
State-Aware Policy Optimization for a Reliable Multi-Turn, Multi-Tool Scientific Agent in Kinetic Biological Models

Anonymous Authors¹

Abstract

Large language model (LLM)-based agents show promise for scientific tasks requiring multi-turn reasoning with multi-tool use, but often fail due to error propagation, where early mistakes cascade through downstream steps. In kinetic biological modeling, including quantitative systems pharmacology (QSP), multi-step workflows involve simulating and analyzing complex dynamical systems. Standard reinforcement learning (RL) assumes stateless input, yet these scientific workflows require stateful simulation backends that preserve environment state across tool calls. We address this limitation through a state-aware Group Relative Policy Optimization (GRPO) framework that decomposes multi-turn interactions into per-turn episodes while reconstructing environment state by replaying prior tool calls on a live simulation backend. Combined with a hybrid reward function, consisting of deterministic verification and a frozen LLM judge, our approach enables verifiable, sample efficient optimization. We apply this to Talk2BioModels (T2B), an open-source agent to interrogate kinetic biological models through natural language. Using only a compact Qwen2.5-3B-Instruct backbone, our optimized agent achieves 98.8% tool correctness, 91.5% argument correctness, and 73.2% task completion, outperforming frontier models across on a 324-turn benchmark spanning 10 multi-turn, multi-tool scenarios derived from 20 biological models.

1. Introduction

The rapid development of large language models (LLMs) has enabled agentic AI frameworks capable of structured reasoning and iterative tool use, allowing them to solve

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. [†]Correspondence to: Anonymous Author <anon.email@domain.com>.

Submitted to the AI for Science workshop (ICML 2026).

complex, multi-step tasks (Singh et al., 2025b; Huang et al., 2025; Gao et al., 2025). In scientific domains, such systems offer the potential to act as AI scientists, providing conversational access to computational analyses. This is particularly relevant in kinetic biological modeling (*e.g.*, quantitative systems pharmacology, QSP), where researchers rely on mathematical models to describe complex biological processes. These models are often formulated as ordinary differential equation (ODE) systems that capture biochemical reaction networks, cellular signaling pathways, and pharmacokinetic–pharmacodynamic relationships. Working with such models requires a high degree of rigor and reliability when interacting with simulation tools (Saini & Farnoud, 2025; Shahin et al., 2025), including precise parameter specification, careful state management, and multi-step analysis, where errors can propagate across steps.

Current agentic systems struggle in such settings, especially over long, multi-turn interactions. While small language models (SLMs) may suffice for repetitive tasks (Belcak et al., 2025), they often exhibit incomplete execution in long-horizon, multi-turn interactions (Wan et al., 2025). Common failure modes include incorrect tool selection, argument parsing errors, insufficient tool documentation, and limited context management, all of which can propagate errors and undermine reasoning capabilities (Xiong et al., 2025; Shah et al., 2026). Addressing these limitations requires improving both reasoning and decision-making during tool use. Supervised fine-tuning (SFT) enables task-specific behavior and improves instruction following (Wei et al., 2022; Chung et al., 2024; Pareja et al., 2024; Gao et al., 2025). Reinforcement learning from human feedback (RLHF), in contrast, optimizes behavior using reward signals, often via actor–critic methods such as proximal policy optimization (PPO) (Christiano et al., 2017; Schulman et al., 2017). Alternative approaches, such as direct preference optimization (DPO), aim to avoid the instability of RL by removing explicit policy optimization (Rafailov et al., 2023; Meng et al., 2024). Despite these differences, RL-based methods remain attractive for complex scientific tasks. They have been shown to encourage emergent behaviors such as reflection, verification, and adaptation, and often improve generalization beyond training data (Guo et al., 2025; Chu et al.,

2025). Recently proposed Group Relative Policy Optimization (GRPO) further improves efficiency by eliminating the critic model while maintaining strong instruction-following performance (Shao et al., 2024).

More recent studies have explored applying GRPO to optimize tool-calling agent policies. Approaches such as Nemotron-Tool-N1 (Zhang et al., 2025) and ToolRL (Qian et al., 2025) focus on reward design for tool calls. Multi-turn extensions such as RC-GRPO (Zhong et al., 2026) and MT-GRPO (Wei et al., 2025) improve rollout diversity and per-turn credit assignment. Frameworks such as RLFac-tory (Chai et al., 2025b) and AgentFly (Wang et al., 2025) incorporate environment interaction by reconstructing state during rollouts. However, existing approaches do not fully address the challenges of persistent, stateful environments. Most methods either assume stateless interactions (Zhang et al., 2025; Qian et al., 2025), reconstruct state without isolating per-turn credit (Chai et al., 2025b; Wang et al., 2025), or decompose interactions into independent steps evaluated against static ground truth (Wei et al., 2025; Feng et al., 2025). These assumptions break down in scientific settings, where each step depends on an evolving simulator state that must be rebuilt through a sequence of tool calls.

To address this gap, we propose a framework for learning robust AI scientist policies for multi-turn, multi-tool interactions using state-aware GRPO. Our approach decomposes interactions into turn-level episodes while reconstructing the environment state before each rollout by replaying prior tool calls against a live simulation backend. This design preserves the benefits of per-turn credit assignment (Wei et al., 2025) while maintaining faithful execution of stateful processes (Chai et al., 2025b). It also reduces error accumulation compared to full-trajectory rollouts. Following the paradigm of RL with verifiable rewards (RLVR) (Lambert et al., 2024; Guo et al., 2025), we adopt a reward-model-free approach that combines deterministic tool-use verification with a frozen LLM-based judge (Zheng et al., 2023). We evaluate our framework on Talk2BioModels (T2B) (Singh et al., 2025b; Wehling et al., 2025), an open-source agent for kinetic biological modeling, using a synthetic multi-turn, multi-tool benchmark derived from its data generation pipeline. Employing a compact Qwen2.5-3B-Instruct backbone (Qwen et al., 2024), our framework achieves 98.8% tool correctness, 91.5% argument correctness, and 73.2% task completion, outperforming evaluated proprietary frontier models, including O3-mini and GPT-5.4-mini. While broadly applicable, this approach is particularly relevant for domains such as biological modeling, where reliable multi-step computational modeling is essential.

In summary, our contributions in this study are as follows:

- **State-aware Per-turn Optimization Framework:** We propose a two-stage framework combining SFT and state-aware GRPO with verifiable rewards, which formulates multi-turn interactions as per-turn episodes by reconstructing environment state and replaying prior tool calls on a live simulation backend. This enables finer-grained credit assignment in stateful scientific tool-use settings while reducing cascading errors in long-horizon rollouts.
- **Optimized AI Scientist for Biological Models:** We apply our proposed framework to T2B, improving its capabilities as an AI scientist by enhancing tool-calling consistency and accuracy for reliable interrogation of kinetic biological models.
- **AI Scientist using a Compact Model:** Using our proposed framework, a 3B-parameter AI scientist achieves strong performance gains over several frontier models in tool correctness, argument correctness, and task completion, revealing the effectiveness of the approach for improving smaller models on scientific tasks.

2. Proposed Framework

We denote a benchmark dataset comprising tool-based question-answer (QA) pairs as $\mathbb{X} = \{(\mathbf{q}_n, (\mathbf{f}_n, \mathbf{a}_n))\}_{n=1}^N$, where N is the total number of samples and $\mathbf{f} \in \{f_1, f_2, \dots, f_K\}$ indicates which of the K available tools is invoked to produce the answer. In a multi-turn setting, each n -th element unfolds as a sequence of alternating questions, tool invocations, and answers, defined as $(\mathbf{q}_n, (\mathbf{f}_n, \mathbf{a}_n)) = [q_n^1, (f_n^1, a_n^1), \dots, q_n^T, (f_n^T, a_n^T)]$, where T denotes the total number of turns. Hereafter, indices are omitted when unambiguous for de-cluttering. Our main objective is to optimize the agent policy π_θ for reliable tool calling, enabling correct answer generation within an RL-based framework, as illustrated in Fig. 1.

2.1. Synthetic Benchmark Construction

Following the self-instruct paradigm (Wang et al., 2023), a prior study (Wehling et al., 2025) introduced a synthetic benchmark construction pipeline to generate QA pairs for evaluating Talk2BioModels (T2B) in terms of task completion. We extend this pipeline to support multi-turn, multi-tool invocation using five T2B tools: (i) `get_modelinfo`, which retrieves descriptive metadata of a biological model; (ii) `simulate_model`, which performs time-course simulation of a biological model; (iii) `steady_state`, which computes steady-state concentrations of model species; (iv) `ask_question`, which queries simulation and steady-state results; and (v) `parameter_scan`, which explores how parameter variations in the model affect concentrations of selected species. We select these tools due to their complex argument structures and design multi-turn scenarios

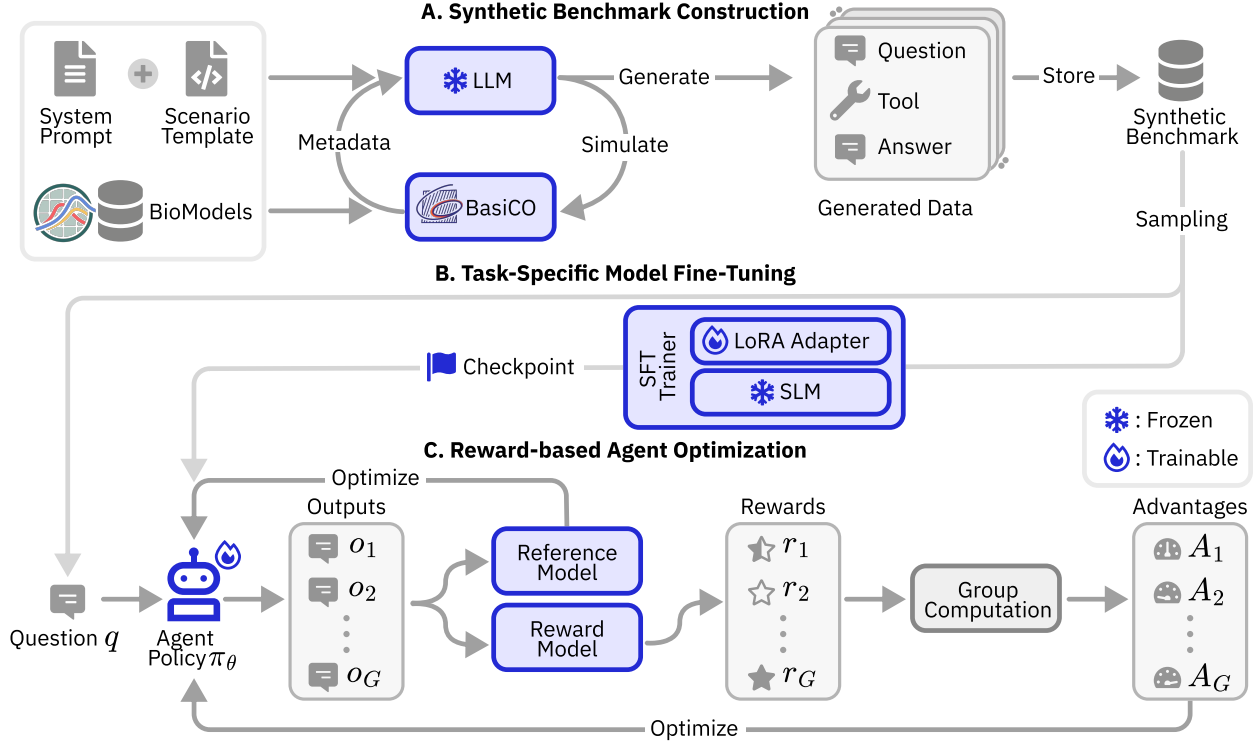


Figure 1. Overview of the proposed framework, including (a) synthetic benchmark construction and (b–c) staged SFT and state-aware GRPO optimization with verifiable rewards for multi-turn, multi-tool calling agent training.

involving their combined usage (see Appendix A). For each multi-turn scenario S , we set a system prompt ϕ (see Appendix A.2) for a frozen LLM π_{QA} to produce a question q together with its answer a , guided by an expected tool f that is iteratively specified after each turn in a scenario $s \in S$ until a predefined number N_s of samples is reached, as formalized as

$$\{(q_n, (f_n, a_n))\}_{n=1}^{N_s} \leftarrow \pi_{QA}(\phi, s), \quad \forall s \in S, N_s \leq N.$$

2.2. Task-Specific Model Fine-tuning

Prior to reward-based optimization, we first ensure that the base model acquires the capabilities required for the target tasks, such as reading simulated time-series data in dataframe format. We apply supervised fine-tuning (SFT) with low-rank adaptation (LoRA) (Hu et al., 2022) on the base model, optimizing the objective $\mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(q,a) \sim \mathbb{X}_{\text{SFT}}} \sum_{t=1}^{|a|} \log \pi_\theta(a_t | q, a_{<t})$, where \mathbb{X}_{SFT} denotes the SFT dataset, q is the prompt, and a is the sequence of answer tokens, including tool calls and final responses.

