PLANNER-R1: REWARD SHAPING ENABLES EFFI-CIENT AGENTIC RL WITH SMALLER LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

We investigated Agentic RL with large language models on the TRAVELPLAN-NER benchmark. Our approach, PLANNER-R1, achieved a 56.9% final-pass rate with only 180 training queries, a $2.7 \times$ improvement over GPT-5's 21.2% baseline and the strongest agentic result on the public leaderboard. A central finding was that smaller models (8B) were highly responsive to reward shaping: with dense process-level signals, they reached competitive performance while being 3.5× more compute-efficient and 1.5× more memory-efficient than 32B models. Larger models were more robust under sparse rewards but exhibited smaller relative gains from shaping and higher variance across runs. While curriculum learning offered no significant benefit, shaped rewards consistently amplified learning dynamics, making 8B models the most efficient setting for agentic RL. Crucially, these gains did not come at the cost of overfitting: fine-tuned models mostly maintained or exceeded baseline performance on out-of-domain tasks, including MULTI-IF, NAT-URALPLAN, and τ -BENCH. These results establish reward shaping as a decisive lever for scaling agentic RL, highlight the competitive strength of smaller models, and demonstrate that efficiency can be achieved without sacrificing generalization.

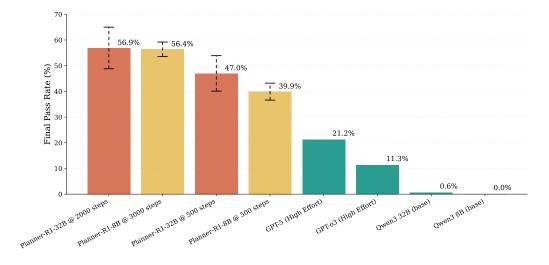


Figure 1: Final-pass rate on the leaderboard test set for tool-use travel planning. Our Planner-R1 models outperformed SOTA LLMs reaching 56.9% average final pass rate.

1 Introduction

Large Language Models (LLMs) have recently posted striking gains in deliberate reasoning and decision making, propelled in part by large-scale reinforcement learning (RL) that trains models to think before they answer (OpenAI et al., 2024a; Guo et al., 2025). Beyond language understanding, LLM agents now demonstrate emerging competence in structured reasoning, tool use, and multistep problem solving across embodied and web environments (Wang et al., 2023a; Huang et al.,

2024; Feng et al., 2025). Yet turning these abilities into *reliable* long-horizon execution under real-world constraints remains challenging: prompting-only agents such as ReAct and Reflexion frequently mis-sequence actions, loop, or hallucinate when tasks demand coordinated tool use and strict constraint satisfaction (Yao et al., 2023b; Shinn et al., 2023).

Planning tasks such as meeting scheduling and multi-day itineraries are demanding: agents must coordinate *heterogeneous tools* (calendars, maps, flights, booking APIs), satisfy *hard, interdependent constraints*, and maintain *global consistency* over long horizons. TRAVELPLANNER makes these difficulties concrete by casting travel itinerary creation as tool-augmented, constraint-driven planning (Xie et al., 2024). The benchmark provides a sandbox with nearly four million records and 1,225 curated intents with reference plans, and evaluates whether an agent can gather evidence via tools and synthesize itineraries that satisfy both explicit user constraints and commonsense feasibility. At release, even strong models struggled—e.g., GPT-4-Turbo with ReAct achieved only a 0.6% *final pass rate* on the 1,000-example test split—underscoring the gap between fluent language modeling and dependable constraint-aware planning (Xie et al., 2024).

To close this gap, researchers have explored different training paradigms. A natural starting point is behavior cloning via supervised fine-tuning (SFT), where a teacher generates "golden" trajectories and a policy maximizes their likelihood, often masking environment observations and tool outputs. While simple and widely used, SFT largely imitates expert behavior and is brittle under distribution shift or suboptimal data. This motivates the search for approaches that directly optimize for end-task success rather than imitation fidelity. RL provides precisely such a mechanism: rewards encode task success, and the policy is updated to increase the likelihood of action sequences that satisfy constraints while suppressing those that fail. Recent work has shown that RL can deliver state-of-the-art gains in model-based reasoning and planning (OpenAI et al., 2024a; Guo et al., 2025), making it a promising direction for tackling long-horizon tool use in TRAVELPLANNER. In addition to model performance, there is growing interest in building efficient agentic systems with smaller models (Belcak et al., 2025). Such models show promising potential for inference and training efficiency, but there remains limited understanding of how agentic RL can best improve their performance without overfitting. Our study addresses this gap by examining how model size, reward shaping, and efficiency interact in agentic RL.

We formulate TRAVELPLANNER as a multi-step, tool-use MDP with constraint-aware planning, where the agent gathers missing facts, reconciles conflicts, and outputs a structured itinerary. Training uses agentic RL with trajectory-level rewards gated by schema validity. Our main focus is the role of *reward density*: we vary feedback from dense, process-level signals to sparse final-pass rewards, and also test a curriculum that transitions between them. All reward variants are *properly shaped*, ensuring they converge to the same optimal policy while revealing how granularity influences learning dynamics. Our contributions are summarized below.

- **SOTA Tool-Use on TravelPlanner** PLANNER-R1-32B achieved a **56.9**% final-pass rate on the official 1,000-query test split, a 2.7× improvement over GPT-5. This is the strongest agentic result on TRAVELPLANNER, demonstrating that RL-tuned models can surpass state-of-the-art proprietary models.¹
- Reward shaping dynamics We find a strong link between reward granularity and policy competence. Smaller models (8B) were especially responsive to shaped, process-level rewards, achieving performance competitive with 32B models while being up to 3.5× more compute-efficient and 1.5× more memory-efficient. Larger models (32B) performed well across reward settings and remained more robust under sparse signals, but exhibited higher variance. In contrast, 8B models depended more heavily on dense shaping. Curriculum learning alone provided no measurable benefit, whereas reward shaping consistently amplified learning dynamics, making the 8B models the most efficient setting for agentic RL.
- Generalization Beyond Training Domain Our agents did not overfit to TRAVELPLANNER: Planner-R1 models mostly maintained or exceeded baseline performance on out-of-domain tasks including MULTI-IF, NATURALPLAN, and τ-BENCH. This demonstrates that the efficiency gains from agentic RL come without sacrificing robustness, supporting transfer to diverse planning and tool-use settings.

¹Hao et al. (2025) achieved 93.9% correctness with external SAT/SMT solvers; our focus is on end-to-end agentic planning without such solvers.

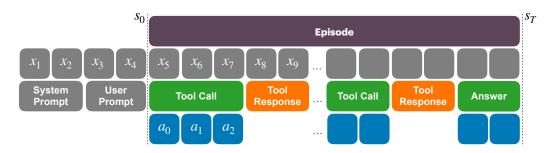


Figure 2: MDP Visualization. x_i represent the *i*th token, while a_t represents the action the agent took at time t. Notice that initial prompts and tool responses contain tokens, but they dont increase the time step t.

