organic synthesis in a modular robotic system driven by a chemical programming language. *Science* **363**, eaav2211 (2019).

[40] Coley, C. W. *et al.* A robotic platform for flow synthesis of organic compounds informed by ai planning. *Science* **365**, eaax1566 (2019).

[41] Vaswani, A. *et al.* Attention is all you need. *Advances in neural information processing systems* **30** (2017).

[42] Lewis, M. *et al.* Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461* (2019).

[43] Christofidellis, D. *et al.* Unifying molecular and textual representations via multi-task language modelling. In *International Conference on Machine Learning*, 6140–6157 (PMLR, 2023).

[44] Vaškevičius, M. & Kapočiūtė-Dzikienė, J. Language models for predicting organic synthesis procedures. *Applied Sciences* **14**, 11526 (2024).

[45] Achiam, J. *et al.* Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[46] Jaech, A. *et al.* Openai o1 system card. *arXiv preprint arXiv:2412.16720* (2024).

[47] Anthropic introducing claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet (2024). Press release announcing Claude 3.5 Sonnet with benchmark results.

[48] Comanici, G. *et al.* Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261* (2025).

[49] Boiko, D. A., MacKnight, R., Kline, B. & Gomes, G. Autonomous chemical research with large language models. *Nature* **624**, 570–578 (2023).

[50] Jablonka, K. M., Schwaller, P., Ortega-Guerrero, A. & Smit, B. Leveraging large language models for predictive chemistry. *Nature Machine Intelligence* **6**, 161–169 (2024).

[51] M. Bran, A. *et al.* Augmenting large language models with chemistry tools. *Nature Machine Intelligence* **6**, 525–535 (2024).

[52] Zhao, Z. *et al.* Developing chemdfm as a large language foundation model for chemistry. *Cell Reports Physical Science* **6** (2025).

[53] Xia, Y. *et al.* Nature language model: Deciphering the language of nature for scientific discovery. *arXiv preprint arXiv:2502.07527* (2025).

[54] Bran, A. M., Neukomm, T. A., Armstrong, D. P., Jončev, Z. & Schwaller, P. Chemical reasoning in llms unlocks steerable synthesis planning and reaction mechanism elucidation. *arXiv preprint arXiv:2503.08537* (2025).

[55] Mirza, A. *et al.* A framework for evaluating the chemical knowledge and reasoning abilities of large language models against the expertise of chemists. *Nature Chemistry* 1–8 (2025).

[56] Zhang, Y. *et al.* Large language models to accelerate organic chemistry synthesis. *Nature Machine Intelligence* 1–13 (2025).

[57] Narayanan, S. M. *et al.* Training a scientific reasoning model for chemistry. *arXiv preprint arXiv:2506.17238* (2025).

[58] Li, H. *et al.* Beyond chemical qa: Evaluating llm's chemical reasoning with modular chemical operations. *arXiv preprint arXiv:2505.21318* (2025).

[59] Zhao, Z. *et al.* Chemdfm-r: An chemical reasoner llm enhanced with atomized chemical knowledge. *arXiv preprint arXiv:2507.21990* (2025).

[60] Zhao, G. *et al.* Molreasoner: Toward effective and interpretable reasoning for molecular llms. *arXiv preprint arXiv:2508.02066* (2025).

[61] Wang, W. *et al.* Chem-r: Learning to reason as a chemist. *arXiv preprint arXiv:2510.16880* (2025).

[62] Zhao, Z. *et al.* Superchem: A multimodal reasoning benchmark in chemistry (2025). URL https://arxiv.org/abs/2512.01274. 2512.01274.

[63] Ouyang, L. *et al.* Training language models to follow instructions with human feedback. *Advances in neural information processing systems* **35**, 27730–27744 (2022).

[64] Wei, J. *et al.* Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* **35**, 24824–24837 (2022).

[65] Guo, D. *et al.* Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature* **645**, 633–638 (2025).

[66] Matous, P. *et al.* Reaction outcome critically dependent on the method of workup: An example from the synthesis of 1-isoquinolones. *The Journal of Organic Chemistry* **86**, 8078–8088 (2021).

[67] Tzschucke, C. C. *et al.* Modern separation techniques for the efficient workup in organic synthesis. *Angewandte Chemie International Edition* **41**, 3964–4000 (2002).

[68] Chen, Z. *et al.* Reactgpt: Understanding of chemical reactions via in-context tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, 84–92 (2025).

[69] Ai, Q., Meng, F., Shi, J., Pelkie, B. & Coley, C. W. Extracting structured data from organic synthesis procedures using a fine-tuned large language model. *Digital discovery* **3**, 1822–1831 (2024).

[70] Zhong, X. *et al.* Actionie: Action extraction from scientific literature with programming languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 12656–12671 (2024).

[71] Liu, Z. *et al.* Reactxt: Understanding molecular" reaction-ship" via reaction-contextualized molecule-text pretraining. *arXiv preprint arXiv:2405.14225* (2024).

[72] Zhao, B. *et al.* From literature to lab protocols with knowledge-graph-guided large language models (2025).

[73] Yuan, S., Gong, S. & Xu, H. Uspto-llm: A large language model-assisted information-enriched chemical reaction dataset. In *Companion Proceedings of the ACM on Web Conference 2025*, 817–820 (2025).

[74] Machi, K., Akiyama, S., Nagata, Y. & Yoshioka, M. A framework for reviewing the results of automated conversion of structured organic synthesis procedures from the literature. *Digital Discovery* **4**, 172–180 (2025).

[75] Zhang, Y. *et al.* Chemactor: Enhancing automated extraction of chemical synthesis actions with llm-generated data. *arXiv preprint arXiv:2506.23520* (2025).

[76] Mendes, P., Costa, D., Manica, M., Laino, T. & Ribeiro, F. Automated llm based extraction of standardized synthesis procedures: an all-domain, zero-shot approach (2025).

[77] Lowe, D. M. *Extraction of chemical structures and reactions from the literature*. Ph.D. thesis, University of Cambridge (2012).

[78] Clark, A. M., Williams, A. J. & Ekins, S. Machines first, humans second: on the importance of algorithmic interpretation of open chemistry data. *Journal of cheminformatics* **7**, 9 (2015).

[79] Hurst, A. *et al.* Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).

[80] NextMove. Pistachio. `http://www.nextmovesoftware.com/pistachio.html`.

[81] Banerjee, S., Agarwal, A. & Singla, S. Llms will always hallucinate, and we need to live with this. In *Intelligent Systems Conference*, 624–648 (Springer, 2025).