2.3. Reward-based Agent Optimization

We adopt GRPO (Shao et al., 2024) to optimize the agent’s policy π_θ using $\mathcal{L}_{\text{GRPO}} = \mathbb{E}_{q \sim \mathbb{X}, \mathbb{O} \sim \pi_{\theta_{old}}} \left(\frac{1}{G} \sum_{g=1}^G \min(\rho_g(\theta) A_g, [\rho_g(\theta)]_{1-\epsilon}^{1+\epsilon} A_g) - \beta \mathbb{D}_{KL}(\pi_\theta \| \pi_{ref}) \right)$, where $\mathbb{O} = \{o_g\}_{g=1}^G$ denotes a group of

G sampled outputs produced by the agent policy for prompt q , and $\rho_g(\theta)$ is the standard importance ratio between the current and behavior policies. On top of GRPO, we further consider two key factors to improve optimization effectiveness: verifiable reward design and state-aware per-turn policy optimization.

Verifiable Reward Design Following evidence that learned reward models are prone to overoptimization and reward hacking in structured settings (Gao et al., 2023; Guo et al., 2025), we adopt a *reward-model-free* design, consistent with the RLVR paradigm (Lambert et al., 2024). This design is motivated by three factors. First, tool invocation correctness is directly verifiable via rule-based checks, removing the need for a learned proxy. Second, prior work shows that rule-based rewards can effectively train strong tool-calling agents (Zhang et al., 2025; Qian et al., 2025; Zhong et al., 2026). Third, cascaded gating ensures that argument and task completion rewards are only computed when tool execution is correct, preventing contamination from invalid tool calls.

In particular, each sampled output is scored using a composite reward r_g combining deterministic heuristics and a frozen LLM judge as

$$r_g = \lambda_{\text{tool}} r_g^{\text{tool}} + \lambda_{\text{arg}} r_g^{\text{arg}} + \lambda_{\text{task}} r_g^{\text{task}}, \quad (1)$$

where r_g^{tool} is a deterministic tool-sequence correctness score, r_g^{arg} is a deterministic argument-schema correctness score with soft numerical tolerance, and r_g^{task} is an LLM-as-a-judge (Zheng et al., 2023) score for task completion using a frozen, pre-specified LLM, with λ_* denoting the corresponding weighting coefficients. The group-relative advantages are then computed as $A_g = (r_g - \mu(\mathbf{r}))/\sigma(\mathbf{r})$, where $\mu(\mathbf{r})$ and $\sigma(\mathbf{r})$ denote the mean and standard deviation of rewards within the group, respectively.

State-Aware Per-Turn Policy Optimization In most multi-turn scientific workflows, later turns strongly depend on the state produced in earlier ones. In T2B as an AI scientist setting, users are able to query simulated data from biological models, where the correctness of each response depends on whether the underlying simulation and its outputs are valid. Standard GRPO treats each prompt as stateless, which is incompatible with such dependencies. Existing multi-turn GRPO methods either (i) use full-trajectory rollouts with live execution (Zhong et al., 2026), which accumulate errors and provide only trajectory-level credit, (ii) decompose into per-turn episodes but evaluate against static ground truth without live execution (Zhang et al., 2025; Qian et al., 2025; Wei et al., 2025; Feng et al., 2025; Chai et al., 2025a), which fails in stateful environments where rewards depend on simulator outputs, or (iii) reconstruct state during rollout (Chai et al., 2025b; Wang et al., 2025), which entangles state reconstruction with generation and prevents clean per-turn credit assignment.

Compared to full-trajectory rollouts, per-turn decomposition avoids cascading errors and enables fine-grained credit assignment without stepwise estimators. Unlike static scoring, live execution on reconstructed state grounds rewards in simulator outputs, while replaying ground-truth tool calls ensures consistent state independent of policy errors, separating state construction from evaluation. An ideal approach should also preserve GRPO’s efficiency by requiring only G rollouts per turn rather than full multi-turn trajectories.

Hence, we propose a per-turn episode formulation with stateful context reconstruction, as shown in Algorithm 1. Specifically, we decompose each multi-turn conversation into turn-level GRPO episodes and initialize a fresh environment for each turn. The environment reconstructs the current state by replaying state-building tool calls from the ground-truth conversation history. These calls are then re-executed against the live simulation backend, ensuring that the state is faithfully restored for the specific turn. We treat this process as a pre-warming phase where intermediate scoring traces are discarded after state reconstruction, so that only newly generated policy rollouts contribute to the reward signal. Finally, the policy generates G rollouts for the current turn, which are evaluated in the reconstructed live environment and optimized accordingly.

Algorithm 1 State-Aware Per-Turn GRPO with verifiable rewards

Require: Conversation $\mathbb{X}_i = [(q_1, f_1, a_1), \dots, (q_T, f_T, a_T)]$

- 1: **for** each turn $t = 1, \dots, T$ **do**
- 2: $\mathcal{E}_t \leftarrow \text{NEWENVIRONMENT}(\pi_\theta)$
- 3: $\mathcal{P}_t \leftarrow \{(f_j, \mathbf{a}_j) \mid j < t, f_j \in \mathcal{F}_{\text{state}}\}$ {Prior state-building calls}
- 4: **for** each $(f_j, \mathbf{a}_j) \in \mathcal{P}_t$ **do**
- 5: $\mathcal{E}_t.\text{EXECUTE}(f_j, \mathbf{a}_j | \pi_\theta)$ {Reconstruct state}
- 6: **end for**
- 7: $\mathcal{E}_t.\text{CLEARSCORINGARRAYS}()$ {Only policy calls count}
- 8: $\mathbf{p}_t \leftarrow \text{BUILD PROMPT}(C_{<t}, q_t)$ {Full history + current query}
- 9: $\mathcal{O}_t = \{o_g\}_{g=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | \mathbf{p}_t, \mathcal{E}_t)$ { G rollouts with live tools}
- 10: $r_g \leftarrow \lambda_{\text{tool}} r_g^{\text{tool}} + \lambda_{\text{arg}} r_g^{\text{arg}} + \lambda_{\text{task}} r_g^{\text{task}}$
- 11: $A_g \leftarrow (r_g - \mu(\mathbf{r}))/\sigma(\mathbf{r})$
- 12: Update π_θ using $\mathcal{L}_{\text{GRPO}}$ with $\{(o_g, A_g)\}_{g=1}^G$
- 13: **end for**

3. Results

3.1. Experimental Setup

Dataset Using our synthetic benchmark pipeline, we construct 10 scenarios (Table 5) derived from 20 ODE-based human biological models (Table 3) encoded in the Systems Biology Markup Language (SBML) covering various biological settings. Each model defines species (quantities of entities such as molecules or populations), reactions (rules describing how species are transformed or interact), and parameters (numerical values controlling reaction rates and system behavior). For each scenario, we first use BasiCO (Bergmann, 2023) to run simulations¹ and generate scientifically grounded model behaviors. These simulations are then used to guide a frozen LLM (GPT-5-nano (Singh et al., 2025a)) in generating multi-turn question-answer trajectories. This process yields 1,200 conversations, which are split in a scenario-stratified manner into 960/120/120 train, validation, and test samples. The test set contains 324 evaluation turns. Additional dataset statistics are provided in Appendix A.3.

Implementation Details We extend T2B, originally built on LangGraph (LangChain, Inc., 2024), to support policy optimization using the Transformer Reinforcement Learning (TRL) library (von Werra et al., 2020). We use Qwen2.5-3B-Instruct (Qwen et al., 2024) as the base model and train it with LoRA adapters in two-stages: (i) SFT on scenario-

¹All floating-point arguments derived from the simulated data are rounded to six significant figures.

Table 1. Quantitative performance on the held-out test split, comparing our optimized models with several OpenAI frontier models. Results for Qwen-SFT and Qwen-SFT+GRPO are based on a single training seed; tool and argument correctness are deterministic (no variance across runs), while task completion may vary slightly due to the frozen LLM judge.

Model	Tool Correctness	Argument Correctness	Task Completion	Latency (s/turn)
GPT-5 mini	0.407	0.446	0.355	28.0
GPT-5.4 nano	0.469	0.423	0.364	11.2
GPT-5 nano	0.556	0.441	0.491	41.1
GPT-4o mini	0.654	0.477	0.543	13.1
O4 mini	0.694	0.481	0.623	19.9
GPT-5.4 mini	0.824	0.503	0.669	8.2
O3 mini	0.843	0.500	0.703	17.0
Qwen2.5-3B (base)	0.682	0.248	0.524	6.9
Qwen-SFT (ours)	0.960	0.856	0.729	8.8
Qwen-SFT+GRPO (ours)	0.988	0.915	0.732	10.8

stratified multi-turn traces, followed by (ii) state-aware GRPO with verifiable rewards. During RL training, we incorporate state-aware online scoring by executing tool calls and reconstructing turn-level context at each step. We compare against seven proprietary frontier models, including GPT-4o mini, GPT-5 mini, GPT-5 nano, GPT-5.4 mini, GPT-5.4 nano, O4 mini, and O3 mini (Singh et al., 2025a; OpenAI, 2025). All models are evaluated under identical conditions using the same stateful execution harness and LLM-based judge. Experiments are conducted on a system with four NVIDIA H100 80GB GPUs. Additional implementation details are provided in Appendix B.

3.2. Performance Comparison

Table 1 summarizes quantitative performance on the held-out test split (324 turns, 120 conversations) across all evaluated models. Our fully optimized T2B agentic AI scientist (Qwen-SFT+GRPO) achieves the best results on all three metrics, outperforming all evaluated proprietary frontier models. Compared to the strongest baseline (*i.e.*, O3-mini), Qwen-SFT+GRPO improves tool correctness by +14.5 percentage points (pp), argument correctness by +41.5 pp, and task completion by +2.9 pp. The largest gains are observed in argument correctness: all proprietary models remain below 51%, whereas our 3B-parameter model reaches 91.5%, indicating that domain-specific fine-tuning with verifiable rewards is particularly effective for structured reasoning. Figure 2 further shows the composite score–latency Pareto frontier, where Qwen-SFT+GRPO achieves the highest overall performance while remaining competitive in latency.

Beyond the reported quantitative results, we analyzed turn-level multi-turn conversations (Figure 3), where baseline models degraded at later turns due to difficulty handling state-dependent interactions requiring prior context, while the GRPO-optimized model remained stable. Based on these findings, we conducted further analyses along two dimensions, including conversation rate and cascading fail-

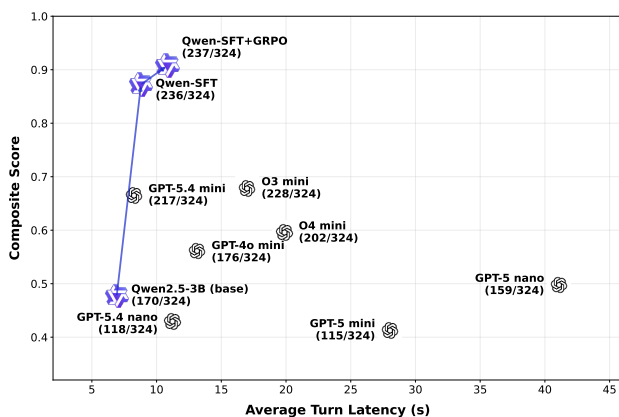


Figure 2. Pareto frontier of composite score vs. average turn latency. Numbers in parentheses denote successfully scored turns (out of 324). Our Qwen-SFT+GRPO model lies on the upper-left frontier, achieving the highest composite score at competitive latency. The arrow shows the Base→SFT→SFT+GRPO trajectory.

ures in stateful interactions. Specifically, we define a perfect conversation as a scenario in which the model correctly performed tool parsing and invocation at every turn and produced the correct final answer for the task. Under this criterion, Qwen-SFT+GRPO achieved 78.3% perfect conversations, compared to 20.8% for the strongest proprietary model (*i.e.*, O3-mini), corresponding to a 3.8× improvement. We found that the primary bottleneck arose in cases where the correct tool was selected but the arguments were incorrect. In these cases, O3-mini produced incorrect arguments on 33.0% of turns, compared to 6.8% for our SFT+GRPO-optimized model. This disparity suggests that tool selection and argument construction are distinct capabilities governed by different failure modes, and that argument construction benefits from domain-specific optimization.