• RL Benchmark Formulation We recast TRAVELPLANNER as a multi-step agentic RL benchmark by leveraging the official sandbox and its seven tools, and we designed verifiable reward functions aligned with the task's success criteria. Policies were trained with VERL(ver, 2024), where our system-level optimizations reduced runtime and memory usage by 20%, enabling efficient large-scale experimentation. (see Appendix A for details)

2 PLANNER RL

2.1 PROBLEM FORMULATION

We cast tool-augmented planning as a Markov Decision Process (MDP) $\mathcal{M}=(\mathcal{S},\mathcal{A},P,r,\gamma)$. Since our MDP is episodic, we set $\gamma=1$. Each episode is initialized with two textual inputs: a *system prompt y*, which defines the agent's role and available tools (see Appendix B.1), and a *user prompt u*, which specifies the task goal and user preferences. At each time step t, the agent interacts with the environment by emitting a token, alternating between natural language and structured tool invocations, until it decides that a complete plan has been formed. Figure 2 illustrates this process. While our instantiation focuses on the Travellanner benchmark (Xie et al., 2024), the formulation is general and extends naturally to other agentic RL tasks. We next describe the individual components of the MDP:

States. $s_t \in \mathcal{S}$ denotes the complete history, including the initial system and user prompt, the agent's partial plan, and all tool calls and responses observed up to step t, beginning from $s_0 = (y, u)$.

Actions. $a_t \in \mathcal{A}$ is the generated token at time t. The agent issues tool calls through tokens to gather the necessary information and then produces the final plan through a text action. Tool calls are realized as seven APIs connected to a sand-box with millions of grounded records: search_flights, search_accommodations, search_restaurants, search_attractions, search_ground_transportation, get_cities, and calculator. Each call takes JSON arguments and is wrapped inside <tool_call>...</tool_call>, returning a structured JSON object: successful calls yield a list of serialized rows, while failures return an error field.

Compared to the original TRAVELPLANNER, we added the calculator API for explicit numeric reasoning and disabled the lightweight semantic memory so that tool responses appear directly in the context. The final text action directly outputs an itinerary enclosed in <answer>...</answer>...</answer>...</answer>...</arrangled the final deliverable unambiguous.

Transitions. The environment appends each action to the state; if a tool call is completed, it is executed and the output o_t is added, otherwise o_t is null. The next state is $s_{t+1} = (s_t, a_t, o_{t+1})$, with older context truncated when exceeding the window. A key difference from the original benchmark (Xie et al., 2024) is that we append tokens chronologically to the state, making our transition more generic, as opposed to moving the tool responses to a specific part of the context.

Reward. In this domain, success is sparse and binary. A plan receives a reward of one only at termination if it is schema-valid and satisfies both commonsense and user-specified constraints. User queries are designed to ensure that at least one feasible plan exists.

To pass schema validation, the plan must be a valid JSON array of day-level objects, each conforming to a fixed schema with fields for days, city, transportation, attraction, accommodation, breakfast, lunch, dinner. Importantly, city and transportation are typed objects with required fields (e.g., transportation must specify mode, origin, destination, and duration), rather than free-form strings. The full schema is provided in Appendix B.2.

Constraints fall into two categories. First, there are $N_{\rm cs}$ commonsense constraints, which are not explicitly given to the agent but must nonetheless be satisfied (e.g., transportation segments cannot overlap). Second, there are $N_{\rm hard}$ hard constraints, explicitly specified in the user prompt, such as departure and return dates. Formal definitions and the complete list of constraints are provided in the work of Xie et al. (2024).

Our objective is to learn a policy $\pi_{\theta}(a \mid s)$ that maximizes the expected cumulative reward, which here reduces to optimizing the terminal reward: $\max_{\theta} \mathbb{E}_{\pi_{\theta}}[r_T]$.

2.2 Multi-Stage Reward

Due to the extreme sparsity of the reward function, we shape it using auxiliary metrics defined in the original paper. In particular,

- $r_{\text{schema}} = \mathbb{I}[\text{plan conforms to schema}]$: indicator of schema compliance,
- $r_{cs}^{micro} = \frac{S_{cs}}{N_{cs}}$: fraction of satisfied commonsense constraints,
- $r_{\text{hard}}^{\text{micro}} = \frac{S_{\text{hard}}}{N_{\text{hard}}}$: fraction of satisfied hard constraints,
- $r_{\rm cs}^{\rm macro} = \mathbb{I}[r_{\rm cs}^{\rm micro} = 1]$: indicator that all commonsense constraints pass,
- $r_{
 m hard}^{
 m macro} = \mathbb{I} ig[r_{
 m hard}^{
 m micro} = 1 ig]$: indicator that all hard constraints pass,
- $r_{\text{pass}} = \mathbb{I}[r_{\text{cs}}^{\text{macro}} \wedge r_{\text{hard}}^{\text{macro}}]$: indicator that both commonsense and hard constraints pass.

Here, \mathbb{I} is the indicator function. The micro rewards are necessary to provide partial credit when all constraints are not met, the macro rewards emphasize satisfying entire categories, and $r_{\rm pass}$ corresponds to the original evaluation metric. The terminal reward in the generic form can then be written as:

$$r = r_{\text{schema}} \left(\lambda_1 r_{\text{cs}}^{\text{micro}} + \lambda_2 r_{\text{hard}}^{\text{micro}} + \lambda_3 r_{\text{cs}}^{\text{macro}} + \lambda_4 r_{\text{hard}}^{\text{macro}} + \lambda_5 r_{\text{pass}} \right). \tag{1}$$

By adjusting $\lambda = [\lambda_1, \dots, \lambda_5]$, we control the reward density. In practice, we consider three stages:

- Stage 1: $\lambda = [1, 1, 1, 1, 1]$ (dense feedback),
- Stage 2: $\lambda = [0, 0, 1, 1, 1]$ (category-level),
- Stage 3: $\lambda = [0, 0, 0, 0, 1]$ (sparse final pass).

This setup defines proper reward shaping: auxiliary terms provide intermediate guidance, while the final-pass reward captures the true objective. Crucially, all of the above weightings preserve the same optimal policy. Building on this, we define a curriculum that schedules λ across training, beginning with dense feedback for partial credit, then shifting to category-level rewards, and finally collapsing to the sparse end reward. Transitions occur at predefined step counts.

2.3 OPTIMIZATION

We used GRPO (Shao et al., 2024), a clipped PPO-style objective without KL regularization. For each planning query $u \in \mathcal{D}$, we sample G trajectories $\mathcal{T} = \{\tau_i\}_{i=1}^G$ with corresponding Returns $\mathbf{r} = \{r_1, r_2, \cdots, r_G\}$ from the behavior policy $\pi_{\theta_{\text{old}}}$, where $\tau_i = (s_0^i, a_0^i, \dots, s_{T_i}^i)$. The loss is

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{u \sim \mathcal{D}, \{\tau_i\} \sim \pi_{\theta_{\text{old}}}} \left[\frac{1}{G} \sum_{i=1}^{G} \frac{1}{T_i} \sum_{t=0}^{T_i-1} \min \left(\rho_{\theta}^{i,t} \, \hat{A}_i, \, \operatorname{clip}(\rho_{\theta}^{i,t}, 1 - \epsilon, 1 + \epsilon) \, \hat{A}_i \right) \right], \quad (2)$$

with clipping hyperparameter $\epsilon>0$. The token-level importance ratio and trajectory-level advantage are defined as

$$\rho_{\theta}^{i,t} = \frac{\pi_{\theta}(a_t^i \mid s_t^i, a_{\leq t}^i)}{\pi_{\theta_{\text{old}}}(a_t^i \mid s_t^i, a_{\leq t}^i)}, \qquad \hat{A}_i = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})}.$$

3 EMPIRICAL RESULTS

3.1 SETUP

In-Domain We fine-tuned Qwen3 8B/32B models across 5 runs with fixed set of seeds on TRAV-ELPLANNER. Due to GPU memory and context budget constraints, and inspired by the recent findings that thinking may not always improve performance (Gema et al., 2025; Shojaee* et al., 2025), we did not enable "thinking-in-context." We named these models Planner-R1. The official 45/180 train-validation split was merged and reshuffled into 180 training and 45 validation queries, while preserving the easy/medium/hard ratio. We evaluated three single-stage reward configurations with 500 steps and a curriculum regime for 8B and 32B models with 100/300/100 and 50/350/100 steps respectively. 8 rollouts were executed in sglang with a standard ReAct-style agent. We capped trajectories at 30 steps, tool responses at 8,192 tokens, and model outputs at 30,500 tokens. All runs used two nodes (16×H200 GPUs). We used learning rate of 10⁻⁶. Although Qwen3 provided an explicit <think>...

 think>...
 ...
 /think> mode (qwe, 2025b;a), the additional reasoning tokens inflated context length and, in pilots, yielded no gains on task metrics. Full hyperparameters and implementation details are given in Appendix B.3, and decoding and sampling presets are summarized in Appendix B.4.