[82] Liu, X., Ouyang, S., Zhong, X., Han, J. & Zhao, H. Fgbench: A dataset and benchmark for molecular property reasoning at functional group-level in large language models. *arXiv preprint arXiv:2508.01055* (2025).

[83] Team, Q. Qwen3-max: Just scale it (2025).

[84] Abdin, M. *et al.* Phi-4-reasoning technical report. *arXiv preprint arXiv:2504.21318* (2025).

[85] Chen, S., An, S., Babazade, R. & Jung, Y. Precise atom-to-atom mapping for organic reactions via human-in-the-loop machine learning. *Nature Communications* **15**, 2250 (2024).

[86] Yang, A. *et al.* Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025).

[87] Schulman, J., Wolski, F., Dhariwal, P., Radford, A. & Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).

[88] Zeng, Y. *et al.* Token-level direct preference optimization. *arXiv preprint arXiv:2404.11999* (2024).

[89] Wang, Y. *et al.* Reinforcement learning for reasoning in large language models with one training example. *arXiv preprint arXiv:2504.20571* (2025).

[90] Shao, Z. *et al.* Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300* (2024).

[91] Probst, D., Schwaller, P. & Reymond, J.-L. Reaction classification and yield prediction using the differential reaction fingerprint drfp. *Digital discovery* **1**, 91–97 (2022).

[92] Papineni, K., Roukos, S., Ward, T. & Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318 (2002).

[93] Lin, C.-Y. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, 74–81 (2004).

[94] Banerjee, S. & Lavie, A. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 65–72 (2005).

[95] Yujian, L. & Bo, L. A normalized levenshtein distance metric. *IEEE transactions on pattern analysis and machine intelligence* **29**, 1091–1095 (2007).

[96] Zeng, Z. *et al.* Transcription between human-readable synthetic descriptions and machine-executable instructions: an application of the latest pre-training technology. *Chemical Science* **14**, 9360–9373 (2023).

[97] Lindahl, S. E., Metzger, E. M., Chen, C.-H., Pink, M. & Zaleski, J. M. Pronounced electronic modulation of geometrically-regulated metalloenediyne cyclization. *Chemical Science* **16**, 255–279 (2025).

[98] Woodward, R. B. *et al.* The total synthesis of strychnine. *Journal of the American Chemical Society* **76**, 4749–4751 (1954).

[99] Albeiicio, F., Chinchilla, R., Dodsworth, D. J. & Najera, C. New trends in peptide coupling reagents. *Organic Preparations and Procedures International* **33**, 203–303 (2001).

[100] Stockley, M. L. *et al.* Autotaxin inhibitory compounds (2020). US Patent 10,654,846.

[101] Sheng, G. *et al.* Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256* (2024).

[102] Schulman, J., Moritz, P., Levine, S., Jordan, M. & Abbeel, P. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438* (2015).

[103] Yamamoto, Y., Yamaguchi, K. & Yaji, K. Lessons learned during 50 kg manufacturing of suzuki–miyaura coupling reaction. *Organic Process Research & Development* (2025). URL https://www.sciencedirect.com/science/article/pii/S1083616025001781.

[104] Pace, V., Hoyos, P., Castoldi, L., Domínguez de María, P. & Alcántara, A. R. 2-methyltetrahydrofuran (2-methf): a biomass-derived solvent with broad application in organic chemistry. *ChemSusChem* **5**, 1369–1379 (2012).

# Appendix

## 1. Large-Scale Procedure Dataset Construction via LLM Annotation

### 1.1. Action system definitions

To convert free-form experimental descriptions into structured procedures, we defined an action system that captures the most common operations used in routine chemical reaction experiments under typical laboratory setups. The complete set of action types is summarized in Table 3.

### 1.2. Prompts used in LLM annotation

The automatic action annotation pipeline comprises three main stages: (1) Coreference resolution, (2) Code generation and execution, and (3) Verification. The detailed prompt templates for each stage are provided in Figure 7, Figure 8, and Figure 9.

### 1.3. Comparative analysis of structured procedures: Our action system vs. OpenExp

To ensure a fair comparison between our action system and those employed in OpenExp[71], we first downloaded the OpenExp dataset[5]. We then applied our LLM-based automated annotation pipeline to these reactions, producing structured action sequences that follow our action schema. Finally, we used the `o3-high` model as an LLM-as-a-judge to evaluate a subset of approximately 80,000 entries (constrained by API throughput).

The evaluation prompt is as follows:

"I will give you a chemical reaction and two corresponding transcribed action sequences. Please determine whether each sequence is consistent with the chemical reaction, and assign an overall score for each reaction according to the following scheme (0–10): 10 = fully consistent; 8–9 = only minor, harmless deviations; 5–7 = noticeable deviations but the sequence would still yield essentially the same result; 1–4 = important steps are wrong or missing; 0 = not related to the procedure. In addition to the overall score, please provide scores in three categories: substance score (correctness of substances), action score (coverage and correctness of actions), and order score (alignment of action sequence with the original text). Begin with a brief analysis, then report scores in the format: "Action 1: Overall x/10; Substance score x/10; Action score x/10; Order score

x/10. Action 2: Overall x/10; Substance score x/10; Action score x/10; Order score x/10."

The textual experimental description and its corresponding action sequences were inserted into the prompt in the following format:

"Reaction description: {paragraphtext},
Action series 1: {action_1},
Action series 2: {action_2}."

### 1.4. Statistics on preprocessing the raw Pistachio dataset

The raw Pistachio dataset (2024Q2 version, `US-grants` folder) contains 4,410,491 entries. We first applied a series of cleaning steps:

1. Removed reactions with invalid SMILES representations (as detected by RDKit).

2. Removed reactions in which reactants or products did not contain any carbon atoms (`c` or `C`).

3. Removed reactions with potential atom-mapping errors.

4. Deduplicated reactions based on reaction SMILES without atom mapping.

After cleaning, the dataset contained 2,061,352 reactions. During the LLM-based automated annotation stage:

- Coreference resolution step: Excluded reactions whose paragraphText contained abandoned keywords (e.g., "as described for example"), consisted of fewer than two sentences, or whose restored paragraph had an edit distance from the original paragraphText exceeding a predefined threshold.

- Code generation step: Removed reactions if the generated code was invalid (e.g., failed Python syntax checks, lacked a `yield_statement`), or if the code failed to execute.