On the other hand, argument errors are compounded in multi-turn scenarios. While the metrics tool and argument correctness are deterministic, task completion is assessed by a frozen LLM judge comparing the model’s textual response to the ground truth (GT) answer. We present three representative cases where proprietary models selected the correct tool but produced incorrect arguments, directly illustrating the 91.5% vs. <51% argument correctness gap (Table 2). In **Case A** (Scenario 4), the simulated user queried an IL-6 signaling model (Dwivedi et al., 2014) via the following sequence of tool-calls: `get_modelinfo` → `simulate_model` + `ask_question` → `steady_state` → `ask_question`. At Turn 3 in particular, the user requested a steady-state computation with three species (*i.e.*, IL6, `sgp130`, `sR_IL6_sgp130`). Here, GPT-5.4-mini appended compartment-annotation suffixes to species identifiers (*i.e.*, `IL6{serum}`) instead of IL6), introducing a semantic error that conflated the model’s display notation with the tool’s input namespace, and

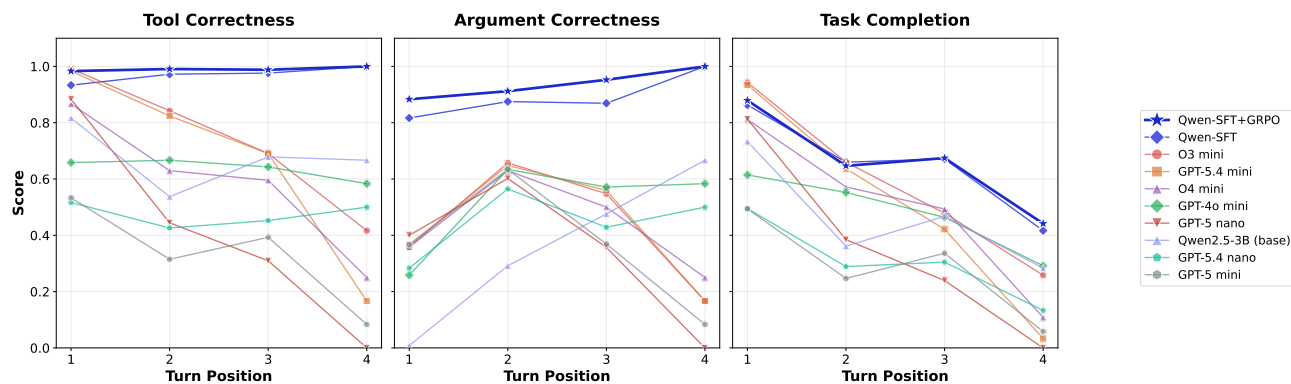


Figure 3. Performance by turn position across all evaluated models. Qwen-SFT+GRPO (dark blue, thick line) maintains near-perfect tool and argument correctness across all turn positions, while proprietary models degrade at later turns where state-dependent interactions dominate. All models show declining task completion at Turn 4, but GRPO degrades least. This robustness at later positions directly reflects the benefit of state-aware training with prior-turn replay.

omitted time data. Meanwhile, O3-mini correctly handled species names but still omitted time data and fabricated the experiment name. Consequently, both argument errors led to downstream failure, as neither model invoked the follow-up `ask_question` tool at Turn 4 (*i.e.*, both models produced free-form text responses instead of calling the expected tool), resulting in zero scores across all three metrics. In comparison, our SFT+GRPO-optimized model achieved correct tool selection and argument accuracy across all turns, but lower task completion scores at Turns 2 and 4. These correspond to `ask_question` steps, where the model was required to summarize raw dataframe outputs into natural-language responses, such that the LLM judge penalized imprecise phrasing even when the underlying tool execution was correct. We discuss this aspect in a dedicated subsection on task completion.

Meanwhile, **Case B** (Scenario 6) followed `simulate_model` \rightarrow `parameter_scan` \rightarrow `ask_question`, involving NFAT activation dynamics simulation (Fisher et al., 2006), a parameter scan over `k7`, and a query over the results. At the second turn, the user requested a parameter scan with “3000s duration and 10s interval”. Both O3-mini and GPT-5.4-mini set the interval to 300 instead, resulting in temporally undersampled scan results. This error subsequently cascaded to the next turn, where both models referenced fabricated experiment names instead of the canonical identifier. In contrast, our optimized model correctly computed the interval and consistently used the correct experiment namespace across turns.

Finally, **Case C** (Scenario 8) involved the sequence `steady_state` \rightarrow `ask_question`, examining steady-state analysis followed by a parameter scan on a calcium oscillation model (Fisher et al., 2006). At the first turn, O3-mini placed the time data field outside the required argument structure instead of nesting it within the schema and set the

duration to 100. Meanwhile, GPT-5.4-mini omitted both the time data field and the required boolean SBML flag. Our optimized model produced correct schema nesting, duration values, and experiment naming across both turns. These results indicate argument construction errors rather than tool-selection failures, suggesting that argument formation is the primary bottleneck in proprietary models. By incorporating domain-specific optimization into our framework, we effectively mitigate this issue. We provide a detailed per-tool breakdown in Appendix D.

4. Discussion

Complementary Roles of SFT and GRPO The ablation results (Table 6) show that SFT and GRPO play complementary roles. SFT yields the largest improvement, increasing argument correctness from 24.8% to 85.6% (+60.8 pp) by learning structured tool-call formats and domain-specific argument schemas. GRPO then further refines the policy, improving tool correctness from 96.0% to 98.8% and argument correctness from 85.6% to 91.5%. The per-tool analysis (Table 7) shows that GRPO’s gains are concentrated on tools with complex argument structures. For example, argument correctness of the `steady_state` tool improves from 77.8% to 100% (+22.2 pp), while tool correctness of the `parameter_scan` tool from 81.2% to 93.8% (+12.6 pp). We also conducted a turn-position analysis (Figure 3), showing that GRPO maintains near-perfect performance at later turns where state-dependent interactions dominate, while baseline models degrade significantly, consistent with the benefit of prior-turn state reconstruction.

Understanding Failure Modes in Tool-Calling Agents

The largest performance difference is observed in argument correctness, where all evaluated proprietary models remain below 51%, whereas our 3B-parameter SFT+GRPO model

State-Aware Policy Optimization for a Reliable Multi-Turn, Multi-Tool Scientific Agent in Kinetic Biological Models

Table 2. Three representative multi-turn argument failure cases. Each block shows a scenario with the user query, ground-truth tool and arguments, and per-turn model behavior. Correct outputs are marked in blue, and incorrect or missing outputs in red. Free-form responses without tool invocation are labeled “skipped.” We include O3-mini and GPT-5.4 as the next-best models based on quantitative analysis.

Case A (Scenario 4): IL-6 signaling (Dwivedi et al., 2014)				
	Turn 1	Turn 2	Turn 3	Turn 4
Query	“What parameters are defined for BIOMD537?”	“Simulate 2016 h, 14 steps, IL6{serum}=2.63e-4. What is CRP{serum}?”	“Please run the steady state... IL6=2.99e-4, sgp130=2.748, sR_IL6_sgp130=0.0605.”	“Concentrations and transition times for the species.”
Ground Truth	get_modelinfo; parameters=true, all others false	simulate+ask_question; duration=2016, interval=144; species=[“IL6”]; experiment=“simulate.BIOMD...537”	steady_state; species=[“IL6”, “sgp130”, “sR_IL6_sgp130”]; time={1000,100}; experiment=“ss.BIOMD...537”	ask_question; experiment=“ss.BIOMD...537”; question_context=“steady_state”
Qwen-SFT+GRPO	all arguments passed correctly	match ground truth	match ground truth	match ground truth
O3-mini	only true fields; 5 false omitted	4 calls (+2 extra)	time.data omitted; exp fabricated	skipped
GPT-5.4	only true fields; no use_uploaded_sbml	simulate only; species=[“IL6{serum}”]	species with {serum} suffix; time.data omitted	skipped
Case B (Scenario 6): NFAT activation (Fisher et al., 2006)				
	Turn 1	Turn 2	Turn 3	
Query	“Simulate BIOMD123, 3000 s, 5 steps, Calcineurin nuc.=2.59e-5 μmol.”	“Please run a parameter scan on the model... Ca Cyt=1.174, NFAT_Nuc=5.37e-4. k7: 0.00821→0.01321, step 5e-4.”	“Concentration of Phospho. NFAT in nucleus?”	
Ground Truth	simulate; duration=3000, int=600; species=[“Active Calcineurin in nucleus”]	param_scan; duration=3000, interval=10; parameter=k7, 11 values; experiment=“param_scan.BIOMD...123”	ask_question; experiment=“simulate.BIOMD...123”; question_context=“simulation”	
Qwen-SFT+GRPO	match ground truth	interval=10; k7 match	match ground truth	
O3-mini	match ground truth	interval=300 (30×); experiment fabricated	experiment fabricated (cascaded)	
GPT-5.4	match ground truth	interval=300; exp=“...k7_scan”	experiment cascaded	
Case C (Scenario 8): Calcium-dependent NFAT dynamics (Fisher et al., 2006)				
	Turn 1	Turn 2		
Query	“Please run the steady state simulation on BIOMD122... Calcineurin cyt=1.35e-5, NFAT_Calc=3.71e-6, Phospho. NFAT=3.03e-4 μmol.”	“Please run a parameter scan on the model by updating the initial... k9 4.334→5.334, step 0.1; NFAT_nuc=4.99e-4, NFAT_Cyt=1.14e-4.”		
Ground Truth	steady_state; time={1000,100} in arg_data; species=[“Active Calc...”, “NFAT Calc...”, “Phospho. NFAT...”]	param_scan; duration=400, int=16; parameter=k9; experiment=“param_scan.BIOMD...122”		
Qwen-SFT+GRPO	match ground truth	match ground truth		
O3-mini	time.data at root (not in arg_data); duration=100	match ground truth		
GPT-5.4	time.data omitted; use_uploaded_sbml omitted	match ground truth		

reaches 91.5%. Case studies in Table 2 indicate that this gap is associated with consistent failure patterns rather than random errors. Such failure patterns include namespace confusion (e.g., appending compartment suffixes to species identifiers in biological models), schema nesting errors (placing fields at incorrect levels of the argument hierarchy), and numerical reasoning failures (such as miscalculating intervals from duration and step count). These patterns suggest that general-purpose models lack structured knowledge of tool schemas, which can be mitigated through domain-specific policy optimization. In contrast, the results of the tool selection suggest that this aspect is comparatively easier. For example, O3-mini achieves 84.3% tool correctness without domain-specific training, but only 50.0% argument correctness. This disparity indicates that tool selection and argument construction are governed by different failure modes, and that argument construction is the primary bottleneck where targeted optimization yields the largest gains.

Task Completion under Mixed Verifiable and LLM-Based Rewards Despite large gains in tool and argument correctness, our optimized agent shows only marginal improvement in task completion, with a +0.3 pp increase from SFT to SFT+GRPO, reflecting a fundamental distinction in what each metric captures. In particular, tool and argument correctness are optimized using deterministic, verifiable rewards, which account for 80% of the total reward signal. In contrast, task completion is evaluated by a frozen LLM-based judge on open-ended responses and contributes the remaining 20%. This weighting reflects a deliberate design choice: in scientific tool-use settings, correct tool execution is a prerequisite for meaningful downstream results, and our framework successfully optimizes this objective. However, the smaller reward weight provides limited signal for improving natural-language summarization, resulting in only marginal gains. Nevertheless, our model still achieves the highest task completion among all evaluated models (73.2% vs. 70.3% for O3-mini), suggesting that accurate tool use

385 provides a strong foundation for answer quality.

386 We also conducted a turn-level performance analysis by
387 breaking it down by scenario (Figure 8), revealing two distinct
388 sources of task completion degradation. First, scenarios
389 involving individual `ask_question` turns (e.g., Turn 3
390 in Scenarios 1, 3 and 9; Turn 2 in Scenario 2) exhibit
391 reduced task completion in all models evaluated, including
392 our optimized agent. This tool requires a downstream
393 LLM to parse large tabular data and extract precise numerical
394 responses, a process that becomes unreliable when the
395 data exceed the effective context capacity of the model.
396 Second, prompts invoking multi-tool combination turns
397 (e.g., `simulate_model` followed by `ask_question`
398 in Scenario 4 Turn 2, and `steady_state` followed by
399 `ask_question` in Scenarios 7 and 10 Turn 2) exhibit
400 more severe degradation, as the agent must both execute a
401 state-building tool and immediately query its results within
402 a single turn. Since our per-turn decomposition treats each
403 turn as an independent episode, it does not explicitly train
404 the agent on intra-turn tool composition, which may compound
405 the tabular data parsing challenge. These observations
406 suggest two correlative directions for future work: (i)
407 replacing the current `ask_question` tool with a dedicated
408 data-analysis agent or a more capable LLM backend for handling
409 large tabular outputs, and (ii) developing finer-grained
410 decomposition strategies that explicitly optimize intra-turn
411 multi-tool orchestration and execution.