Out-of-Domain A central concern with task-specific fine-tuning is whether it harms generalization outside the target domain. To probe this, we evaluated our trained models on three complementary suites. All were unseen during training, with evaluation limited to task instructions. (i) NATURAL PLAN (Zheng et al., 2024) (Trip Planning, Meeting Planning, Calendar Scheduling), where tool outputs were provided as context and accuracy was scored by *Exact Match*; we followed the official five-shot prompting protocol. (ii) MULTI-IF (He et al., 2024) (English), a multi-turn instruction-following benchmark derived from IFEval, where the input at turn t concatenated all prior turns ($\leq t-1$); we reported the mean of turn-wise scores. (iii) τ -BENCH (Yao et al., 2024) (retail, function-calling), which measured goal completion against a simulated backend and policy documents; we reported pass@1.

3.2 EVALUATIONS

Table 1 depicts the TravelPlanner results based on Qwen3 (Yang et al., 2025), GPT (OpenAI et al., 2024b; OpenAI, 2025b;a), and our Planner-R1 models using four reward models across five metrics defined in 2. Numbers after \pm indicates 95% confidence intervals.

Base models showed partial competence but struggled with full constraint satisfaction. While stronger base models achieved 99%+ delivery rates and moderate commonsense and hard-constraint coverage, they did not perform well end-to-end. For instance, GPT-5 and GPT-03 achieved final pass rates of 21.2% and 11.3%, respectively. In contrast, the open-weight Qwen3 series performed substantially worse: the 8B model failed entirely, and the 32B model achieved only 0.6% final pass despite a 41.9% delivery rate. This stark disparity underscored that the challenge lay not in planning individual items, but in coordinating tool calls and enforcing all constraints jointly. Prior work (Yao et al., 2023b; Nakano et al., 2021) suggested that prompting alone often underutilized tool feedback, whereas robustness emerged when models interleaved reasoning with actions to query, observe, and update plans. Our findings, as we will see in Section 4, align with this view: base models were able to generate fluent itineraries, but their failures centered on tool sequencing and constraint bookkeeping rather than basic retrieval.

Agentic RL delivered large gains; smaller models were reward-sensitive. RL fine tuning improved both 8B and 32B Qwen3 substantially. Specifically PLANNER-R1-8B using Stage 1 reward and PLANNER-R1-32B using Curriculum reward reached 39.9% and 47% final pass rates respectively. Compared to the 32B model, the 8B model was more sensitive towards sparser rewards: using Stage 2 and Stage 3 rewards resulted in 3/5 and 5/5 model collapses respectively, showing the importance of reward shaping (Ng et al., 1999; dos Santos et al., 2024; Qian et al., 2025b). These collapses explain the large confidence intervals of Stage 2 results. 32B models were more robust. All

²Given the strong Stage 3 performance of larger models, we advanced them more quickly from Stage 1.

reward resulted in 42%+ final pass rate, although increased sparsity resulted in increased variance. 32B model learnt best with Curriculum reward yet the difference were not statistically significant.

Smaller models delivered superior GPU efficiency compared to larger ones. Given the strong performance of Stage 1 training, we extended experiments with high-capacity settings, training the 8B model for 3,000 steps and the 32B model for 2,000 steps. Figure 3 reports results from five independent runs. The left panel, plotted against training steps, shows that both models achieved broadly similar performance trajectories. The right panel, however, plots final pass rate against estimated FLOPs (see Appendix B.6 for details) and reveals a clear efficiency gap. While the 32B model reached 90% of its peak performance (52.3%) at 7.6×10^{20} FLOPs, the 8B model achieved the same level at only 2.1×10^{20} FLOPs, a $3.5 \times$ improvement in efficiency. Although 32B models attained a slightly higher peak accuracy (56.9% vs. 56.4%), this difference was not statistically significant and came with higher variance. A complementary *memory-efficiency* analysis appears in Appendix B.5, where we discuss GPU memory footprint and its implications for agentic RL with long multi-turn contexts. Overall, the FLOPs-based comparison highlights that smaller models are substantially more GPU-efficient for agentic RL training using shaped rewards when data generation is not a limiting factor.

Table 1: Results on the TRAVELPLANNER test set. For Planner-R1 models, we report mean performance with 95% confidence intervals over five runs at 500 training steps. Stage 1–3 denote runs trained exclusively on one stage for 500 steps each. Curriculum uses three phases: for 8B, 100/300/100 steps; for 32B, 50/350/100 steps across Stages 1–3.

Method	Delivery	Commonsense		Hard Co	Final	
	Rate	Micro	Macro	Micro	Macro	Pass Rate
	(%)	(%)	(%)	(%)	(%)	(%)
Qwen3-8B	0.0	0.0	0.0	0.0	0.0	0.0
Qwen3-32B	41.9	27.5	1.7	11.4	7.2	0.6
GPT-o3 (high)	99.6	74.2	14.3	57.7	48.0	11.3
GPT5 (high)	99.8	81.0	23.4	75.4	71.1	21.2
Planner-R1-8B						
Stage1	$99.5 {\pm} 0.8$	$94.8 {\pm} 1.2$	69.0 ± 6.9	61.0 ± 2.6	$46.2 {\pm} 2.5$	39.9 ± 4.3
Stage2	99.9 ± 0.2	80.6 ± 18.2	30.2 ± 51.9	63.4 ± 13.8	$48.6 {\pm} 16.3$	13.3 ± 23.2
Stage3	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
Curriculum	$99.7 {\pm} 0.8$	92.7 ± 3.1	57.9 ± 18.6	53.9 ± 5.7	$38.2 {\pm} 4.2$	27.1 ± 12.6
Planner-R1-32B						
Stage1	99.3 ± 1.6	$95.2 {\pm} 1.6$	70.4 ± 13.4	$74.2 {\pm} 1.4$	$56.4 {\pm} 2.9$	42.3 ± 8.0
Stage2	$91.1 {\pm} 0.5$	87.7 ± 2.2	69.1 ± 14.5	70.0 ± 5.6	55.0 ± 7.6	$44.1 {\pm} 9.4$
Stage3	$99.4 {\pm} 0.9$	$94.7 {\pm} 2.5$	71.9 ± 15.2	60.8 ± 16.6	$48.2 {\pm} 15.1$	$44.3 {\pm} 14.1$
Curriculum	99.1 ± 1.7	95.9 ± 2.5	78.5 ± 7.9	72.1 ± 5.0	55.1 ± 6.2	47.0 ± 6.9

RL fine-tuned models generalized beyond the training domain. Table 2 shows that RL fine-tuned models performed mostly on par with, and often surpassed, their pretrained counterparts across NAT-URAL PLAN (Zheng et al., 2024), MULTI-IF (He et al., 2024), and τ -BENCH (Yao et al., 2024). Blue and red indicate significant improvements and degradations, respectively. After 2,000 steps, both models improved on most metrics, and even at 3,000 steps the 8B model outperformed baselines on five of seven metrics with marginal regressions on two metrics. We attribute this robustness to the JSON-gated output structure, which couples semantics with format and reinforces tool-conditioned behaviors, consistent with prior findings that structured generation improves reliability (Oestreich et al., 2025) and supports generalization to unseen schemas (Liu et al., 2019).