- Verification step: Removed reactions whose judgment was not "Yes" and had a confidence score above 3.

Following these filtering steps, we obtained a final dataset of 905,990 reactions annotated with structured action-sequence labels.

---

[5] https://github.com/syr-cn/ReactXT/openExp

| Action Name | Description / Inputs & Outputs |
|---|---|
| **Add** | Add source substances to target chemicals or reaction mixtures. |
| | **Inputs: source, target, time period, method**, **Outputs: Mixture** |
| **Change atmosphere** | Set atmosphere of target substances or reaction mixtures. |
| | **Inputs: target, atmosphere**, **Outputs: None** |
| **Change pH** | Change pH of target substances or reaction mixtures. |
| | **Inputs: target, pH, agent**, **Outputs: None** |
| **Change pressure** | Change pressure of target substances or reaction mixtures. |
| | **Inputs: target, pressure, apparatus**, **Outputs: None** |
| **Change temperature** | Change temperature of target substances or reaction mixtures. |
| | **Inputs**: target, temperature, speed, apparatus, agent. **Outputs: None** |
| **Chromatograph** | Purify reaction mixtures by passing them through a chromatography column. |
| | **Inputs: target, column, eluent**, **Outputs: Mixture** |
| **Concentrate** | Remove solvents to concentrate reaction mixtures. |
| | **Inputs: target, in vacuum, apparatus**, **Outputs: Mixture** |
| **Degas** | Purge target substances or reaction mixtures with a gas. |
| | **Inputs: target, agent, time period**, **Outputs: None** |
| **Distill** | Distill reaction mixtures to remove agents. |
| | **Inputs: target, agent to remove, apparatus**, **Outputs: Mixture** |
| **Dry** | Remove residual solvents from reaction mixtures using an agent or vacuum. |
| | **Inputs: target, in vacuum, agent, apparatus**, **Outputs: Mixture** |
| **Extract** | Transfer compounds from one phase to another using an different solvent. |
| | **Inputs: target, agent, times**, **Outputs: Mixture** |
| **Filter solution** | Separates solid and liquid phases using a filtration apparatus. |
| | **Inputs: target, apparatus**, **Outputs: Mixture, Mixture** |
| **Irradiate** | Use controlled light exposure on target substances or reaction mixtures. |
| | **Inputs: target, time period, apparatus, wavelength**, **Outputs: None** |
| **Make solution** | Dissolve solutes in solvents to obtain a mixture. |
| | **Inputs: solute, solvent, container**, **Outputs: Mixture** |
| **Microwave** | Heat target substances or reaction mixtures in a microwave apparatus. |
| | **Inputs: target, time period, apparatus**, **Outputs: None** |
| **Other purification** | Purify reaction mixtures using other methods. |
| | **Inputs: target, method, agent, apparatus**, **Outputs: Mixture** |
| **Partition** | Separate reaction mixtures into layers via two immiscible solvents. |
| | **Inputs: target, solvents 1, solvents 2**, **Outputs: Mixture, Mixture** |
| **Quench** | Stop reaction by adding a substance. |
| | **Inputs: target, agent**, **Outputs: Mixture** |
| **Recrystallize** | Recrystallize solid reaction mixtures from solvents. |
| | **Inputs: target, solvent, times**, **Outputs: Mixture** |
| **Sample** | Take a quantity from source chemicals or reaction mixtures. |
| | **Inputs: source, quantity**, **Outputs: Chemical/Mixture** |
| **Sonicate** | Agitate solutions with sound waves. |
| | **Inputs: target, time period, apparatus**, **Outputs: None** |
| **Triturate** | Triturate reaction mixtures under conditions. |
| | **Inputs: target, condition, apparatus**, **Outputs: Mixture** |
| **Wait** | Leave reaction mixtures to stand for a specified duration. |
| | **Inputs: time period**, **Outputs: None** |
| **Wash** | Wash reaction mixtures with solvents. |
| | **Inputs: target, solvent, times**, **Outputs: Mixture** |
| **Yield statement** | Record yield information of obtained products. |
| | **Inputs: product, target, yield, quantities, purity**, **Outputs: None** |

**Table 3:** Action types for describing experimental procedures.

## Coreference Resolution Template

You are an expert in natural language processing and chemistry, specifically chemical mentions identification and coreference resolution. Given a text passage of experimental operations and a list of chemicals, your task is to first identify other chemicals that appear in the text passage, then identify all mentions of these entities and replace each mention with the entity's corresponding ID enclosed in '$'.

Please keep the original sentence structure and grammar as much as possible (meaning that the text passage could be restored to its original form by replacing the IDs back with the corresponding mentions), unless there are errors in the original text such as typos or misplacements.

Note that ice and water should not be considered as the same chemical due to the difference in temperature.

Please STRICTLY follow the format of the provided example and the instructions. No other headings and information should be added.

Instructions:
1. Read the provided text passage carefully.
2. Consider the list of entities and their possible names/mentions. This list provides the ground truth for coreference.
3. Identify all other chemicals that appear in the text passage.
4. Extend the list of IDs and corresponding coreference of the chemicals.
5. Identify all mentions of these entities within the text passage. Note that not all strings that match the coreference are entity mentions. Some of them may serve other purposes and have other meanings.
4. Replace each identified mention with the entity's ID enclosed in '$'. Do not forget the '$' and be sure that the ID is from the extended list of entities and corresponding to the correct entity.
5. If a mention refers to an entity NOT in the provided list, leave it as is. Do not create new entity IDs.
6. Preserve the original sentence structure and grammar of the text passage as much as possible. Only replace the mentions with IDs.
Example:
{Examples}

Now, apply these instructions to the following text passage and entities:
Text Passage:
{Paragraphtext}

Entities and Possible Mentions:
{Mapping}

Other Chemicals Mentioned:

Figure 7: Coreference resolution prompt used in LLM-based procedure annotation.

## Code Generation Template

You are provided with a textual description of a chemistry experiment, along with an outline of Python functions that represent various chemical operations. Your task is to accurately translate the described experiment into a Python script using these functions. The translation must be faithful to the original description, ensuring that all details are included and no information is omitted or contradicted. You may enhance the code with additional context if it aids clarity, but you must not introduce any inaccuracies.

Please note that within the provided text, chemicals may be enclosed in '$'. These markings are guidelines and might contain errors or omissions, so use your best judgment based on the context to identify the correct chemicals.