412
413
414 **Limitations** Despite encouraging results, several limitations
415 remain. First, performance varies across scenarios. For
416 example, in Scenario 5 (Appendix E), a single-turn task (invoking
417 the `parameter_scan` tool), our model achieved
418 83.3% tool correctness and 58.3% argument correctness,
419 compared to 100% and 66.7% for O3-mini. Analysis shows
420 that the model often omits a required field (e.g., the argument
421 of experiment name). Because this omission is not penalized
422 by the current argument-level reward, it reveals a gap in the
423 reward design: unverified fields can still affect downstream
424 behavior. Second, our evaluation includes only OpenAI frontier
425 models due to cost and time constraints. Expanding the set of
426 baselines would provide a more comprehensive comparison.
427 Third, performance may improve further with models pretrained
428 on biomedical corpora (Labrak et al., 2024; Bolton et al., 2024),
429 which offer stronger domain-specific priors. Fourth, the synthetic
430 benchmark, while grounded in BasicO simulations, may reflect
431 biases introduced by the generator LLM and may not capture
432 the full diversity of expert queries. While we provide a preliminary
433 out-of-distribution evaluation in Appendix F, showing that the
434 optimized agent retains meaningful performance on unseen tool
435 compositions and novel tool interfaces, performance in this
436 controlled setting does not guarantee robustness under ambiguous
437 or underspecified
438
439

inputs. Fifth, our current proposed framework focuses on a single-agent setting. Extending it to multi-agent systems, such as those combining simulation and knowledge-graph agents (Singh et al., 2025b), is a promising direction for future work, particularly in the context of hierarchical reward optimization (Geng et al., 2024; 2025).

Finally, a key limitation of our work is that it assumes biological models are already fully specified, with predefined species, parameters, and initial conditions. In practice, however, models must often be initialized from unstructured literature, requiring the extraction and grounding of entities and interventions from free-text sources. This upstream step remains a major bottleneck (Tiwari et al., 2021; Karr et al., 2022) and is not addressed by our framework. Future work could integrate agentic extraction and mapping pipelines to convert literature-derived information into structured, simulation-ready model configurations, enabling a more complete end-to-end workflow.

5. Conclusion

In this study, we propose a state-aware policy optimization framework for improving tool-calling behavior in AI scientists. We apply this approach to the open-source Talk2BioModels (T2B) agent, enabling reliable tool selection and argument construction when interacting with kinetic biological models. Our method introduces state-aware GRPO, which decomposes multi-turn interactions into per-turn optimization steps while reconstructing environment state through prior-turn replay. This design enables fine-grained credit assignment in stateful scientific workflows. Combined with verifiable rewards, our approach trains a compact Qwen2.5-3B-Instruct model, achieving 98.8% tool correctness, 91.5% argument correctness, and 73.2% task completion. The model outperforms several frontier baselines while maintaining competitive latency, demonstrating that targeted policy optimization can yield strong gains even at small scale. A preliminary out-of-distribution evaluation further shows that the optimized policy generalizes to unseen tool compositions and novel tool interfaces to some extent. While our current experiments focus on synthetic scenarios, future work will extend this framework to real-user queries, improve task completion by replacing the current `ask_question` tool with a dedicated data-analysis agent, develop finer-grained decomposition strategies for intra-turn multi-tool coordination, and investigate hierarchical optimization in multi-agent settings.

References

Belcak, P., Heinrich, G., Diao, S., Fu, Y., Dong, X., Muralidharan, S., Lin, Y. C., and Molchanov, P. Small Language Models are the Future of Agentic AI. *arXiv preprint*

- 440 *arXiv:2506.02153*, 2025.
- 441
- 442 Bergmann, F. T. BASICO: A simplified Python interface to
443 COPASI. *Journal of Open Source Software*, 8(90):5553,
444 2023.
- 445 Bolton, E., Venigalla, A., Yasunaga, M., Hall, D., Xiong, B.,
446 Lee, T., Daneshjou, R., Frankle, J., Liang, P., Carbin,
447 M., et al. BioMedLM: A 2.7B Parameter Language
448 Model Trained On Biomedical Text. *arXiv preprint*
449 *arXiv:2403.18421*, 2024.
- 450
- 451 Chai, H., Cao, Z., Ran, M., Yang, Y., Lin, J., Peng, X.,
452 Wang, H., Ding, R., Wan, Z., Wen, M., Liu, W., Zhang,
453 W., Huang, F., and Wen, Y. PARL-MT: Learning to
454 Call Functions in Multi-Turn Conversation with Progress
455 Awareness. *arXiv preprint arXiv:2509.23206*, 2025a.
- 456
- 457 Chai, J., Yin, G., Xu, Z., Yue, C., Jia, Y., Xia, S., Wang, X.,
458 Jiang, J., Li, X., Dong, C., He, H., and Lin, W. RLFactory:
459 A Plug-and-Play Reinforcement Learning Post-Training
460 Framework for LLM Multi-Turn Tool-Use. *arXiv preprint*
461 *arXiv:2509.06980*, 2025b.
- 462
- 463 Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg,
464 S., and Amodei, D. Deep Reinforcement Learning from
465 Human Preferences. *Advances in neural information*
466 *processing systems*, 30, 2017.
- 467
- 468 Chu, T., Zhai, Y., Yang, J., and Ma, Y. SFT Memorizes, RL
469 Generalizes: A Comparative Study of Foundation Model
470 Post-training. 267:10818–10838, 2025.
- 471
- 472 Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fe-
473 dus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., et al.
474 Scaling instruction-finetuned language models. *Journal*
475 *of Machine Learning Research*, 25(70):1–53, 2024.
- 476
- 477 Dwivedi, G., Fitz, L., Hegen, M., Martin, S., Harrold, J.,
478 Heatherington, A., and Li, C. A Multiscale Model of
479 Interleukin-6-Mediated Immune Regulation in Crohn’s
480 Disease and Its Application in Drug Discovery and De-
481 velopment. *CPT: Pharmacometrics & Systems Pharma-*
482 *cology*, 3(1):1–9, 2014.
- 483
- 484 Feng, L., Xue, Z., Liu, T., and An, B. Group-in-Group Pol-
485 icy Optimization for LLM Agent Training. In *Advances*
486 *in Neural Information Processing Systems*, 2025. *arXiv*
487 *preprint arXiv:2505.10978*.
- 488
- 489 Fisher, W. G., Yang, P.-C., Medikonduri, R. K., and Jafri,
490 M. S. NFAT and NF κ B activation in T lymphocytes: a
491 model of differential activation of gene expression. *An-*
492 *nals of biomedical engineering*, 34(11):1712–1728, 2006.
- 493
- 494 Gao, L., Schulman, J., and Hilton, J. Scaling Laws for Re-
ward Model Overoptimization. 202:10835–10866, 2023.
- Gao, S., Zhu, R., Kong, Z., Noori, A., Su, X., Ginder, C.,
Tsiligkaridis, T., and Zitnik, M. TxAgent: An AI agent
for therapeutic reasoning across a universe of tools. *arXiv*
preprint arXiv:2503.10970, 2025.
- Geng, M., Pateria, S., Subagdja, B., and Tan, A.-H. Hi-
SOMA: A hierarchical multi-agent model integrating
self-organizing neural networks with multi-agent deep re-
inforcement learning. *Expert Systems with Applications*,
252:124117, 2024.
- Geng, M., Pateria, S., Subagdja, B., Li, L., Zhao, X., and
Tan, A.-H. L2M2: A Hierarchical Framework Integrating
Large Language Model and Multi-agent Reinforcement
Learning. In *International Joint Conference on Artificial*
Intelligence, 2025.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R.,
Zhu, Q., Ma, S., Wang, P., Bi, X., et al. DeepSeek-R1:
Incentivizing Reasoning Capability in LLMs via Rein-
forcement Learning. *arXiv preprint arXiv:2501.12948*,
2025.
- Hoops, S., Sahle, S., Gauges, R., Lee, C., Pahle, J., Simus,
N., Singhal, M., Xu, L., Mendes, P., and Kummer, U. CO-
PASI — a COMplex PATHway SIMulator. *Bioinformatics*,
22(24):3067–3074, 2006. doi: 10.1093/bioinformatics/
btl485.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang,
S., Wang, L., and Chen, W. LoRA: Low-Rank Adaptation
of Large Language Models. In *International Conference*
on Learning Representations, 2022.
- Huang, K., Zhang, S., Wang, H., Qu, Y., Lu, Y., Roohani,
Y., Li, R., Qiu, L., Zhang, J., Di, Y., et al. Biomni:
A General-Purpose Biomedical AI Agent. *bioRxiv*, pp.
2025–05, 2025.
- Karr, J. R., Malik-Sheriff, R. S., Osborne, J. M., González-
Parra, G., Forgoon, E., Bowness, R., Liu, Y., Thompson,
R. N., Garira, W., Barhak, J., Rice, J. R., Torres, M.,
Dobrovolny, H. M., Tang, T., Waites, W., Glazier, J. A.,
Faeder, J. R., and Kulesza, A. Model integration in
computational biology: The role of reproducibility,
credibility and utility. *Frontiers in Systems Biology*, 2,
March 2022. doi: 10.3389/fsysb.2022.822606. URL
[https://www.frontiersin.org/articles/
10.3389/fsysb.2022.822606/pdf](https://www.frontiersin.org/articles/10.3389/fsysb.2022.822606/pdf).
- Labrak, Y., Bazoge, A., Morin, E., Gourraud, P.-A., Rouvier,
M., and Dufour, R. BioMistral: A Collection of Open-
Source Pretrained Large Language Models for Medical
Domains. *arXiv preprint arXiv:2402.10373*, 2024.
- Lambert, N., Morrison, J., Pyatkin, V., Huang, S., Ivison,
H., Brahman, F., Miranda, L. J. V., Liu, A., Dziri, N., Lyu,