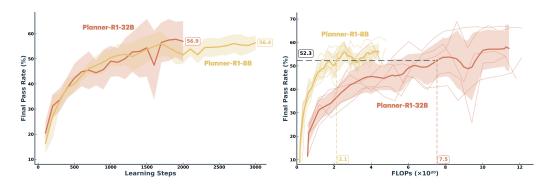


Figure 3: Performance of 8B and 32B Planner-R1 during training based on learning steps (left) and training FLOPS (right). The horizontal dashed line highlights 90% of the maximum average performance of 32B models, while vertical dashed lines show the required FLOPs to reach that performance by both 8B and 32B models.

Table 2: Transferability to external benchmarks without target-domain training (percent). Models are evaluated on NATURAL PLAN, MULTI-IF, and τ -BENCH. Blue = significant improvement over the base model; red = significant degradation from the base model.

Method (Training Steps)	NATURAL PLAN			Multi-IF			au-bench
	Trip	Meeting	Calendar	1st-Turn	2nd-Turn	3rd-Turn	Pass@1
Qwen3-8B	12.9 ± 0.2	82.0 ± 0.0	22.7 ± 0.3	88.9 ± 0.6	82.8 ± 0.6	75.4±0.5	9.5±2.1
Planner-R1-8B (500)	14.0 ± 0.9	83.2 ± 1.1	24.3 ± 0.9	89.4 ± 0.4	$83.5 \!\pm\! 0.7$	76.9 ± 0.6	11.1 ± 0.9
Planner-R1-8B (2000)	14.0 ± 2.1	84.0 ± 0.6	23.2 ± 2.1	89.8 ± 0.4	84.0 ± 0.5	77.2 ± 0.4	12.1 ± 2.3
Planner-R1-8B (3000)	10.7 ± 1.8	84.5 ± 1.3	20.1 ± 2.0	89.8 ± 0.1	83.9 ± 0.4	76.7 ± 0.4	15.1 ± 3.1
Qwen3-32B	11.3 ± 0.0	77.0 ± 0.0	32.2 ± 0.0	89.1 ± 0.3	83.1 ± 0.3	77.1 ± 0.4	28.0 ± 2.2
Planner-R1-32B (500)	15.7 ± 2.2	79.8 ± 1.6	33.2 ± 0.5	88.7 ± 0.2	$83.4 \!\pm\! 0.6$	77.7 ± 0.6	28.7 ± 2.1
Planner-R1-32B (2000)	19.5 ± 1.2	$80.2 \!\pm\! 1.1$	34.4 ± 1.4	89.8 ± 0.3	$84.1 \!\pm\! 0.3$	78.5 ± 0.4	33.9 ± 3.8

4 QUALITATIVE ANALYSIS

To illustrate the effects of RL and model scale, we present a qualitative analysis of Planner-R1 8B and 32B models across training checkpoints, with GPT-5 included as a reference point. We highlight progression in failure modes, tool use, and subreward acquisition for the trained models, and report failure patterns for GPT-5.

Failure Progression Figure 4 shows the progression of the top five failure categories for Planner-R1 8B (left) and 32B (right) models during training³. For hallucination detection, we verify whether the origin city, destination city, attractions, accommodations, and restaurants are present in the corresponding databases. Both models began with high failure rates, particularly on accommodation and cost constraints. For the 8B model, hallucination and cost remain persistent challenges, while all other failures fall below 10% after 800 steps. For the 32B model, accommodation and cost remain dominant errors, with all other failures dropping below 10% by 600 steps. Notably, the 32B model exhibits substantially fewer hallucinations but struggles more with finding accommodations that qualify, for example when the chosen accommodation has a minimum-night requirement and the planned stay must meet this constraint. Another stark observation is the spike at 1,600 steps which can be also observed in Figure 3. In general, we found 32B models with more variations, specially in one run the pass rate dropped from 44% to 26% to 51% impacting the average.

Tool-Use Progression We observed clear improvements in tool-use behavior as training progressed. Early checkpoints of both the 8B and 32B models exhibited poor sequencing, often looping on

³Top categories selected based on their AUC during training. For another lens with top 3 failures at each learning step see Figure A.2

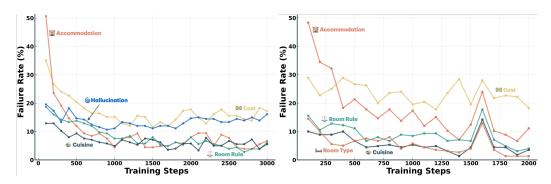


Figure 4: Progression of top 5 failures for 8B (left) and 32B (right) Planner-R1 during training

repetitive calls (e.g., repeatedly invoking the calculator or restaurant tools), which led to incoherent or incomplete plans. As training progressed, both model failures shifted from syntactic to semantic failures: they returned schema valid plans but often failed to call necessary tools to meet the required constraints. With more training, models could often return valid plans. For further details, see the visualizations of tool-call trajectories in Appendix Figures A.3-A.12.

Sub-Reward Progression For both 8B and 32B models, the initial ranking of subrewards from highest to lowest was consistent (see Figure A.1 in Appendix): (1) Schema, (2) Commonsense Micro, (3) Hard Micro, (4) Commonsense Macro, (5) Hard Macro, and (6) Final Pass. As training progressed, success rates increased across all categories, yet this relative ordering remained largely unchanged. This pattern aligned with the λ values defined in Section 2.2, reinforcing our intuition about the relative difficulty of these subrewards and underscoring the role of reward shaping in guiding models through progressively harder objectives.

GPT-5 Behavior Across multiple scenarios, GPT-5 exhibits several recurring error patterns. These include repetition errors, such as selecting the same restaurant or revisiting a city multiple times in violation of commonsense constraints; incomplete plans, where the model fails to return to the departing city or omits key itinerary elements; constraint violations, such as booking fewer than the required minimum hotel nights; and hallucinations, including inventing non-existent hotels or skipping required meals. Illustrative examples of these failure modes are provided in Appendix A.4–A.8.