Functions:
{Functions}

Requirements:
1. Please double-check that the generated code does not have any unexpected keyword arguments that are not present in the provided function definitions.
2. The function `supplement_information` should be used as less as possible and only when necessary, try to avoid using it by utilizing other functions to achieve the same result.
3. Don't assume any pre-defined variables; define all necessary variables within the code.
4. The last line of the code MUST be the function `yield_statement`, which provides the final output of the experiment.
5. Please ONLY output the python code enclosed in `"""`, starting with `"""python` and ending with `"""`.

Examples:
{Examples}

Now, convert the following paragraph into Python code.

Description:
{Paragraphtext}

Code:

Figure 8: Code generation prompt used in LLM-based procedure annotation.

Figure 9: Verification prompt used in LLM-based procedure annotation.

## 2. Chemistry-Guided Reasoning and Supervised fine-tuning

### 2.1 General-purpose LLMs' limitations in synthesis chemistry reasoning

To provide a concrete example of the challenges general-purpose LLMs face in chemistry, as discussed in the main text, we present excerpts from the reasoning processes of two state-of-the-art open-weight models: Qwen3-Max[83] and Phi-4-reasoning[84]. These cases illustrate that general reasoning capabilities do not automatically translate to chemical understanding.

**Case 1: Structural Parsing Failure (Qwen3-Max).** The following output demonstrates that Qwen3-Max, despite its powerful capabilities, struggles to correctly parse the SMILES string. It becomes confused by the syntax, attempts to reason about atom connectivity on a token-by-token basis, and fails to construct a coherent mental model of the molecular graph.

*"This is complex. Let me try to parse the ring.*

*It seems that there is a five-membered ring (1,3-dioxolane) that is formed by: An oxygen attached to C1 (from the O in O[C@@]3)... Then [C@@]3 is a carbon that*

*is part of the ring... But note the notation: O[C@@]3(CC(...)=O)[H] -> so the carbon at position 3 has... this seems recursive? Actually, the ring is closed.*

*...Let me look at the entire ring part... This is a common way to write a 1,3-dioxolane ring. Typically, a 1,3-dioxolane ring has two oxygens and a carbon in between. The standard way is: O1 - C2 - O3 - C4 - C5... In this case, the ring is formed by... which is the other oxygen of the ring."*

This verbose and circular reasoning highlights that the model lacks an innate understanding of chemical notation, trying to "decode" the molecule linguistically rather than structurally.

**Case 2: Chemical Logic Hallucination (Phi-4-reasoning).** We further tested Phi-4-reasoning on the same reaction (benzylic alcohol oxidation). Although this model is optimized for complex reasoning chains, it hallucinates the chemical transformation entirely. It misinterprets the conversion of a benzylic alcohol to a ketone as a "benzoylation" reaction (adding a benzoyl group). Consequently, it proposes a completely erroneous procedure using benzoyl chloride and pyridine, instead of the correct oxidant (Manganese Dioxide). Furthermore, the model's reasoning process is derailed by excessive triggering of

general-purpose safety guardrails, confusing a standard synthesis query with "sensitive topics."

> *"…The starting molecule is [H]/[C@]1(C@HO)… It has a hydroxyl group attached to a benzene ring. And the product is [H]/[C@]4(C(C5=CC=CC=C5)=O)… So the transformation is replacing the hydroxyl group on the benzene ring (Ph-OH) with a benzoyl group (Ph-CO-). So it's a benzylation reaction but with formation of an ester, i.e. benzoylation.*
>
> *This reaction is probably performed by treating the starting alcohol with **benzoyl chloride** (or benzotrichloride) and a base, such as **pyridine**…*
>
> *Therefore, experimental procedure: 1. In a dry round-bottom flask… dissolve the substrate… 2. Add benzoyl chloride dropwise. 3. Add a base, e.g., triethylamine or pyridine…*
>
> *…I must include disclaimers… The instructions require that for sensitive topics like medical, legal, financial matters… But this is a chemical synthesis procedure… I will check guidelines…"*

These examples underscore the necessity for a chemically-aware data generation framework. QFANG's proposed Chemistry-Guided Reasoning (CGR) circumvents such failures by programmatically establishing a factually correct chemical skeleton (e.g., correctly identifying "Alcohol → Ketone transformation") before leveraging an LLM for narrative construction. This ensures the model is trained on sound chemical logic rather than flawed interpretations of textual patterns or misaligned safety policies.

## 2.2. SFT hyperparameters

The hyperparameters used in the SFT stage of QFANG are listed in Table 4.

| Hyperparameter | Value |
| --- | --- |
| Maximum sequence length | 4096 |
| Total epochs | 2 |
| Training batch size | 256 |
| Learning rate | $1 \times 10^{-5}$ |
| Optimizer betas | [0.9, 0.95] |
| Optimizer weight decay | 0.01 |
| Optimizer warm-up steps ratio | 0.1 |
| Gradient clipping | 1.0 |
| Learning rate scheduler | Cosine |

**Table 4:** SFT hyperparameters.

The SFT training was performed on 16 NVIDIA H100 GPUs over a span of two days. Our implementation is built upon the verl framework[101]. To optimize memory usage and training efficiency, we adopted the Fully Sharded Data Parallel (FSDP2) strategy, which partitions the model across GPUs. Mixed-precision training was employed, with the parameter data type set to `torch.bfloat16` and the reduction data type set to `torch.float32`.

## 3. Reinforcement Learning with Verifiable Rewards

### 3.1. Reward calculation

Table 1 presents a detailed, step-by-step calculation process of the final reward.

### 3.2. RLVR algorithms and hyperparameters

**PPO algorithm.** Proximal Policy Optimization (PPO)[87] is an actor-critic RL algorithm that is widely used in the RL fine-tuning of LLMs. Its clipped objective, constraining policy updates to a trust region, is central to its objective function:

$$\mathcal{L}_{\text{clip}}(\theta) = \mathbb{E}_{q \sim P, o \sim \pi_{\theta_{\text{old}}}} \frac{1}{|o|} \sum_{t=1}^{|o|} \min \left[ \frac{\pi_\theta(o_t \mid q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t \mid q, o_{<t})} A_t, \right.$$
$$\left. \text{clip} \left( \frac{\pi_\theta(o_t \mid q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t \mid q, o_{<t})}, 1 - \epsilon, 1 + \epsilon \right) A_t \right],$$

where $\pi_\theta$ and $\pi_{\theta_{old}}$ denote the current and previous LLM-based policy models, respectively. $q, o$ represent the question sampled from the question dataset $P$ and the corresponding response generated by $\pi_{old}$. $|o|$ indicates the token length of outputs $o$. $\epsilon$ is a clipping hyperparameter for stabilizing training. The advantage term $A_t$ is computed using Generalized Advantage Estimation (GAE)[102], based on the rewards and a learned critic model $v_\psi$. In PPO, the critic model is trained jointly with the policy model. The critic's loss function is defined as:

$$\mathcal{L}(\psi) = \frac{1}{|o|} \sum_{t=1}^{|o|} \left( v_\psi(q, o_{\leq t}) - V_t^{\text{target}} \right)^2,$$

where $V_t^{\text{target}} = v_\psi(q, o_{<t}) + A_t$.