- 495 S., et al. TULU 3: Pushing Frontiers in Open Language
496 Model Post-Training. *arXiv preprint arXiv:2411.15124*,
497 2024.
- 498
499 LangChain, Inc. LangGraph: Build Resilient Language
500 Agents as Graphs, 2024. URL <https://github.com/langchain-ai/langgraph>.
- 501
502 Malik-Sheriff, R. S., Glont, M., Nguyen, T. V., Tiwari, K.,
503 Roberts, M. G., Xavier, A., Vu, M. T., Men, J., Maire,
504 M., Kananathan, S., et al. Biomodels—15 years of shar-
505 ing computational models in life science. *Nucleic acids*
506 *research*, 48(D1):D407–D415, 2020.
- 507
508 Meng, Y., Xia, M., and Chen, D. SimPO: Simple Preference
509 Optimization with a Reference-Free Reward. *Advances*
510 *in Neural Information Processing Systems*, 37:124198–
511 124235, 2024.
- 512
513 OpenAI. OpenAI o3 and OpenAI o4-mini Sys-
514 tem Card. Technical report, OpenAI, 4 2025.
515 URL [https://cdn.openai.com/pdf/](https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf)
516 [2221c875-02dc-4789-800b-e7758f3722c1/](https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf)
517 [o3-and-o4-mini-system-card.pdf](https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf). System
518 card.
- 519
520 Pareja, A., Nayak, N. S., Wang, H., Killamsetty, K., Su-
521 dalairaj, S., Zhao, W., Han, S., Bhandwadar, A., Xu,
522 G., Xu, K., et al. Unveiling the secret recipe: A guide
523 for supervised fine-tuning small LLMs. *arXiv preprint*
524 *arXiv:2412.13337*, 2024.
- 525
526 Qian, C., Acikgoz, E. C., He, Q., Wang, H., Chen,
527 X., Hakkani-Tür, D., Tur, G., and Ji, H. ToolRL:
528 Reward is All Tool Learning Needs. *arXiv preprint*
529 *arXiv:2504.13958*, 2025.
- 530
531 Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng,
532 B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin,
533 H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou,
534 J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu,
535 L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin,
536 R., Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y.,
537 Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z.,
538 and Qiu, Z. Qwen2.5 Technical Report. *arXiv preprint*
539 *arXiv:2412.15115*, 2024.
- 540
541 Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D.,
542 Ermon, S., and Finn, C. Direct Preference Optimiza-
543 tion: Your Language Model is Secretly a Reward Model.
544 *Advances in neural information processing systems*, 36:
545 53728–53741, 2023.
- 546
547 Saini, A. and Farnoud, A. QSP-Copilot: An AI-Augmented
548 Platform for Accelerating Quantitative Systems Pharma-
549 cology Model Development. *CPT: Pharmacometrics &*
Systems Pharmacology, 2025. doi: 10.1002/psp4.70127.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and
Klimov, O. Proximal Policy Optimization Algorithms.
arXiv preprint arXiv:1707.06347, 2017.
- Shah, M. B., Morovati, M. M., Rahman, M. M., and Khomh,
F. Characterizing Faults in Agentic AI: A Taxonomy
of Types, Symptoms, and Root Causes. *arXiv preprint*
arXiv:2603.06847, 2026.
- Shahin, M. H., Goswami, S., Lobentanzer, S., and Corrigan,
B. W. Agents for change: Artificial intelligent workflows
for quantitative clinical pharmacology and translational
sciences. *Clinical and Translational Science*, 2025. doi:
10.1111/cts.70188.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang,
H., Zhang, M., Li, Y., Wu, Y., et al. DeepSeekMath:
Pushing The Limits of Mathematical Reasoning in Open
Language Models. *arXiv preprint arXiv:2402.03300*,
2024.
- Singh, A., Fry, A., Perelman, A., Tart, A., Ganesh, A.,
El-Kishky, A., McLaughlin, A., Low, A., Ostrow, A.,
Ananthram, A., et al. OpenAI GPT-5 System Card. *arXiv*
preprint arXiv:2601.03267, 2025a.
- Singh, G., Wehling, L., Mulyadi, A. W., Sreenath,
R. H., Klabunde, T., Andreani, T., and McCloskey, D.
Talk2Biomodels and Talk2KnowledgeGraph: AI agent-
based application for prediction of patient biomarkers and
reasoning over biomedical knowledge graphs. In *ICLR*
2025 Workshop on Machine Learning for Genomics Ex-
plorations, 2025b.
- Tiwari, K., Kananathan, S., Roberts, M. G., Meyer, J. P.,
Shohan, M. U. S., Xavier, A., Maire, M., Zyoud, A.,
Men, J., Ng, S., Nguyen, T. V. N., Glont, M., Hermjakob,
H., and Malik-Sheriff, R. S. Reproducibility in systems
biology modelling. *Molecular systems biology*, 17(2),
2021. doi: 10.15252/msb.20209982. URL [https://](https://pubmed.ncbi.nlm.nih.gov/33620773/)
pubmed.ncbi.nlm.nih.gov/33620773/.
- von Werra, L., Belkada, Y., Tunstall, L., Beeching,
E., Thrush, T., Lambert, N., Huang, S., Rasul, K.,
and Gallouédec, Q. TRL: Transformers Reinforce-
ment Learning, 2020. URL [https://github.com/](https://github.com/huggingface/trl)
[huggingface/trl](https://github.com/huggingface/trl).
- Wan, G., Ling, M., Ren, X., Han, R., Li, S., and Zhang,
Z. Compass: Enhancing agent long-horizon reasoning
with evolving context. *arXiv preprint arXiv:2510.08790*,
2025.
- Wang, R., Genadi, R. A., El Bouardi, B., Wang, Y., Koto, F.,
Liu, Z., Baldwin, T., and Li, H. AgentFly: Extensible and
Scalable Reinforcement Learning for LM Agents. *arXiv*
preprint arXiv:2507.14897, 2025.

- 550 Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A.,
551 Khashabi, D., and Hajishirzi, H. Self-Instruct: Aligning
552 Language Models with Self-Generated Instructions. In
553 *Proceedings of the 61st Annual Meeting of the Association
554 for Computational Linguistics*, 2023.
- 555
556 Wehling, L., Singh, G., Mulyadi, A. W., Sreenath, R. H.,
557 Hermjakob, H., Nguyen, T. V., Rückle, T., Mosa, M. H.,
558 Cordes, H., Andreani, T., et al. Talk2Biomodels: AI
559 Agent-Based Open-Source LLM Initiative for Kinetic Bi-
560 ological Models. *BMC bioinformatics*, 26(1):276, 2025.
- 561
562 Wei, J., Bosma, M., Zhao, V., Guu, K., Yu, A. W., Lester,
563 B., Du, N., Dai, A. M., and Le, Q. V. Finetuned lan-
564 guage models are zero-shot learners. In *International
565 Conference on Learning Representations*, 2022.
- 566
567 Wei, Q., Zeng, S., Li, C., Brown, W., Frunza, O., Deng,
568 W., Schneider, A., Nevmyvaka, Y., Zhao, Y. K., Garcia,
569 A., and Hong, M. Reinforcing Multi-Turn Reasoning
570 in LLM Agents via Turn-Level Reward Design. *arXiv
571 preprint arXiv:2505.11821*, 2025.
- 572
573 Xiong, Q., Huang, Y., Jiang, Z., Chang, Z., Zheng, Y., Li,
574 T., and Li, M. Butterfly effects in toolchains: A com-
575 prehensive analysis of failed parameter filling in LLM
576 tool-agent systems. In *Findings of the Association for
577 Computational Linguistics: EMNLP 2025*, pp. 16712–
16729. Association for Computational Linguistics, 2025.
- 578
579 Zhang, S., Dong, Y., Zhang, J., Kautz, J., Catanzaro, B.,
580 Tao, A., Wu, Q., Yu, Z., and Liu, G. NemoTron-Research-
581 Tool-N1: Exploring Tool-Using Language Models with
582 Reinforced Reasoning. *arXiv preprint arXiv:2505.00024*,
583 2025.
- 584
585 Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z.,
586 Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. P., Zhang,
587 H., Gonzalez, J. E., and Stoica, I. Judging LLM-as-a-
588 Judge with MT-Bench and Chatbot Arena. In *Advances
589 in Neural Information Processing Systems*, volume 36,
590 2023.
- 591
592 Zhong, H., Zhai, J., Song, L., Bian, J., Liu, Q., and Tan, T.
593 RC-GRPO: Reward-Conditioned Group Relative Policy
594 Optimization for Multi-Turn Tool Calling Agents. *arXiv
595 preprint arXiv:2602.03025*, 2026.
- 596
597
598
599
600
601
602
603
604

A. Synthetic Benchmark Construction

A.1. Synthetic Data Generation Pipeline

We constructed a synthetic dataset using 20 models selected from the BioModels repository (Malik-Sheriff et al., 2020) (Table 3), spanning diverse areas of biology. For each model, we manually extracted the required simulation times from the corresponding publications. Scientifically grounded question-answer pairs were generated using the BasiCO (Bergmann, 2023) library, a Python interface to COPASI (Hoops et al., 2006), by extracting model species, parameters, and metadata, such as names, descriptions, units, and compartment information, using the `get_species()`, `get_parameters()`, and `get_model_info()` functions.

The `get_model_info` tool retrieves descriptive metadata about a model, including its name, compartments, species, parameters, and other relevant information. For this tool, we provided the LLM with metadata obtained using the `get_model_info()` function from BasiCO. The LLM was instructed to generate natural-language questions that a user might plausibly ask when exploring a new biological model (e.g., requesting a description or name of the model), and the corresponding answers. The system prompt used for this tool was provided as Prompt S4 in Section A.2.

The `simulate_model` tool performs a time-course simulation of a model under specified initial conditions. The LLM was provided with the list of species obtained using `get_species()` function from BasiCO. The LLM was instructed to generate a question that selected up to three species whose initial concentrations could be modified prior to simulation and up to three species whose concentrations should be measured after the simulation. We then extracted the original initial concentrations of the selected species using the BasiCO library `get_species()` function and perturbed them by randomly scaling their values to lie between half and twice the original concentration. The question explicitly described these requested modifications and measurements in natural language. The answer associated with this tool invocation confirmed whether the simulation completed successfully. The corresponding system prompt was provided as Prompt S2 in the Appendix. We also executed a time-course simulation using the `run_time_course()` function from BasiCO, and saved the results for later use by the `ask_question` tool.

The `steady_state` tool computes the steady-state concentrations of model species under specified initial conditions. For this tool, the LLM was provided with the list of model species obtained from `get_species()`, and instructed to generate a question specifying up to three species for concentration modification prior to the steady-state computation and up to three species for post-simulation measurement. The initial concentrations of the selected species were perturbed using the same random scaling strategy described above. The question explicitly described these requested modifications and measurements in natural language. The answer to this tool invocation indicated whether the model successfully reached a steady state. The system prompt for this tool was provided as Prompt S3 in the Appendix. We also computed the steady-state solution using the `run_steadystate()` function from BasiCO, and saved the results for later use by the `ask_question` tool.

The `ask_question` tool queries results produced by the `simulate_model` and `steady_state` tools. For this tool, we generated questions that asked for the final concentrations of specific species identified earlier in the interaction. Since the species to be measured were already selected during the generation of the simulation or steady-state questions, the follow-up questions were constructed to explicitly query these species. To generate answers, we extracted the final concentrations of the requested species from the stored simulation results. For time-course simulations, the answer reported the final concentrations of the queried species. For steady-state simulations, if the model failed to reach steady state, the answer did not report concentrations and instead indicated that steady state was not achieved. The system prompt used for this tool is provided in Section A.2.

The `parameter_scan` tool explores how varying a model parameter affected the concentrations of selected species under specified initial conditions. For this tool, the LLM was provided with the list of species obtained from `get_species()` and was instructed to generate a question that selected up to three species for concentration modification and up to three species for post-scan measurement. In addition, the LLM was provided with a list of parameters and instructed to select one parameter for scanning. Using the `get_parameters()` function from BasiCO, we retrieved the initial value of the selected parameter and varied it over a range from half to twice its original value. The scan was performed using a randomly chosen number of steps between 5 and 10. Species concentrations were perturbed prior to the scan using the same random scaling strategy described earlier. The question explicitly described these requested modifications and measurements in natural language. The answer to this tool invocation confirmed that the parameter scan completed successfully. The corresponding system prompt was provided as Prompt S6 in the Section A.2.

A.2. System Prompts

PROMPT S1: SELECT MODIFIABLE AND MEASURABLE SPECIES FOR SIMULATE_MODEL AND STEADY_STATE

- Given:** The list of species in a systems biology model.
- Tasks:** Identify a list of 1-2 species whose initial concentrations can be modified. Identify a list of 1-2 different species whose concentrations can be measured.
- Constraints:** The two lists must not contain the same species. If the model appears to be a QSP model, prioritize pharmacologically relevant species (e.g., drug, biomarker, disease marker). Prefer species that are clinically measurable for the measurement list and species that can be pharmacologically modulated for the modification list. Do not label or name the lists in the response.

PROMPT S2: PREPARE Q&A FOR SIMULATE_MODEL

- Given:** A list of species with updated initial concentrations. A list of species with their expected concentrations after simulation. A model ID (unless this is a follow-up question). Model units.
- Tasks:** Write a natural-language question requesting a simulation of the model for a specified duration and time step, explicitly stating the new initial concentrations. Ask for the concentrations of the specified species at the end of the simulation. Write a corresponding natural-language answer containing the expected concentrations.
- Guidelines:** Species names and numerical values must exactly match those provided. Units must appear inline with numerical values and nowhere else. Write as an expert in systems biology and pharmacology.

PROMPT S3: PREPARE Q&A FOR STEADY_STATE

- Given:** A list of species with updated initial concentrations. A list of species with their expected steady-state concentrations and/or transition times. A model ID (unless this is a follow-up question). Model units.
- Tasks:** Write a natural-language question requesting that the model be run to steady state, explicitly stating the new initial concentrations. Ask for the steady-state concentrations and/or transition times of the specified species. Write a corresponding natural-language answer containing the expected results.
- Guidelines:** Species names and numerical values must exactly match those provided. Units must appear only inline with numerical values. Write as an expert in systems biology and pharmacology.

PROMPT S4: RETRIEVE MODEL INFORMATION

- Given:** Information from one or more of the following model components: species, parameters, compartments, units, description, model name.
- Tasks:** Write a natural-language question asking to retrieve one or more of these components from a systems biology model with a given ID. Write a natural-language answer containing exactly the provided information.
- Guidelines:** Species names and numerical values must exactly match those provided. Units must appear inline with numerical values and nowhere else. Write as an expert in systems biology and pharmacology.

PROMPT S5: SELECT SPECIES/PARAMETERS FOR PARAMETER_SCAN

- Given:** The species and parameters of a systems biology model.
- Tasks:** Randomly select 1-3 species whose initial concentrations will be modified before running a parameter scan. Select 1 species or parameter whose value will be scanned. Select 1-3 species whose concentrations will be measured after the scan. Return only the three lists based on the provided species and parameters.

PROMPT S6: PREPARE Q&A FOR PARAMETER_SCAN

- Given:** A list of species with new initial concentrations. A species or parameter to be scanned. A species to be observed after the scan. A model ID (unless this is a follow-up question). Model units.
- Tasks:** Write a natural-language question requesting a parameter scan. Specify the new initial concentrations. Define the scan range with a start value, end value, and step size for the selected species or parameter. State the simulation duration and time step. Specify the species to be observed. Write a natural-language answer confirming that the parameter scan has been completed.
- Guidelines:** Species names and numerical values must exactly match those provided. Units must appear inline with numerical values and nowhere else. Write as an expert in systems biology and pharmacology.