5 RELATED WORK

Planning Early *chain-of-thought* prompting showed that writing out intermediate steps boosts LLM performance on complex QA and math (Wei et al., 2022; Kojima et al., 2022). Subsequent variants—most notably self-consistency and structured schemes such as Least-to-Most and Plan-and-Solve—further reduce errors by decomposing problems and aggregating diverse solution paths (Wang et al., 2023c; Zhou et al., 2023; Wang et al., 2023b). To address the brittleness of linear chains, search-based methods recast reasoning as combinatorial exploration with lookahead and backtracking, operating over trees (Tree of Thoughts) and graphs (Graph of Thoughts) (Yao et al., 2023a; Besta et al., 2024). Multi-agent formulations extend this idea via division of labor: Chain-of-Agents partitions long inputs among workers while a manager aggregates their outputs (Chen et al., 2024). Decoupling planning from execution further improves robustness: *Plan-and-Act* pairs a planner with an executor and scales supervision via synthetic trajectories, while Iterative Programmatic Planning treats planning as code synthesis (Erdogan et al., 2025; Aravindan et al., 2025). Formal methods offer another angle: Hao et al. (2025) translate planning queries into SAT/SMT specifications solved by external verifiers, achieving rigorous correctness guarantees; in contrast, we keep planning internal to the agent and optimize policies end-to-end with RL. Finally, to reach beyond the context window, recent systems interleave reasoning with targeted search: Search-o1 triggers agentic retrieval under uncertainty and distills evidence via a Reason-in-Documents step, while AI-SearchPlanner trains a lightweight RL planner to trade off query utility and cost, yielding cross-model gains (Li et al., 2025a; Mei et al., 2025).

Agentic RL RL is increasingly used to make tool-use strategic and long-horizon: Search–R1 learns to issue multi-turn web queries during reasoning (Jin et al., 2025), SkyRL trains multi-turn agents inside real software environments (Cao et al., 2025), and ReTool interleaves Python execution within the reasoning loop under outcome-based rewards (Feng et al., 2025). Complementing these, the Tool-Integrated Reasoning line embeds tools directly into the RL objective: ToRL scales toolintegrated RL from base models and reports emergent selective tool invocation with strong math gains (Li et al., 2025b), while ToolRL systematically studies reward design for tool selection and parameterization, showing that carefully shaped rewards with GRPO yield robust improvements over SFT (Qian et al., 2025a). Biomni applies end-to-end reinforcement learning, creating rewards and RL environments tailored to biomedicine, scalably training the agent to carry out research tasks more effectively (Huang et al., 2025). In parallel, large-scale RL fine-tuning (Kimi k1.5, DeepSeek-R1) boosts general reasoning, and Qwen3 introduces dynamic "thinking" vs. "non-thinking" modes to balance depth and latency (Kimi Team et al., 2025; Guo et al., 2025; Yang et al., 2025). Building on this momentum, Kimi K2 emphasizes open agentic intelligence with agentic data synthesis and a joint RL stage (Kimi Team, 2025); GLM-4.5 proposes ARC (Agentic, Reasoning, Coding) foundation models with hybrid thinking/direct modes and RL post-training (GLM-4.5 Team, 2025); and Microsoft's rStar2-Agent explores reliable Python tool use with a Resample-on-Correct strategy for agentic RL (Shang et al., 2025). Most closely related, Chen et al. (2025) introduce LOOP, a data- and memory-efficient variant of PPO that enables reinforcement learning for interactive digital agents directly within stateful, multi-domain environments such as AppWorld.

DISCUSSION AND LIMITATIONS

We showed that agentic RL can substantially enhance planning with tool use, using TRAVELPLANNER as a testbed. Our method achieved state-of-the-art performance among open-weight models and outperformed GPT-5 baselines by 2.7×. A key finding is that smaller models (8B) are especially responsive to shaped, process-level rewards: with only 180 training queries, they reached competitive performance with 32B models while operating at up to 3.5× higher compute-efficiency and 1.5x higher memory-efficiency. Larger models were more robust under sparse signals but gained less from shaping and exhibited greater variability and higher compute demand, underscoring reward design as a key lever for scaling efficiency.

These efficiency gains at 2,000 steps did not come at the expense of robustness. Fine-tuned models maintained or exceeded baseline performance on out-of-domain benchmarks such as MULTI-IF, NATURALPLAN, and τ -BENCH, showing no evidence of overfitting. At 3,000 steps, the 8B model still improved 5 of 7 metrics but regressed on 2, highlighting a potential drawback of excessive fine-tuning. Although our study followed leaderboard rules and avoided prompt engineering, both baselines and our method may benefit from future prompt optimization.

Our study also has limitations. We focus on TRAVELPLANNER, a constrained benchmark, and smaller models may not remain competitive on more complex or open-ended tasks. While 8B models are more FLOP-efficient, larger models can reach higher peak accuracy, which may be necessary in applications where absolute performance is critical. Finally, we explored curriculum learning only in a simple staged form, leaving richer scheduling strategies for future work. Overall, our results highlight reward shaping as central to agentic RL and position smaller models as an efficient, generalizable path forward.

REFERENCES

verl: Volcano engine reinforcement learning for llms. https://github.com/volcengine/verl, 2024. Open-source implementation of the HybridFlow paper.

Qwen3-8b: Switching between thinking and non-thinking mode. https://huggingface.co/Qwen/Qwen3-8B#switching-between-thinking-and-non-thinking-mode, 2025a. Accessed: 2025-09-24.

Qwen3: Think deeper, act faster. https://qwenlm.github.io/blog/qwen3/, 2025b. Accessed: 2025-09-24.

Ashwath Vaithinathan Aravindan, Zhisheng Tang, and Mayank Kejriwal. Code-driven planning in grid worlds with large language models. *arXiv preprint arXiv:2505.10749*, 2025.

- Peter Belcak, Greg Heinrich, Shizhe Diao, Yonggan Fu, Xin Dong, Saurav Muralidharan, Yingyan Celine Lin, and Pavlo Molchanov. Small language models are the future of agentic ai, 2025. URL https://arxiv.org/abs/2506.02153.
 - Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michał Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: solving elaborate problems with large language models. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence*, pp. 17682–17690, 2024.
 - Shiyi Cao, Sumanth Hegde, Dacheng Li, Tyler Griggs, Shu Liu, Eric Tang, Jiayi Pan, Xingyao Wang, et al. Skyrl-v0: Train real-world long-horizon agents via reinforcement learning, 2025. arXiv:2502.02789.
 - Kevin Chen, Marco Cusumano-Towner, Brody Huval, Aleksei Petrenko, Jackson Hamburger, Vladlen Koltun, and Philipp Krähenbühl. Reinforcement learning for long-horizon interactive llm agents. arXiv preprint arXiv:2502.01600, 2025. URL https://arxiv.org/abs/2502.01600.
 - Tianjun Chen et al. Chain of agents: Toward long context reasoning in language models via multiagent cooperation. *arXiv preprint arXiv:2406.02818*, 2024.
 - Joao dos Santos et al. Revisiting sparse rewards for goal-reaching reinforcement learning. *arXiv*:2407.00324, 2024. URL https://arxiv.org/abs/2407.00324.
 - Lutfi Eren Erdogan, Nicholas Lee, Sehoon Kim, Suhong Moon, Hiroki Furuta, Gopala Anumanchipalli, Kurt Keutzer, and Amir Gholami. Plan-and-act: Improving planning of agents for long-horizon tasks. *arXiv preprint arXiv:2503.09572*, 2025.
 - Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, et al. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*, 2025.
 - Aryo Pradipta Gema, Alexander Hägele, Runjin Chen, Andy Arditi, Jacob Goldman-Wetzler, Kit Fraser-Taliente, Henry Sleight, Linda Petrini, Julian Michael, Beatrice Alex, Pasquale Minervini, Yanda Chen, Joe Benton, and Ethan Perez. Inverse scaling in test-time compute, 2025. URL https://arxiv.org/abs/2507.14417.
 - GLM-4.5 Team. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models, 2025.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.
 - Yilun Hao, Yongchao Chen, Yang Zhang, and Chuchu Fan. Large language models can solve real-world planning rigorously with formal verification tools, 2025. URL https://arxiv.org/abs/2404.11891.
 - Yun He, Di Jin, Chaoqi Wang, Chloe Bi, Karishma Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu Xu, Hongjiang Lv, et al. Multi-if: Benchmarking llms on multi-turn and multilingual instructions following. *arXiv preprint arXiv:2410.15553*, 2024.
 - Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani, Ryan Li, Lin Qiu, Gavin Li, Junze Zhang, Di Yin, Shruti Marwaha, Jennefer N. Carter, Xin Zhou, Matthew Wheeler, Jonathan A. Bernstein, Mengdi Wang, Peng He, Jingtian Zhou, Michael Snyder, Le Cong, Aviv Regev, and Jure Leskovec. Biomni: A general-purpose biomedical ai agent. bioRxiv, 2025. doi: 10.1101/2025.05.30.656746. URL https://www.biorxiv.org/content/early/2025/06/02/2025.05.30.656746.
 - Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey, 2024. URL https://arxiv.org/abs/2402.02716.

- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan O. Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.
- Kimi Team. Kimi k2: Open agentic intelligence, 2025.

547

551

552

553

554

555

556

558

559

561

562 563

564 565

566

567

568

569

570

571 572

573

574

575

576

577 578

579

580

581

582

583

584

585

586

588

589

- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, et al. Kimi k1.5: Scaling reinforcement learning with large language models. *arXiv preprint arXiv:2501.12599*, 2025.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems* (NeurIPS), pp. 22199–22213, 2022.
 - Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*, 2025a.
 - Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl. *arXiv preprint arXiv:2503.23383*, March 2025b. doi: 10.48550/arXiv.2503.23383. URL https://arxiv.org/abs/2503.23383.
 - Tianyu Liu, Furu Wei, and Ming Zhou Wang. Table-to-text generation with unseen schemas. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
 - Lang Mei, Zhihan Yang, and Chong Chen. Ai-searchplanner: Modular agentic search via pareto-optimal multi-objective reinforcement learning. *arXiv* preprint arXiv:2508.20368, 2025.
 - Reiichiro Nakano et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv:2112.09332*, 2021. URL https://arxiv.org/abs/2112.09332.
 - Andrew Y. Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, pp. 278–287, 1999. URL https://people.eecs.berkeley.edu/~pabbeel/cs287-fa09/readings/NgHaradaRussell-shaping-ICML1999.pdf.
 - Julian Oestreich, Lydia Müller, et al. Evaluating structured decoding for text-to-table generation: Evidence from three datasets. *arXiv* preprint arXiv:2508.15910, 2025.
 - OpenAI. Gpt-5 system card. https://cdn.openai.com/gpt-5-system-card.pdf, 2025a. Accessed: 2025-09-20.
 - OpenAI. Openai o3 and o4-mini system card. https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf, 2025b. Accessed: 2025-09-20.
 - OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng,

595

596

597

600

601

602

603

604

605

607

608

609

610

612

613

614

615

616

617 618

619

620

621

622

623

625

626

627

630

631

632

633

634

635

636

637

638

639

640

641

642

645

646

Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. Openai ol system card, 2024a. URL https://arxiv.org/abs/2412.16720.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl,

653

655

657

661

663

664

666 667

668

669 670

671

672

673

674

675

676

677

678

679

680 681

682

683

684

685

686

687

688

689

690

691

692 693

694

696

697

699

700

648 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra 649 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, 650 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, 652 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-654 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan 656 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-658 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming 659 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao 660 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024b. URL https://arxiv.org/abs/2303.08774. 662

- Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. arXiv preprint arXiv:2504.13958, April 2025a. doi: 10.48550/arXiv.2504.13958. URL https://arxiv.org/abs/2504. 13958.
- Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tür, and Heng Ji. Toolrl: Reward is all tool learning needs. arXiv:2504.13958, 2025b. URL https://arxiv.org/abs/2504.13958.
- Ning Shang, Yifei Liu, Yi Zhu, Li Lyna Zhang, Weijiang Xu, Xinyu Guan, Buze Zhang, Bingcheng Dong, Xudong Zhou, Bowen Zhang, Ying Xin, Ziming Miao, Scarlett Li, Fan Yang, and Mao Yang. rstar2-agent: Agentic reasoning technical report, 2025. Microsoft Research.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402. 03300.
- Noel Shinn, Federico Cassano, Ashwin Gopinath, Karthik R. Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In Advances in Neural Information Processing Systems (NeurIPS), 2023.
- Parshin Shojaee*, Iman Mirzadeh*, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity, 2025. URL https://ml-site.cdn-apple. com/papers/the-illusion-of-thinking.pdf.
- Guanhua Wang, Yuzhuo Xie, Yuanzhi Jiang, Ajay Mandlekar, Chongjie Xiao, Yuke Zhu, Lifeng Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. arXiv preprint arXiv:2305.16291, 2023a.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In Annual Meeting of the Association for Computational Linguistics (ACL), pp. 2609– 2634, 2023b.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In International Conference on Learning Representations (ICLR), 2023c.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems (NeurIPS), pp. 24824–24837, 2022.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: a benchmark for real-world planning with language agents. In *Proceedings* of the 41st International Conference on Machine Learning, pp. 54590-54613, 2024.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025. https://arxiv.org/abs/2505.09388.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023a.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023b.

Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. *τ*-Bench: A Benchmark for Tool–Agent–User Interaction in Real-World Domains. *arXiv preprint arXiv:2406.12045*, 2024.

Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang, Xinyun Chen, Minmin Chen, Azade Nova, Le Hou, Heng-Tze Cheng, Quoc V Le, Ed H Chi, et al. Natural plan: Benchmarking llms on natural language planning. *arXiv preprint arXiv:2406.04520*, 2024.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H. Chi. Least-to-most prompting enables complex reasoning in large language models. In *International Conference on Learning Representations (ICLR)*, 2023.

A SYSTEM-LEVEL OPTIMIZATIONS

Overview. RL training with large-scale LLMs requires co-locating both training and inference engines on the same set of GPUs. This dual demand creates severe memory pressure, often leading to out-of-memory (OOM) errors when switching between training and rollout phases. To address this, we integrated memory management techniques into our RL Pipelines.

Multi-Stage Awake Memory Management. In VERL, reinforcement learning (RL) training is conducted in Colocate Mode, where both the training engine (e.g., FSDP) and inference engine (e.g., SGLang) share the same GPU resources. A major bottleneck arises when transferring weights from the training engine to the inference engine: model parameters must be copied from FSDP into SGLang, often causing out-of-memory (OOM) failures under high memory pressure.

To address this, we extended the *Sleep/Awake* mechanism in SGLang and introduced the *Multi-Stage Awake* strategy for fine-grained memory management during rollouts. Instead of a single monolithic resume, memory resumption is divided into multiple stages:

- 1. Load training model weights into GPU memory.
- 2. Resume inference model weights at preserved virtual addresses.
- 3. Synchronize weights between training and rollout engines.
- 4. Offload training model weights back to CPU.
- 5. Resume the KV cache region for rollout execution.