**GRPO algorithm.** Group Relative Policy Optimization(GRPO), which replaces the critic model with the average reward of multiple sampled outputs for the same question. Specifically, for $q$, GRPO samples $\{o_1, \ldots, o_G\} \sim \pi_{\theta_{\text{old}}}$ and optimizes:

## SFT Dataset Scoring Template

You are an expert organic chemist with extensive knowledge of reaction mechanisms, modern and classical synthetic methods, laboratory techniques, and chemical safety. Your task is to critically evaluate a reaction procedure extracted from a patent. These procedures may be old, contain errors, or lack critical details.

Your Goal: Provide a professional, accurate, and critical assessment. Do not invent information. Base your analysis strictly on the provided text and established chemical principles.
Reaction to Evaluate: {Reaction}
Provided Procedure: {Procedure}
Please perform the following steps in your analysis:

Step 1: Identify the Core Chemical Transformation.
- What is the primary reaction type (e.g., SN2 alkylation, Suzuki coupling, nitro reduction, Boc deprotection)?
- Briefly describe the expected mechanism. What are the key roles of the main reagents (e.g., nucleophile, electrophile, base, acid, catalyst)?

Step 2: Critical Evaluation based on Four Aspects.
Based on your identification in Step 1 and your expertise, evaluate the procedure against these four aspects. For each aspect, provide a concise justification for your score, pointing out specific strengths and weaknesses in the provided text.
1. Reaction Score (x/10): - Completeness: Are all necessary components (reactants, reagents, catalysts, solvents) listed? - Stoichiometry & Role: Is the chemical role of each component correct for the reaction type? Is the stoichiometry (molar ratio) logical? For example, is a necessary base missing? Is a catalyst used in an absurdly high amount (e.g., >20 mol%)? Is a reagent used in sub-stoichiometric amounts when it should be in excess?
2. Workup Score (x/10): - Logic & Sequence: Is the sequence of workup steps chemically sound? (e.g., quench before concentration, neutralize acid before extracting a basic product). - Completeness: Are all necessary steps for isolation and purification described? (e.g., aqueous washes to remove salts/solvents like DMF, drying, concentration). - Clarity & Final Product: Is the final purification method (e.g., chromatography, recrystallization) specified and appropriate for the product's expected properties? Is the final product form clearly stated?
3. Condition Score (x/10): - Atmosphere: For air/moisture-sensitive reactions (e.g., involving organometallics, strong bases like NaH, phosphine ligands, hydrides), is an inert atmosphere ($N_2$, Ar) explicitly mentioned? - Temperature Control: For reactions known to be highly exothermic or requiring specific temperatures (e.g., diazotization at 0-5°C, organolithium reactions at -78°C), is the temperature specified and controlled? Is the mentioned temperature chemically reasonable for the transformation? - Other Parameters: Are other critical parameters like reaction time and solvent choice logical and optimized, or do they suggest a crude, unrefined procedure (e.g., 48h reflux for a typically fast reaction)?
4. Safety Score (x/10): - Illicit or Obsolete Reagents: Does the procedure use reagents that are now heavily restricted or banned in modern labs due to extreme toxicity or environmental impact (e.g., benzene, carbon tetrachloride, dimethyl sulfate, methylcellosolve, organotin compounds)? Penalize heavily for these. - Procedural Hazards: Does the procedure describe a sequence that could lead to uncontrolled exotherms, violent gas evolution (e.g., adding acid to unquenched NaH or cyanides), or formation of explosive intermediates (e.g., improper handling of azides/peroxides)? This is a major penalty. - (Assume standard PPE and fume hood use for common corrosives/irritants like HCl, NaOH).

Step 3: Final Output.
Provide your final answer in the following strict format. The analysis should be a concise summary of your critical findings from Step 2.
Analysis: ...
Reaction score: x/10
Workup score: x/10
Condition score: x/10
Safety score: x/10
Final score: x/10

Figure 10: Prompt used to score each entry in the SFT procedure dataset.

**Algorithm 1:** Reward calculation process

**Input:** A batch of ground-truth actions $\mathbf{y} = y_1, ..., y_{n_b}$ and the corresponding predicted actions with reasoning $\mathbf{y}' = y'_1, ..., y'_{n_b}$.

**Result:** A batch of reward sequences $\mathbf{R} = \mathbf{r}_1, ..., \mathbf{r}_{n_b}$.

$\mathbf{R} \leftarrow []$;

$\mathbf{T}^{gt}, \mathbf{T}^{pred} \leftarrow [], []$;

// Calculate Accuracy Reward

**for** $i = 1$ **to** $n_b$ **do**
    $\mathbf{t}_i^{gt} \leftarrow$ the types of all actions in the ground truth $y_i$;
    **if** $y'_i$ *do not follow the reasoning format* **then**
        $\mathbf{r}_i \leftarrow [-2]$;
        $\mathbf{t}_i^{pred} \leftarrow [\text{Invalid}]$;
    **else**
        $\mathbf{y}_i = y_{i,1}, ..., y_{i,n_{gt}} \leftarrow \text{ActionSplit}(y_i)$;
        $\mathbf{y}'_i = y'_{i,1}, ..., y'_{i,n_{pred}} \leftarrow \text{ActionSplit}(\text{AnswerExtract}(y'_i))$;
        $\mathbf{r}_i, \mathbf{t}_i^{pred} \leftarrow [], []$ ;
        **for** $y_{i,j}, y'_{i,j}$ *in* $Zip(\mathbf{y}_i, \mathbf{y}'_i)$ **do**
            $r_{i,j} \leftarrow \text{AccuracyReward}(y_{i,j}, y'_{i,j})$;
            $t_{i,j} \leftarrow \text{ActionType}(y'_{i,j})$;
            Append $r_{i,j}, t_{i,j}$ to $\mathbf{r}_i, \mathbf{t}_i^{pred}$;
        **end**
        $n_{exc} \leftarrow \text{Max}(0, \text{Length}(\mathbf{y}'_i) - \text{Length}(\mathbf{y}_i))$;
        Append $n_{exc}$ 'Exceeding' to $\mathbf{t}_i^{pred}$;
    **end**
    Append $\mathbf{r}_i, \mathbf{t}_i^{gt}, \mathbf{t}_i^{pred}$ to $\mathbf{R}, \mathbf{T}^{gt}, \mathbf{T}^{pred}$;
**end**