A.3. Additional Information

Table 3. Manually curated BioModels with ground-truth simulation time points.

Model ID	Model Name	Simulation Time	PubMed ID
BIOMD0000000112	Clarke2006_Smad_signalling	480 min	17186703
BIOMD0000000163	Zi2007_TGFbeta_signaling	8 hrs	17895977
BIOMD0000000186	Ibrahim2008 - Mitotic Spindle Assembly Checkpoint - Dissociation variant	4000 sec	18253502
BIOMD0000000187	Ibrahim2008 - Mitotic Spindle Assembly Checkpoint - Convey variant	4000 sec	18253502
BIOMD0000000424	Faratian2009 - Role of PTEN in Trastuzumab resistance	60 min	19638581
BIOMD0000000357	Lee2010_ThrombinActivation_OneForm_reduced	900 sec	20435402
BIOMD0000000364	Lee2010_ThrombinActivation_OneForm	900 sec	20435402
BIOMD0000000399	Jenkinson2011_EGF_MAPK	60 min	21548948
BIOMD0000000328	Bucher2011_Atorvastatin_Metabolism	600 min	21548957
BIOMD0000000463	Heldt2012 - Influenza Virus Replication	12 hrs	22593159
BIOMD0000000362	Butenas2004_BloodCoagulation	1200 sec	15039440
BIOMD0000000120	Chan2004_TCell_receptor_activation	70 sec	15255048
BIOMD0000000335	Hockin2002_BloodCoagulation	700 sec	11893748
BIOMD0000000528	Fribourg2014 - Dynamics of viral antagonism and innate immune response (H1N1 influenza A virus - Cal/09)	600 days	17532343
BIOMD0000000529	Fribourg2014 - Dynamics of viral antagonism and innate immune response (H1N1 influenza A virus - NC/99)	8 hrs	24594370
BIOMD0000000889	Fribourg2014 - Model of influenza A virus infection dynamics of viral antagonism and innate immune response.	8 hrs	24594370
BIOMD0000000122	Fisher2006_Ca_Oscillation_dpdt_NFAT_dynamics	400 sec	17031595
BIOMD0000000123	Fisher2006_NFAT_Activation	3000 sec	17031595
BIOMD0000000537	Dwivedi2014 - Crohns IL6 Disease model - Anti-IL6R Antibody	2016 hrs	24402116
BIOMD0000000788	Schropp2019 - Target-Mediated Drug Disposition Model for Bispecific Antibodies	150 days	30480383

Table 4. List of tools used to optimize the Talk2BioModels (T2B) agent.

Tool Name	Description	Expected Returns
get_modelinfo	Retrieve model metadata (species, parameters, compartments, units, description, name).	Model metadata summary.
simulate_model	Run time-course simulation with specified initial conditions and time settings.	Simulation results (time-series).
steady_state	Compute steady state given initial conditions.	Steady-state results (concentrations, transition time).
parameter_scan	Scan a parameter over values and observe species responses.	Parameter scan results.
ask_question	Answer a question over simulation or steady-state results using the dataframe.	Numeric/text answer derived from results.

Table 5. Scenario templates and synthetic benchmark generated by our pipeline.

Scenario	Tool chain (multi-turn)	Turns/Conv	#Conv	#Turns	Example
1	get_modelinfo → simulate_model → ask_question	3	120	360	Units query for BIOMD0000000112.
2	steady_state → ask_question	2	120	240	Steady state with new initial conditions.
3	simulate_model → ask_question → ask_question	3	120	360	Simulation with two follow-up queries.
4	get_modelinfo → simulate_model, ask_question → steady_state → ask_question	4	120	480	Compartments and species query.
5	parameter_scan	1	120	120	Scan a parameter and confirm completion.
6	simulate_model → parameter_scan	3	120	360	Simulate then scan a parameter.
7	get_modelinfo → steady_state, ask_question → parameter_scan	3	120	360	Info then steady state and scan.
8	steady_state → parameter_scan	2	120	240	Steady state followed by scan.
9	simulate_model → simulate_model → ask_question	3	120	360	Two simulations then query.
10	get_modelinfo → steady_state, ask_question → ask_question	3	120	360	Info then steady state with follow-ups.
Total			1200	3240	

B. Implementation Details

We trained the policy model using a two-stage pipeline with LoRA (Hu et al., 2022) adapters: SFT followed by GRPO with verifiable rewards. For SFT, we trained Qwen2.5-3B-Instruct with TRL’s (von Werra et al., 2020) `SFTTrainer` on scenario-stratified multi-turn splits, using assistant-only teacher-forcing over tool calls and final answers. We used a tool-calling chat template with structured response formatting to preserve schema-valid invocation patterns. LoRA settings were $r = 16$, $\alpha = 32$, and dropout = 0.05, with adaptation on attention and feed-forward projections.

For GRPO, we initialized from the SFT adapter and trained at turn level. In each update, we sampled $G = 8$ completions per prompt at temperature 0.9 and used a composite verifiable reward combining three components: (i) deterministic tool correctness via exact-match sequence comparison ($\lambda_{\text{tool}} = 0.4$), (ii) deterministic argument correctness with soft numeric tolerance bands ($\lambda_{\text{arg}} = 0.4$), and (iii) task-completion scoring from a frozen LLM judge ($\lambda_{\text{task}} = 0.2$). No learned reward model was used. Cascaded gating ensures argument and task scores are zeroed when tool correctness is zero, preventing gradient signal from incorrect tool invocations.

The argument scoring function assigns graduated credit based on relative numeric error: exact matches receive 1.0, errors within 1% receive 0.9, within 10% receive 0.7, within 100% receive 0.3, and larger errors receive 0.0. Non-numeric arguments use exact string matching. Unverifiable fields (experiment names, free-text questions, ambiguous interval semantics) are stripped before comparison.

Note that training uses these soft tolerance bands while evaluation uses near-exact matching (relative tolerance 10^{-9}). This is a deliberate design choice, not a misalignment: binary argument scoring on numeric floats would provide near-zero gradient signal to a 3B model (most outputs would score 0.0), whereas graduated credit creates a smooth optimization landscape that rewards progressively more precise argument generation. Crucially, soft argument scoring is safe because it is gated behind binary tool correctness the model only receives argument reward when it has already called the correct tool, preventing partial-credit exploitation of the kind that motivated our switch from soft to binary tool scoring (see Section 2.3).

The task-completion judge is GPT-4o mini (Singh et al., 2025a), a frozen proprietary model used to score open-ended task completion. When the policy ends generation on a tool call without producing a final text answer, actual tool execution outputs are used as a proxy answer for judging.

State-aware online scoring replays prior-turn expected tool calls from the dataset to reconstruct turn-local context (e.g., populating simulation dataframes before scoring an `ask_question` turn). The KL penalty coefficient was $\beta = 0.1$, with a generation batch size of 8.

As for evaluation metrics, we used a stateful execution harness and computed argument correctness, tool correctness, and task completion on held-out conversations. For *argument correctness*, we evaluated whether each tool call included the correct input parameters given the user query and expected tool usage, with

$$\text{Argument Correctness} = \frac{1}{|\mathcal{I}|} \sum_{i=1}^{|\mathcal{I}|} m(\tilde{\mathbf{a}}_i, \mathbf{a}_i), \quad (2)$$

where \mathcal{I} indexes the expected tool calls within a turn, $\tilde{\mathbf{a}}_i$ and \mathbf{a}_i denote the predicted and ground-truth arguments, and $m(\cdot, \cdot) \in \{0, 0.5, 1\}$ returns 1 if all fields match (using near-exact numeric comparison with relative tolerance 10^{-9}), 0.5 if a strict subset matches, and 0 otherwise. We estimated *tool correctness* by comparing the executed tool calls against the expected ones for each turn using exact matching with ordering, defined as

$$\text{Tool Correctness} = \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbb{I}[\hat{t}_i = t_i], \quad (3)$$

where \hat{t}_i and t_i denote the predicted and expected tool names at position i in the ordered sequence for a given turn. Task completion was computed using an LLM judge over stateful traces, and aggregated over turns as

$$\text{Task Completion} = \frac{1}{|\mathcal{U}|} \sum_{u=1}^{|\mathcal{U}|} s_u, \quad (4)$$

where $s_u \in [0, 1]$ is the judge score for turn u and \mathcal{U} is the set of evaluated turns.

C. Ablation Studies

Table 6 reports progressive training-stage comparisons: the untuned baseline, SFT-only, and the final SFT+GRPO-optimized policy with verifiable rewards and state-aware online scoring. SFT provides the largest absolute improvement, lifting argument correctness from 24.8% to 85.6% (+60.8 pp) by teaching the model structured tool-call formatting and domain-specific argument schemas. GRPO further refines the policy, improving tool correctness from 96.0% to 98.8% and argument correctness from 85.6% to 91.5%. The marginal task completion gain (+0.3 pp) from GRPO is consistent with the reward design, which allocates 80% of the signal to deterministic tool and argument verification; task completion, assessed by a frozen LLM judge, receives only 20% of the reward weight (see Section D).

Table 6. Stage-wise ablation. GRPO uses deterministic tool/argument rewards with cascaded gating and a frozen LLM judge for task completion (no learned reward model).

Model Variant	Tool Correctness	Argument Correctness	Task Completion
Qwen2.5-3B baseline	0.682	0.248	0.524
+ SFT	0.960	0.856	0.729
+ SFT+GRPO (verifiable rewards)	0.988	0.915	0.732

D. Additional Analyses

Table 7 presents a per-tool breakdown across the three training stages. Key observations: (i) The base model can often identify the correct single tool (e.g., 93.3% for `simulate_model`) but catastrophically fails at argument generation (6.7%), confirming that tool selection and argument construction are distinct capabilities. (ii) SFT provides the largest absolute gains, particularly on argument correctness for complex tools like `steady_state` (0.0 \rightarrow 0.778) and `parameter_scan` (0.0 \rightarrow 0.646). (iii) GRPO further refines argument accuracy, with the most notable improvement on `steady_state` (0.778 \rightarrow 1.000, +22.2 pp). (iv) On multi-tool turns (`simulate_model+ask_question`), GRPO achieves perfect tool and argument correctness (1.000), while task completion (30.8%) is lower due to the open-ended nature of the LLM judge metric on these complex turns; this pattern is consistent across all evaluated models, indicating a metric-level property rather than a framework limitation.

Table 7. Per-tool breakdown on the held-out test split. Multi-tool entries denote turns requiring multiple tool invocations. N is the number of evaluated turns per tool type.

Tool(s)	Model	N	Tool	Arg	Task
<code>ask_question</code>	Base	96	0.750	0.719	0.435
	SFT	96	1.000	1.000	0.639
	SFT+GRPO	96	1.000	1.000	0.642
<code>get_modelinfo</code>	Base	48	0.833	0.021	0.789
	SFT	48	0.979	0.729	0.949
	SFT+GRPO	48	1.000	0.833	0.973
<code>simulate_model</code>	Base	60	0.933	0.067	0.816
	SFT	60	1.000	0.933	0.985
	SFT+GRPO	60	1.000	0.967	0.982
<code>steady_state</code>	Base	36	0.944	0.000	0.739
	SFT	36	1.000	0.778	0.869
	SFT+GRPO	36	1.000	1.000	0.861
<code>parameter_scan</code>	Base	48	0.292	0.000	0.245
	SFT	48	0.812	0.646	0.503
	SFT+GRPO	48	0.938	0.688	0.490
<code>simulate_model+ask_question</code>	Base	12	0.167	0.250	0.150
	SFT	12	1.000	0.958	0.350
	SFT+GRPO	12	1.000	1.000	0.308
<code>steady_state+ask_question</code>	Base	24	0.125	0.146	0.042
	SFT	24	0.875	0.833	0.442
	SFT+GRPO	24	0.958	0.896	0.496