This staged approach minimizes memory waste and prevents fragmentation. Our empirical results show that it provides two key benefits:

- Enables training of larger models: With the same KV cache ratio, our approach reduces peak GPU memory by 20–23%, which unblocks stable training of a 32B-parameter model on 8×H200 GPUs even at higher cache ratios (0.8, 0.85, and up to 0.9). Without Multi-Stage Awake, training consistently ran out of memory beyond 0.7.
- **Improves throughput**: For the same model size, our method allows a larger KV cache ratio to be used, directly improving inference throughput. While throughput gains are workload-dependent and not easily comparable across setups, our experiments show that increasing the ratio from 0.7 to 0.9 leads to significant improvements in rollout efficiency.

B IMPLEMENTATION DETAILS

B.1 SYSTEM PROMPT

756

758

759

760

761

762

763

764 765 766

767 768

769 770

771

772773774

775

776

777 778

779

781

782

783

784 785

786

787 788

789

791

792 793 794

795

796

797 798

799 800

801

802

803

804

805

806

807

808

We include the jinja template of our full system prompt used for Planner-R1 during training/evaluation.

```
You are a helpful travel assistant that plans detailed travel
  itineraries by calling external functions (tools). You have access
   to the following tools and must use them as needed to gather
   accurate, up-to-date information.
# Behavior Guidelines
- If a task requires multiple steps or tools, proceed step by step,

→ calling ONE TOOL per turn.

- Never assume details-always verify all information using tools.
- When you have gathered sufficient information to finalize the plan,
\hookrightarrow respond with an <answer> block with the final itinerary in valid
   JSON format.
# Tool Usage Rules
- Do not repeat the same tool call with identical arguments.
- Always provide complete and correct function arguments.
# Final Plan Format
Once all necessary information is collected, respond with the final
→ plan:
<answer>
      // Day 1 plan following schema
    },
      // Day 2 plan following schema
      ... additional days
</answer>
**IMPORTANT CONSTRAINTS**
- The <answer> must contain ONLY valid JSON, strictly following the
→ plan_schema.
- Do not include any explanatory text inside the <answer> block.
- Do not output <answer> until all needed tool calls are completed.
# Final Plan Schema
Each element in the <answer> JSON array should represent a single day
\rightarrow of the trip and follow this schema exactly:
···json
```

```
{{ plan_schema }}
```

B.2 PLAN JSON SCHEMA

810

816 817

818

819

820

The final itinerary must be a JSON *array* of per-day objects. Each day object is validated against the schema below. This structured contract doubles as a checklist (ensuring coverage of all required fields) and enables automatic reward gating.

```
821
822
             "type": "object",
823
             "required": [
824
                 "days", "city", "transportation", "attraction",
                 → "accommodation", "breakfast", "lunch", "dinner"
825
826
             "properties": {
                 "days": {
828
                     "description": "The day number of the plan starting from
829
                     "type": "integer"
830
831
                 "city": {
832
                     "description": "Can be a city name string if no transfer
833
                     → is needed, or an dict with 'from' and 'to' keys that
834
                     \rightarrow indicates the origin and destination city.",
                     "oneOf": [
835
                         {"type": "string"},
836
837
                              "type": "object",
838
                              "required": ["from", "to"],
839
                              "properties":
                                  "from": {"type": "string"},
840
                                  "to": {"type": "string"}
841
842
                              "additionalProperties": false
843
                         }
844
                     1
845
                 "transportation": {
846
                     "description": "Either '-' if no transportation is needed,
847
                     \rightarrow or an object describing the transportation details.
848
                     \hookrightarrow Instead of total cost, use per person price for flight
849
                      → and per vehicle cost for taxi/self-driving as the
                      ⇔ cost.",
850
                     "oneOf": [
851
                         {
852
                              "type": "string",
853
                              "const": "-"
854
                          },
855
                              "type": "object",
856
                              "required": ["mode", "from", "to", "duration",
857
                              → "distance", "cost"],
858
                              "properties": {
859
                                  "mode": {
                                      "type": "string",
860
                                      "enum": ["flight", "taxi",
861
                                      862
                                      "description": "Type of transportation."
                                  },
```

```
864
                                  "from": {"type": "string", "description":
865
                                   → "Origin city"},
                                  "to": {"type": "string", "description":
867
                                  → "Destination city"},
                                  "duration": {"type": "string", "description":
868
                                  → "Transportation duration"},
869
                                  "distance": {"type": "string", "description":
870
                                  → "Distance of the trip"},
871
                                  "cost": {"type": "integer", "description":
872
                                  → "Cost of the transportation"},
873
                                  "flight_number": {"type": "string",
874
                                  → "description": "Flight number (for flights
875
                                   \hookrightarrow only)"},
876
                                  "departure_time": {"type": "string",
877
                                   → "description": "Flight departure time"},
878
                                  "arrival_time": {"type": "string",
                                  → "description": "Flight arrival time"}
879
880
                              "additionalProperties": false
881
                          }
882
                     ]
883
                 },
                 "attraction": {
884
                     "description": "A list of attraction names planned for the
885

→ day, or '-' if no attractions are planned.",
886
                     "oneOf": [
887
                          {"type": "string", "const": "-"},
888
                              "type": "array",
889
                              "items": {"type": "string"},
890
                              "minItems": 1
891
                          }
892
                     ]
893
                 "accommodation": {
894
                     "description": "The name of the accommodation for today.
895
                      \rightarrow '-' if no accommodation is needed.",
896
                     "type": "string"
897
898
                 "breakfast": {
                     "description": "The name of the breakfast restaurant for
899

→ today. '-' if no breakfast is planned.",
900
                     "type": "string"
901
902
                 "lunch": {
903
                     "description": "The name of the lunch restaurant for

→ today. '-' if no lunch is planned.",
904
                     "type": "string"
905
906
                 "dinner": {
907
                     "description": "The name of the dinner restaurant for

→ today. '-' if no dinner is planned.",
908
                     "type": "string"
909
910
911
             "additionalProperties": false
912
         }
913
914
```

919 920

921

922

923

924

B.3 TRAINING, VALIDATION, AND EVALUATION SETUP

RL framework and resources. We train with VERL using **GRPO** on **2 nodes** with **8 GPUs/node** (16 H200 GPUs total). Rollouts use **sglang** with a multi-turn, tool-augmented agent (ReAct-style).