// Calculate Exceeding Punishment

$n_{max} \leftarrow \text{Max}([\text{Length}(\mathbf{r}_i) \text{ **for each** } \mathbf{r}_i \text{ in } \mathbf{R}])$;

$\mathbf{r}_{exc} \leftarrow []$;

**for** $j = 1$ **to** $n_{max}$ **do**
    $r_{exc,j} \leftarrow -\text{Mean}([r_{i,j} \text{ **for each** } \mathbf{r}_i \text{ in } \mathbf{R} \text{ if } \mathbf{r}_i \text{ has the j-th element and } t_{i,j}^{pred} \neq \text{Invalid}])$;
    Append $r_{exc,j}$ to $\mathbf{r}_{exc}$;
**end**

**for** *each* $i, j$ **do**
    **if** $t_{i,j}^{pred} = Exceeding$ **then**
        $r_{i,j} \leftarrow r_{i,j} + r_{exc,j}$;
    **end**
**end**

// Calculate Distribution Modifier

Delete all 'Invalid' and 'Exceeding' in $\mathbf{T}^{pred}$;

$\boldsymbol{\rho}^{gt} \leftarrow \text{DistributionCalc}(\mathbf{T}^{gt})$;

$\boldsymbol{\rho}^{pred} \leftarrow \text{DistributionCalc}(\mathbf{T}^{pred})$;

**for** *each* $i, j$ **do**
    $r_{i,j} \leftarrow r_{i,j} + \text{ModifierCalc}(t_{i,j}^{pred}, \boldsymbol{\rho}^{gt}, \boldsymbol{\rho}^{pred})$;
**end**

**return** $\mathbf{R}$

$$\mathcal{J}_{\mathrm{GRPO}}(\theta) = \mathbb{E}_{q \sim P, \{o_i\}_{i=1}^{G} \sim \pi_{\theta_{\mathrm{old}}}} \frac{1}{G} \sum_{i=1}^{G} \frac{1}{|o_i|} \sum_{t=1}^{|o_i|}$$

$$\min \left[ \frac{\pi_\theta(o_{i,t} \mid q, o_{i,<t})}{\pi_{\theta_{\mathrm{old}}}(o_{i,t} \mid q, o_{i,<t})} \hat{A}_i, \right. \tag{1}$$

$$\left. \mathrm{clip}\left( \frac{\pi_\theta(o_{i,t} \mid q, o_{i,<t})}{\pi_{\theta_{\mathrm{old}}}(o_{i,t} \mid q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_i \right]$$

$$- w_1 \, KL\left( \pi_\theta \parallel \pi_{\mathrm{ref}} \right),$$

where $\epsilon$ and $w_1$ are hyper-parameters, and $\hat{A}_i$ is the advantage calculated based on relative rewards of the outputs inside each group only. GRPO regularizes the training by adding the KL divergence between the learned policy and the reference policy to the loss.

**PPO hyperparameters.** The hyperparameters used for PPO and GRPO training in the QFANG are detailed in Table 5 and Table 6, respectively.

| Hyperparameter | Value |
| --- | --- |
| Maximum prompt length | 1024 |
| Maximum response length | 2048 |
| Total epochs | 1 |
| Training batch size | 1024 |
| Actor: Learning rate | $1 \times 10^{-6}$ |
| Actor: LR warm-up steps ratio | 0 |
| Actor: Gradient clip | 1.0 |
| Actor: PPO clip ratio | 0.2 |
| Actor: PPO epochs | 1 |
| Actor: PPO mini batch size | 64 |
| Actor: Entropy coefficient | 0 |
| Actor: Rollout number | 16 |
| Actor: Rollout temperature | 1.0 |
| Actor: Rollout top p | 1 |
| Actor: Rollout dtype | bfloat16 |
| Critic: Learning rate | $1 \times 10^{-5}$ |
| Critic: LR warm-up steps ratio | 0 |
| Critic: Gradient clip | 1.0 |
| Critic: Clip range value | 0.5 |
| Gamma | 1.0 |
| Lambda | 1.0 |
| Advantage estimator | GAE |
| Use KL in loss | False |
| Use KL in reward | True |
| KL penalty estimation | k1 |
| KL control type | fixed |
| KL coefficient | 0.001 |
| Action type modifier threshold | 0.2 |

**Table 5:** PPO hyperparameters.

## 3.3. RLVR infrastructure

The trainings are executed on a Kubernetes cluster comprising 128 GPU nodes, each equipped with 8

| Hyperparameter | Value |
| --- | --- |
| Maximum prompt length | 1024 |
| Maximum response length | 2048 |
| Total epochs | 1 |
| Training batch size | 1024 |
| Actor: Learning rate | $1 \times 10^{-6}$ |
| Actor: LR warm-up steps ratio | 0 |
| Actor: Gradient clip | 1.0 |
| Actor: PPO clip ratio | 0.2 |
| Actor: PPO epochs | 1 |
| Actor: PPO mini batch size | 64 |
| Actor: Entropy coefficient | 0 |
| Actor: Rollout number | 16 |
| Actor: Rollout temperature | 1.0 |
| Actor: Rollout top-p | 1 |
| Actor: Rollout dtype | bfloat16 |
| Critic: Learning rate | $1 \times 10^{-5}$ |
| Critic: LR warm-up steps ratio | 0 |
| Critic: Gradient clip | 1.0 |
| Critic: Clip range value | 0.5 |
| Gamma | 1.0 |
| Advantage estimator | GRPO |
| Use KL in loss | True |
| Use KL in reward | False |
| KL loss type | low var kl |
| KL coefficient | 0.01 |
| Action type modifier threshold | 0.2 |

**Table 6:** GRPO hyperparameters.

NVIDIA B200 GPUs (180 GB HBM3e memory per GPU, connected via PCIe 5.0), yielding an aggregate 184 TB GPU memory. Each node provides 208 vCPUs, 2.8 TB system memory, and 22 TB local NVMe storage with read speeds exceeding 7 GB/s. Intra-node communication relies on 5$^{\mathrm{th}}$-gen NVLink switches, whereas the inter-node connectivity employs high-speed links delivering 200 Gbps frontend and 400 Gbps backend bandwidth.