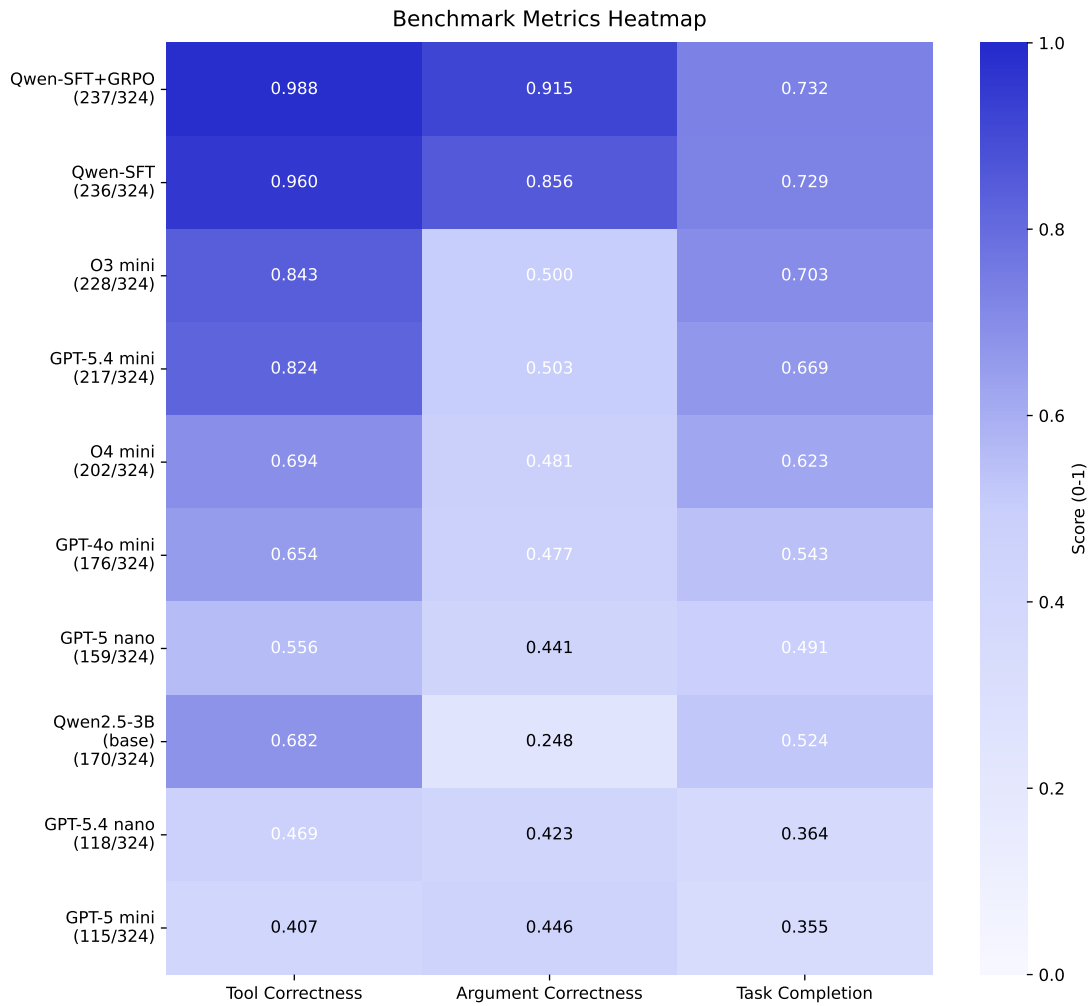


Figure 4. Benchmark metrics heatmap across all evaluated models. Models are sorted by composite performance (top to bottom). Numbers in parentheses indicate successfully scored turns out of 324 total. Qwen-SFT+GRPO (top) dominates on tool and argument correctness; all proprietary models score below 0.51 on argument correctness.

E. Per-Scenario Analysis

Figures 5-7 present per-scenario heatmaps for each evaluation metric across all 10 scenarios and 10 models.

Tool Correctness (Fig. 5) Qwen-SFT+GRPO achieves perfect tool correctness (1.000) on 8 out of 10 scenarios, with the only imperfect scores on Scenario 5 (`parameter_scan` only, 0.933) and Scenario 7 (0.944). In contrast, proprietary models show high variance across scenarios: O3-mini ranges from 0.333 (Scenario 10) to 1.000 (Scenario 9), while GPT-5 mini scores below 0.500 on 7 of 10 scenarios. The base Qwen2.5-3B model struggles most on multi-tool scenarios (Scenarios 4, 7, 10) where it must chain multiple tool invocations.

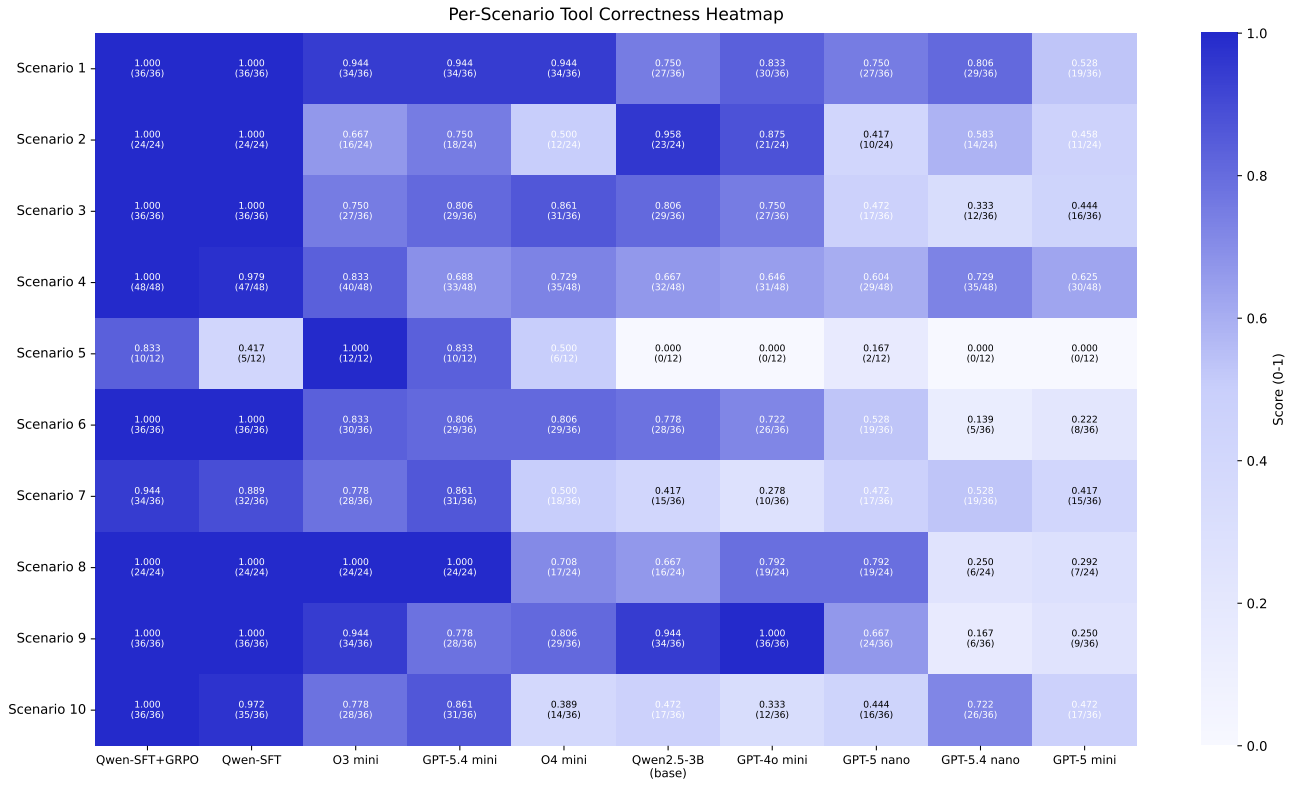


Figure 5. Per-scenario tool correctness heatmap. Each cell shows the mean score with count in parentheses. Qwen-SFT+GRPO achieves perfect or near-perfect tool selection across nearly all scenarios.

Argument Correctness (Fig. 6) The argument correctness gap between our optimized models and proprietary baselines is most pronounced. Qwen-SFT+GRPO exceeds 0.80 on 8 of 10 scenarios, whereas the best proprietary model (GPT-5.4-mini) never exceeds 0.60 on any scenario. Scenarios involving `parameter_scan` (5, 6, 7, 8) are consistently harder for all models due to the complex argument schema (scan ranges, step sizes, multiple species). The SFT→GRPO gain is most visible on Scenarios 2, 8, and 10 where argument structures involve steady-state and multi-tool chains.

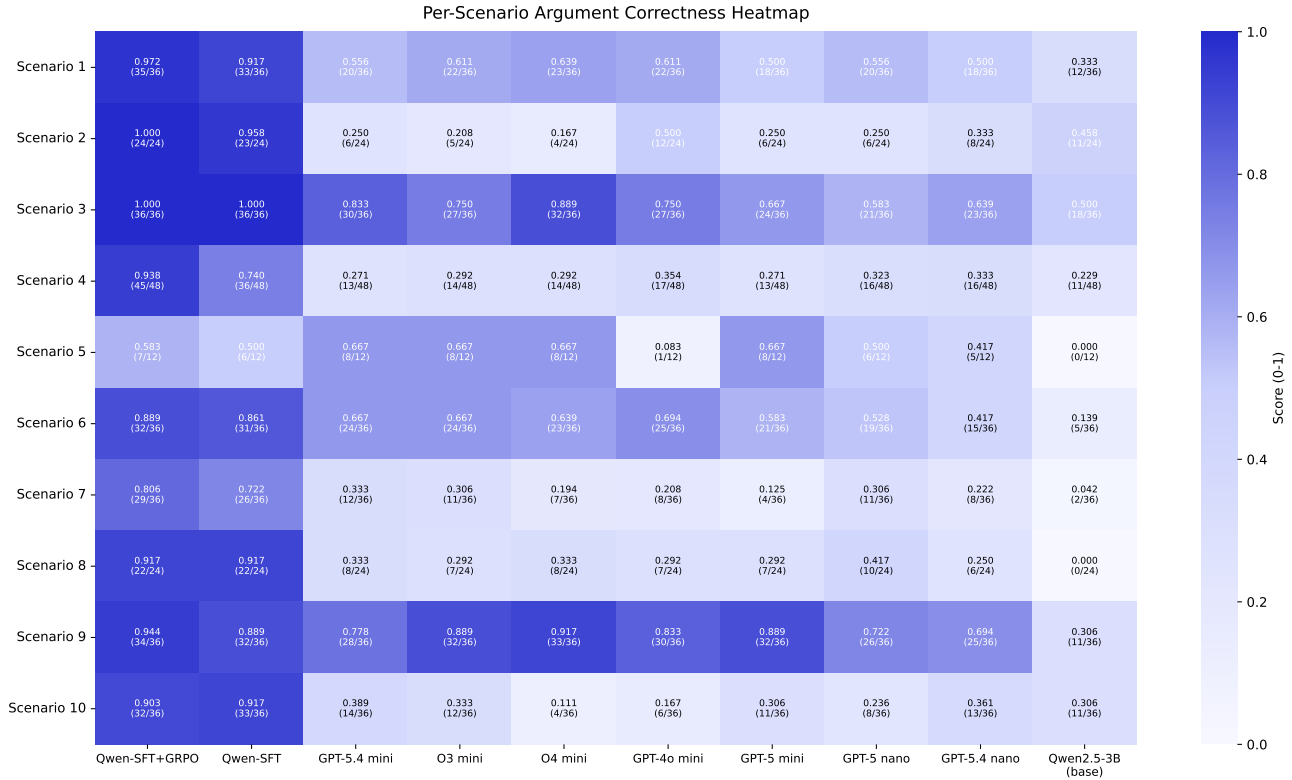


Figure 6. Per-scenario argument correctness heatmap. The gap between Qwen-SFT+GRPO/SFT and proprietary models is largest here, demonstrating the value of domain-specific fine-tuning for structured argument generation.

Task Completion (Fig. 7) Task completion scores are more uniform across models, reflecting the open-ended nature of the LLM judge metric. Qwen-SFT+GRPO leads on most scenarios but the margins are smaller, particularly on Scenarios 5-8 involving parameter_scan and multi-tool chains. The narrower gaps on this metric are consistent with the reward design, which allocates 80% of the signal to deterministic tool and argument verification; task completion reflects natural-language answer quality, where all models produce lower scores on complex tool outputs.