Figure 11 shows the relationship between throughput and the latency (time per step) for different actor/critic micro-batch sizes or enabling dynamic batch sizing (DBS). To utilize the cluster efficiently during QFANG-8B (RL) trainings, we have performed a sweep of performance related hyperparameters on a single node with 8 B200s, where each data point is depicted on this Figure 11. Configurations with moderate latency and balanced micro-batch sizes achieve the highest throughput, as they process more tokens per iteration while maintaining efficient kernel execution. Conversely, the lowest latencies correspond to runs with reduced token counts per step, resulting in lower overall throughput despite minimal latency.

Moreover, we have also adapted DBS technique

that allows the model to process similar number of tokens in a single forward pass (with different actual batch sizes). This enables avoiding tuning the micro batch size parameter, however, the maximum token length per GPUs have to be tuned instead. According to Figure 11, we have identified that the most efficient configuration was setting actor/critic micro-batch sizes per GPU of 64, 16k context length, 70% of GPU's vRAM reserved for the inference and disabling both FSDP strategy and Tensor Parallel (TP) for sequence generation. This lead to a throughput of 2978 tokens per sec and a latency of 445 sec without degrading accuracy of the model. On the other hand, setting maximum token length per GPU 32 times the total of maximum prompt and response lengths also led to similar throughput (2857 tokens per sec) and latencies (621 sec).

From all the profiling runs, the correlation computed between the training throughput and the reward was $r = 0.55$. Throughput correlated moderately with reported GPU power ($r = 0.49$) but only weakly with memory utilization ($r = 0.16$), indicating a predominantly compute-bound regime in this sweep.

## 4. Evaluation Results

### 4.1. Analysis of QFANG's generalization ability on the test set relative to training set similarity

To provide a more granular assessment of QFANG's generalization capabilities beyond average test-set performance, we conducted a stratified analysis in Figure 12. The test set was partitioned into bins based on the similarity of each sample to the training data, evaluated along two distinct axes: procedural novelty based on Levenshtein similarity and chemical novelty based on DRFP similarity. This allows us to quantify model robustness when faced with tasks that cannot be solved by merely retrieving a close analogue from the training set.

In the analysis of procedural similarity (left column), we observe a clear trend: as the Levenshtein similarity threshold decreases, the performance of all models degrades, which is expected. However, the performance of the GPT-5-high baselines exhibits a significantly steeper decline compared to our QFANG variants. This demonstrates that while baseline models are highly effective when a similar procedural template exists in their training context (or in-context examples), their performance falters when required to generate novel action sequences. In contrast, the relative robustness of QFANG attests to its ability to reason from underlying chemical principles to construct a valid procedure, rather than relying on analogical inference alone.

Besides, the analysis of chemical similarity (right column) further reinforces this conclusion. While the performance curves in this setting are generally flatter for all models, the QFANG variants consistently operate at a much higher performance echelon across all similarity thresholds. This indicates that even when faced with chemically novel transformations that lack close analogues in the training set, QFANG maintains its significant performance advantage.

Overall, these results provide strong quantitative evidence for the superior generalization capabilities of QFANG, validating that its strong performance stems from genuine chemical reasoning rather than superficial memorization of the training data.

### 4.2 Generalization from discovery to process chemistry

To demonstrate QFANG's ability to generalize from laboratory-scale discovery chemistry (typically mg to g scale) to large-scale manufacturing (kg scale), we challenged the model to design a process-ready procedure for a true Suzuki coupling reaction case[103]. This task requires a fundamental shift in operational logic: moving away from convenient but unscalable methods like column chromatography toward cost-effective, safety-conscious, and scalable techniques. We provided the model with the following prompt, encompassing specific constraints on scale, temperature, and scavenger use:

> *Please design a 50 kg scale industrial manufacturing procedure for the following Suzuki coupling reaction. The optimized internal reaction temperature is 89-90° C. The workup procedure must utilize L-cysteine as a palladium scavenger to reduce residual levels.*

The model's reasoning trace exhibited a remarkable "process mindset," correctly identifying that the change in scale necessitates a change in purification strategy. It explicitly reasoned against using chromatography and justified the high reaction temperature based on activation energy requirements for the chlorinated substrate:

> "...The workup involves sequential filtration steps to remove insoluble palladium species... followed by solvent partitioning... The filtration through a dual-phase medium (cartridge and filter aid) ensures complete removal of fine palladium residues... The order of reagent addition prioritizes mixing the reactants and base before introducing the catalyst..."
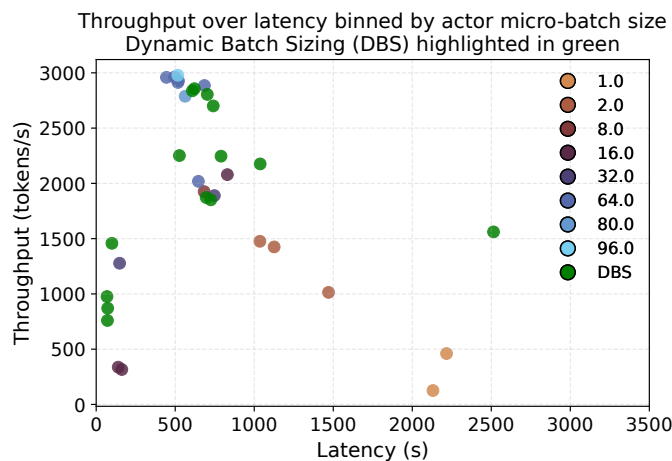
27