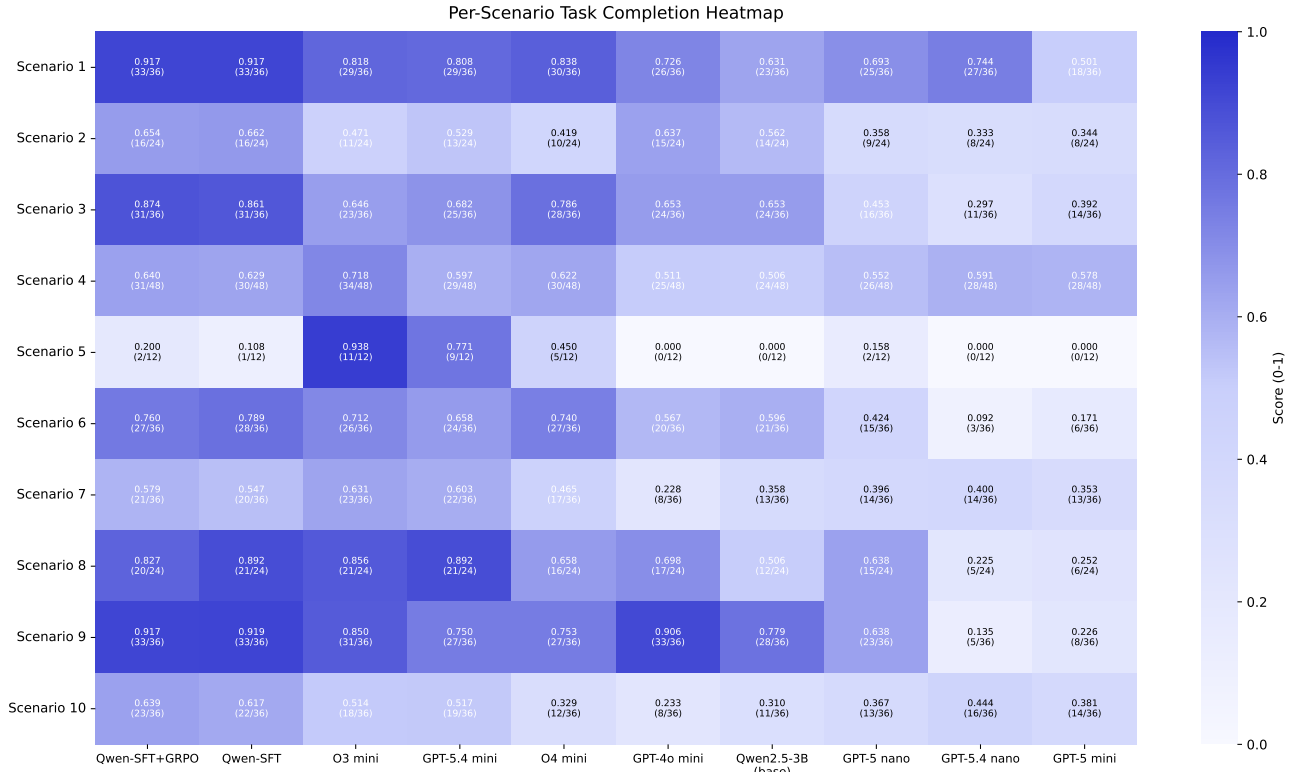


Figure 7. Per-scenario task completion heatmap. Scores are more uniform across models on this metric, reflecting the open-ended nature of LLM judge evaluation. Qwen-SFT+GRPO leads on most scenarios, with narrower margins on parameter scan scenarios (5-8) where all models produce lower scores.

1100 **Turn-Level Breakdown by Scenario (Fig. 8)** Figure 8 decomposes Figure 3 into individual scenarios, revealing where
1101 task completion degrades and why. Two patterns emerge. First, standalone `ask_question` turns consistently show reduced
1102 task completion across all models. In Scenario 1 (Turn 3), Scenario 2 (Turn 2), Scenario 3 (Turns 2 and 3), and Scenario 9
1103 (Turn 3), task completion drops even when tool and argument correctness remain high. These turns require parsing large
1104 simulation or steady-state dataframes to extract precise numeric answers, a task that degrades when the tabular output
1105 exceeds the answering model’s effective capacity. Second, multi-tool combination turns exhibit more severe degradation.
1106 In Scenario 4 (Turn 2: `simulate_model+ask_question`), Scenario 7 (Turn 2: `steady_state+ask_question`),
1107 and Scenario 10 (Turn 2: `steady_state+ask_question`), all models including our optimized agent show sharper
1108 drops in task completion. These turns require executing a state-building tool and immediately querying its results within a
1109 single turn, compounding the dataframe parsing challenge with intra-turn tool coordination. Notably, tool and argument
1110 correctness remain near-perfect for Qwen-SFT+GRPO on these same turns, confirming that the task completion bottleneck
1111 arises from answer generation rather than tool execution. This breakdown motivates two future directions: replacing the
1112 current `ask_question` tool with a dedicated data-analysis agent capable of handling large tabular outputs, and developing
1113 decomposition strategies that explicitly optimize intra-turn multi-tool interactions.

1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154

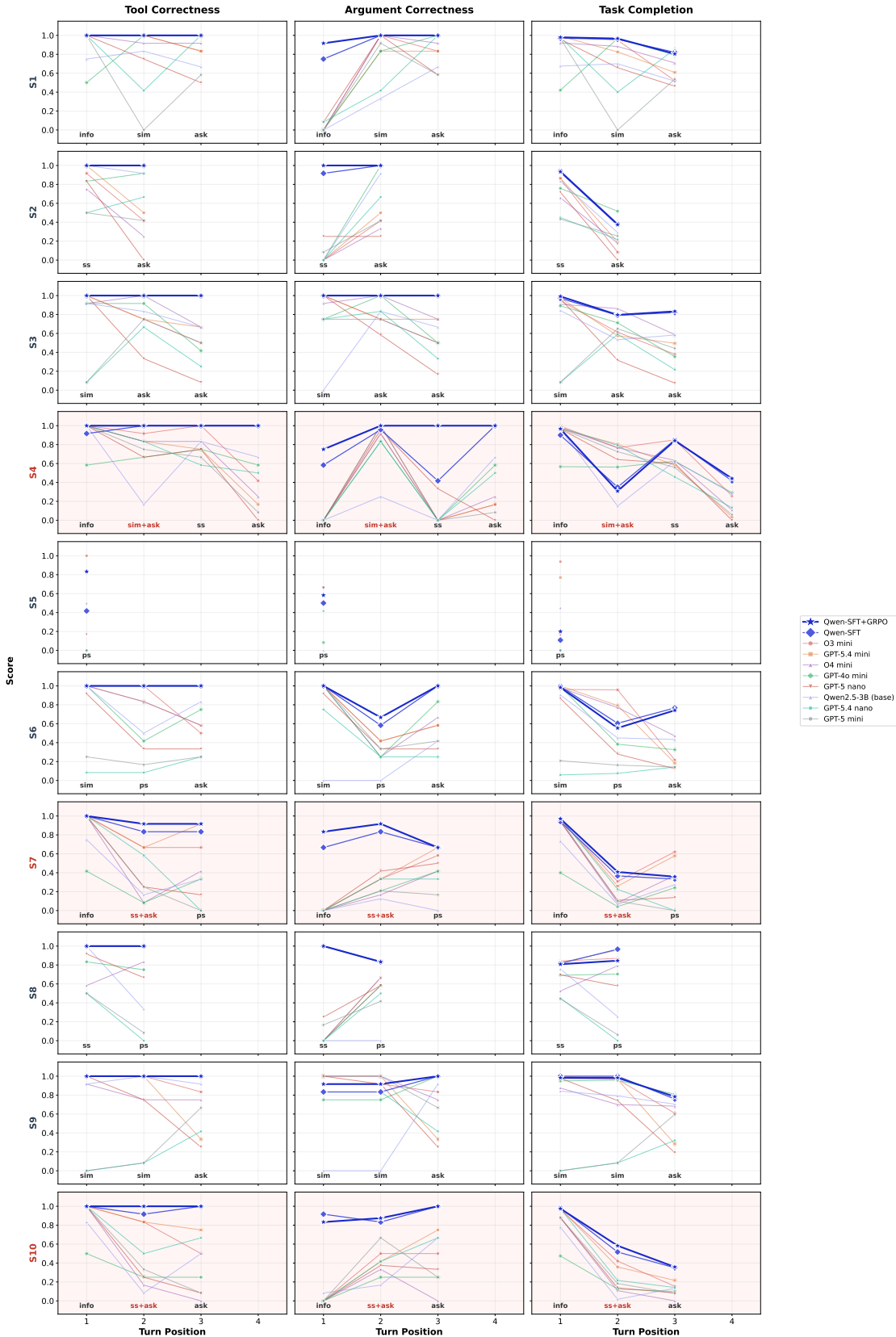


Figure 8. Per-scenario turn-level performance breakdown across all evaluated models and three metrics. Each row corresponds to one scenario (S1-S10, see Table 5 for full definitions); columns show tool correctness, argument correctness, and task completion. Tool labels below each turn indicate which tool(s) are invoked at that position; multi-tool turns are highlighted in red. Standalone ask_{question} turns and multi-tool combination turns (highlighted rows: Scenarios 4, 7, 10) exhibit the sharpest task completion drops, even when tool and argument correctness remain high.

F. Out-of-Distribution Generalization

To further assess whether the optimized policy generalizes beyond the training distribution, we evaluate Qwen-SFT+GRPO on a separate out-of-distribution (OOD) benchmark comprising 9 scenarios (340 conversations, 571 turns) that differ from the 10 in-distribution training scenarios in tool composition, tool novelty, or both. In particular, the OOD scenarios are organized into three categories (Table 8):

1. **Novel composition of training tools** (Scenario 18): Tool chains that use only tools seen during training but in combinations never encountered (*e.g.*, `parameter_scan` \rightarrow `steady_state` \rightarrow `ask_question`).
2. **Novel tools only** (Scenarios 11, 16): Single-turn invocations of tools that have zero training examples (`search_models`, `get_annotation`). The model must rely entirely on schema-following capabilities acquired through SFT and GRPO.
3. **Mixed training and novel tools** (Scenarios 12-15, 17, 19): Multi-turn chains that interleave training tools with novel tools (*e.g.*, `simulate_model` \rightarrow `custom_plotter` \rightarrow `ask_question`), testing compositional generalization across both familiar and unseen tool interfaces.

Furthermore, Table 9 presents the per-category and overall results. Overall, the optimized agent achieves 86.3% tool correctness, 76.9% argument correctness, and 68.9% task completion on the OOD benchmark. While these scores represent a decrease from the in-distribution results (98.8%/91.5%/73.2%), the degradation is moderate and non-uniform across categories. The mixed category, where novel tools are interleaved with training tools, shows the strongest transfer (tool=0.967, arg=0.842, task=0.818), indicating that familiarity with a subset of the tool chain provides sufficient scaffolding for the model to generalize to unseen tools via schema following. The novel tools category (S11, S16), despite having zero training examples, achieves 95.7% tool correctness and 91.3% argument correctness, demonstrating that the SFT+GRPO pipeline equips the model with robust schema-following capabilities that transfer to entirely new tool interfaces. The weakest performance occurs in the novel composition category (S18: `parameter_scan` \rightarrow `steady_state` \rightarrow `ask_question`), where all individual tools were seen during training but their specific combination was not. This scenario achieves only 70.6% tool correctness and 66.3% argument correctness. Inspection reveals that the primary failure mode is on `parameter_scan` turns, consistent with the in-distribution finding that this tool has the most complex argument schema (Table 7). When combined with the reversed ordering relative to training scenarios (where `steady_state` typically precedes or follows `parameter_scan` in different configurations), the model struggles with argument construction, particularly for scan range specification and interval computation.

These results suggest that our optimization framework provides meaningful generalization to out-of-distribution tool compositions and novel tool interfaces, with the strongest transfer occurring when at least part of the tool chain overlaps with training experience.

Table 8. Out of distribution evaluation scenarios. Tools marked with * are novel (zero training examples). Category labels: **NC** = novel composition of training tools, **NT** = novel tools only, **MX** = mixed training and novel tools.

Scenario	Category	#Turns	Tool Chain
18	NC	3	<code>parameter_scan</code> \rightarrow <code>steady_state</code> \rightarrow <code>ask_question</code>
11	NT	1	<code>search_models</code> *
16	NT	1	<code>get_annotation</code> *
12	MX	2	<code>search_models</code> * \rightarrow <code>get_modelinfo</code>
13	MX	2	<code>simulate_model</code> \rightarrow <code>custom_plotter</code> *
14	MX	2	<code>get_modelinfo</code> \rightarrow <code>get_annotation</code> *
15	MX	3	<code>simulate_model</code> \rightarrow <code>custom_plotter</code> * \rightarrow <code>ask_question</code>
17	MX	3	<code>simulate_model</code> \rightarrow <code>ask_question</code> \rightarrow <code>custom_plotter</code> *
19	MX	3	<code>search_models</code> * \rightarrow <code>simulate_model</code> \rightarrow <code>parameter_scan</code>

Table 9. Out of distribution generalization results for Qwen-SFT+GRPO, grouped by category. In-distribution (ID) results are shown for reference. N denotes the number of evaluated turns.

Category	N	Tool	Arg	Task
Novel composition (S18)	163	0.706	0.663	0.468
Novel tools (S11, S16)	23	0.957	0.913	0.700
Mixed (S12-15, 17, 19)	368	0.967	0.842	0.818
OOD overall	571	0.863	0.769	0.689
ID reference (Table 1)	324	0.988	0.915	0.732

1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319