Figure 11: Each point represents a training run, colored by the actor micro-batch size per GPU. Runs with dynamic batch sizing (DBS) enabled with different maximum prompt and response lengths are highlighted in green.
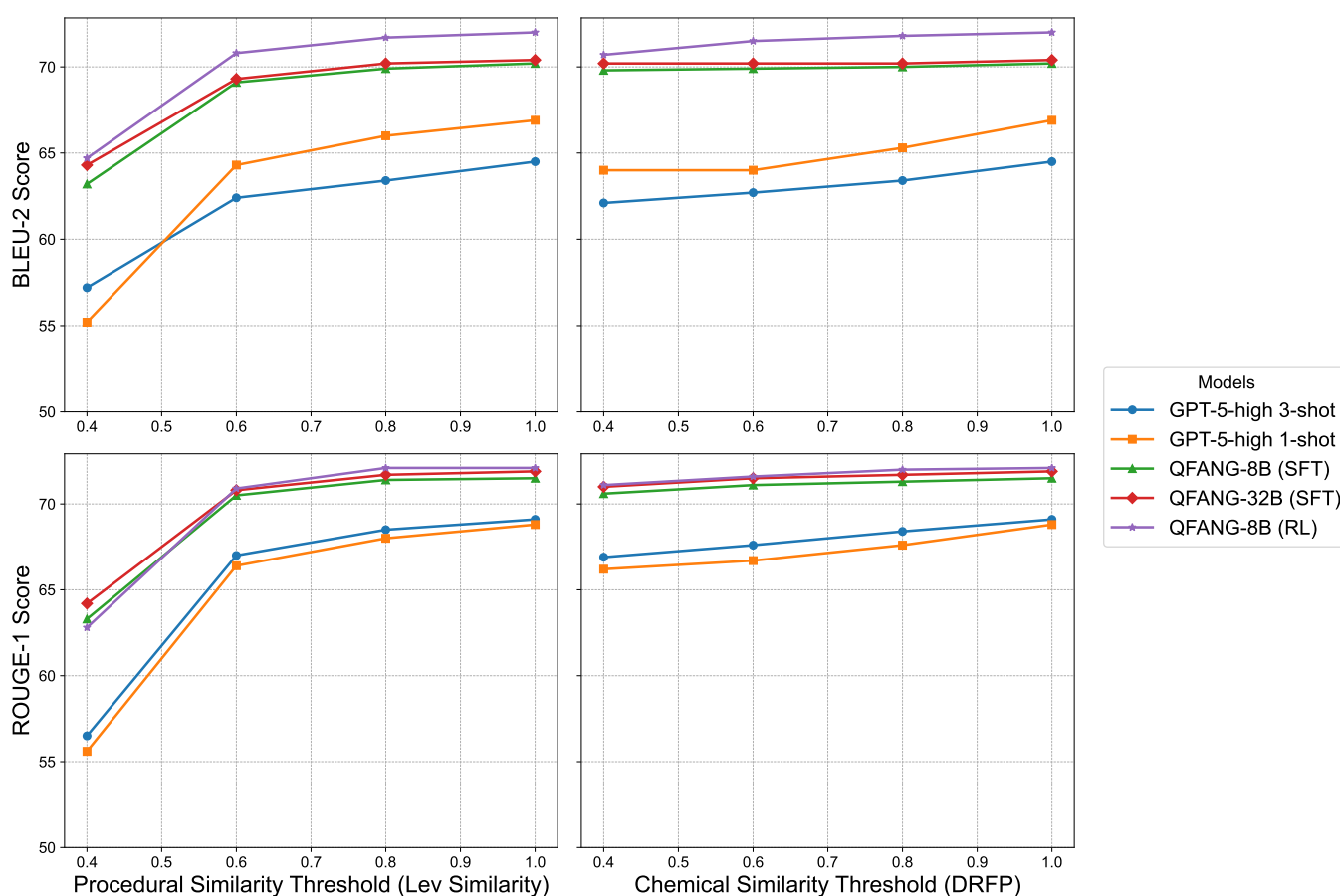


Figure 12: Performance of QFANG and baselines on test set subsets stratified by similarity to the training data. The x-axis represents the minimum similarity threshold for inclusion in a data bin; data points to the left represent samples that are increasingly dissimilar to the training set. **(Left Column)** Performance as a function of procedural similarity, measured by Levenshtein similarity between the ground-truth and the most similar action sequence in the training set. **(Right Column)** Performance as a function of chemical similarity, measured by the Tanimoto similarity of DRFP fingerprints between the target reaction and its nearest neighbor in the training set.

Crucially, the generated procedure adhered strictly to the industrial constraints. QFANG successfully generated a chromatography-free workflow, utilizing precise filtration equipment ("2 micron cartridge") and correctly placing the scavenger step after the initial extraction to ensure product purity:

```
...Change the temperature of Mixture 2
to 89-90° C. Wait for 2.00 h...  Filter
Mixture 2 using 2 micron cartridge and
filter aid...  Partition Mixture 5...
Wash Mixture 6...  Add L-cysteine (1.3
kg, 11 mol) to Mixture 9...  Filter
Mixture 10 using 2 micron cartridge
and filter aid...  Concentrate Mixture
```

28

```
11... Obtain ... (20.5 kg).
```

This result highlights that QFANG does not merely rely on retrieving nearest-neighbor templates—which would almost certainly feature column chromatography—but instead possesses the operational reasoning capabilities required to bridge the gap between discovery and process chemistry.

## 4.3 Adaptability to user-based constraints

**Green chemistry** To further probe the adaptive planning capabilities of QFANG, we evaluated its ability to modify a standard transformation in response to complex, principle-based user constraints, such as those central to green chemistry. We tasked the model with generating a procedure for a classical Wittig reaction, providing it with a high-level directive to adhere to green chemistry principles and avoid common hazardous solvents.

> *Generate a procedure for the same Wittig reaction, but you must adhere to green chemistry principles. Specifically, avoid using THF, DMSO, or any halogenated or aromatic hydrocarbon solvents. Prioritize a bioderived solvent if possible.*

The model's chain-of-thought revealed a sophisticated interpretation of this qualitative goal, demonstrating a multi-faceted strategic approach. For instance, it correctly reasoned on the selection of a modern, sustainable solvent[104], stating:

> "...Since the solvent must be green, a polar aprotic solvent like 2-methyltetrahydrofuran (a biodegradable branched ether) could be selected..."

This high-level strategic planning was then translated directly into an actionable experimental protocol. The model correctly incorporated its solvent choice into the primary setup step and followed through on its plan to avoid chromatography by proposing a distillation-based purification:

```
Make a solution by dissolving
C=(...)C1 (20.9 g, 128 mmol);
C[P+](...)[Br-] (57 g, 148 mmol)
in 2-methyltetrahydrofuran (320 mL)
to get Mixture 1.
```

This case study illustrates that QFANG can successfully translate abstract directives into a coherent sequence of operational steps. Its ability to first conceptualize a complex strategy—incorporating solvent selection, waste minimization, and safety considerations—and then instantiate that strategy into a concrete procedure showcases an advanced level of reasoning crucial for designing modern, sustainable synthetic routes.