Training configuration. Stage 1/2/3 share the same VERL configuration; only the reward weights differ by stage (see Sec. 2). File paths are pseudonymized for readability.

```
925
926
         # verl + GRPO. Stage-agnostic; change reward weights per stage.
927
        actor_rollout_ref:
928
          actor:
             strategy: fsdp
929
             ppo_mini_batch_size: 8
930
             ppo_micro_batch_size: null
931
             ppo_micro_batch_size_per_gpu: 1
932
             use_dynamic_bsz: false
             ppo_max_token_len_per_gpu: 16384
933
             clip_ratio: 0.2
934
             clip ratio low: 0.2
935
             clip_ratio_high: 0.2
936
             policy_loss:
937
               loss_mode: vanilla
               clip_cov_ratio: 0.0002
938
               clip_cov_lb: 1.0
939
               clip_cov_ub: 5.0
940
               kl_cov_ratio: 0.0002
941
               ppo_kl_coef: 0.1
942
             clip_ratio_c: 3.0
             loss_agg_mode: token-mean
943
             entropy_coeff: 0
944
             use_kl_loss: false
945
             use_torch_compile: true
946
             kl_loss_coef: 0.001
947
             kl_loss_type: low_var_kl
             ppo_epochs: 1
948
             shuffle: false
949
             optim:
               lr: 1.0e-06
951
               lr_warmup_steps_ratio: 0.0
952
               {\tt total\_training\_steps:} \ -1
               weight_decay: 0.01
953
               lr_warmup_steps: -1
954
               min_lr_ratio: 0.0
955
               num_cycles: 0.5
956
               warmup_style: constant
957
             grad clip: 1.0
             ulysses_sequence_parallel_size: 1
958
             entropy_from_logits_with_chunking: false
959
             entropy_checkpointing: false
960
             fsdp_config:
961
               wrap policy:
962
                 min_num_params: 0
               param_offload: true
963
               optimizer_offload: true
964
               offload policy: false
965
               reshard_after_forward: true
966
               fsdp_size: -1
967
               forward_prefetch: false
           rollout:
968
             name: sglang
969
             mode: async
970
             temperature: 1.0
             top_k: -1
```

```
972
             top_p: 1
973
             prompt_length: 2268
974
             response_length: 30500
975
             dtype: bfloat16
976
             gpu_memory_utilization: 0.6
             ignore_eos: false
977
             enforce_eager: true
978
             free_cache_engine: true
979
             tensor_model_parallel_size: 4
980
            max_num_batched_tokens: 8192
981
            max_model_len: null
             max_num_seqs: 1024
982
             log_prob_micro_batch_size: null
983
             log_prob_micro_batch_size_per_gpu: 32
984
             log_prob_use_dynamic_bsz: false
985
             log_prob_max_token_len_per_gpu: 16384
986
             disable_log_stats: true
             do_sample: true
987
            n: 8
988
            multi stage wake up: false
989
             val_kwargs:
990
               top_k: -1
991
               top_p: 1.0
               temperature: 0
992
               n: 1
993
               do_sample: false
994
             multi turn:
995
               enable: true
996
               max_assistant_turns: 30
               tool_config_path: ${PROJ_ROOT}/config/tool_config.yaml
997
               max_user_turns: 30
998
               max_parallel_calls: 1
999
               max tool response length: 8192
1000
               tool_response_truncate_side: right
1001
               interaction_config_path: null
               completion_callback: null
1002
               use_inference_chat_template: false
1003
               tokenization_sanity_check_mode: strict
1004
               format: hermes
1005
             calculate_log_probs: false
1006
             agent:
               num_workers: 8
1007
               agent_loop_config_path: ${PROJ_ROOT}/config/agent_loops.yaml
1008
               custom_async_server:
1009
                 path: null
1010
                 name: null
1011
             update_weights_bucket_megabytes: 512
             enable_chunked_prefill: true
1012
             load_format: dummy_dtensor
1013
             layered_summon: false
1014
             enable thinking: false
1015
          hybrid_engine: true
1016
          model:
            path: Qwen/Qwen3-{8B|32B} # base model
1017
             custom chat template: null
1018
             use_shm: false
1019
             external_lib: null
1020
             override_config: {}
1021
             enable_gradient_checkpointing: true
             enable_activation_offload: false
1022
             use_remove_padding: true
1023
             target_modules: all-linear
1024
             exclude_modules: null
1025
             use_liger: false
```

```
1026
             use_fused_kernels: false
1027
             fused kernel options:
1028
               impl_backend: torch
1029
             trust_remote_code: false
1030
        trainer:
          balance batch: true
1031
           total_epochs: 300
1032
          total_training_steps: 3000
1033
          profile_steps: null
1034
          logger:
1035
             - mlflow
           log_val_generations: 0
1036
          rollout_data_dir: null
1037
          nnodes: 2
1038
          n_gpus_per_node: 8
1039
           save_freq: 100
1040
           esi_redundant_time: 0
          resume_mode: auto
1041
          val before train: true
1042
          val_only: false
1043
          test_freq: 50
1044
           critic_warmup: 0
          default_hdfs_dir: null
1045
           del_local_ckpt_after_load: false
1046
          max_actor_ckpt_to_keep: null
1047
          max_critic_ckpt_to_keep: null
1048
          ray_wait_register_center_timeout: 300
1049
          device: cuda
1050
          use_legacy_worker_impl: auto
        data:
1051
          tokenizer: null
1052
          use_shm: false
1053
          train_files: ${PROJ_ROOT}/data/train.parquet
1054
          val_files: ${PROJ_ROOT}/data/test.parquet
1055
          prompt_key: prompt
          reward_fn_key: data_source
1056
          max_prompt_length: 2268
1057
          max_response_length: 30500
1058
          train_batch_size: 16
1059
          val_batch_size: 64
1060
          return_raw_input_ids: false
           return_raw_chat: true
1061
           return_full_prompt: false
1062
          shuffle: true
1063
          dataloader num workers: 8
1064
          validation shuffle: false
1065
           filter_overlong_prompts: true
          filter_overlong_prompts_workers: 1
1066
          truncation: error
1067
          image_key: images
1068
          video key: videos
1069
          trust_remote_code: false
1070
        custom_reward_function:
          path: ${PROJ_ROOT}/rewards_v3.py
1071
          name: compute_score
1072
        algorithm:
1073
           gamma: 1.0
1074
           lam: 1.0
1075
          adv_estimator: grpo
          norm_adv_by_std_in_grpo: true
1076
          use_kl_in_reward: false
1077
          kl_penalty: kl
1078
          kl_ctrl:
1079
             type: fixed
```

1080
1081
 kl_coef: 0.001
 horizon: 10000
 target_kl: 0.1
1083
 use_pf_ppo: false
 pf_ppo:
 reweight_method: pow
 weight_pow: 2.0

Listing 1: VERL/GRPO configuration (paths pseudonymized). Stage 2/3 reuse this config with stage-specific reward weights.

Evaluation protocol. We evaluate on the TRAVELPLANNER *official test set* by reusing the VERL *validation* pipeline to keep decoding/sampling consistent with validation:

- 1. Point the VERL validation loader to the official test split (same sampler settings as validation).
- 2. Run validation to dump trajectories locally (tool calls, responses, final answers) as JSONL.
- 3. **Post-process** each final answer: validate against the plan schema (App. B.2), enforce JSON-gated output, and **convert** to the leaderboard's submission format.
- 4. **Upload** the converted file to the TRAVELPLANNER leaderboard; report Delivery, micro/macro commonsense, micro/macro hard, and Final.

B.4 DECODING AND SAMPLING SETTINGS

We standardize decoding across training and evaluation to isolate the effect of learning. Table 3 summarizes the presets we use for different contexts; "Common runtime limits" apply to all scenarios unless noted. For Validation/Test we reuse VERL's validation path on the official TP test split (Sec. B.3).

Common runtime limits.

• Max response tokens: 30,500 (response_length)

• Max tool response tokens: 8,192 (max_tool_response_length)

• Agent turns cap: 30 assistant turns; 30 tool turns

• Tool-call cap: 30 calls

Table 3: Decoding presets by context.

Context	do_sample	Temp	Top- <i>p</i>	Top-k	n
Training	true	1.0	1.0	-1	8
Validation / Test	false	0.0	1.0	-1	1
NaturalPlan	true	1.0	1.0	-1	1
Multi-IF	false	0.7	0.8	20	1
au-bench	false	0.6	0.95	20	1

B.5 GPU MEMORY FOOTPRINT AND PRACTICAL EFFICIENCY

In the same 2-node/16-GPU training setup, PLANNER-R1-8B uses approximately ∼60 GB of GPU memory per device, whereas PLANNER-R1-32B requires ≥90 GB per device. This difference has practical consequences: the 8B configuration runs comfortably on H100s, while the 32B configuration necessitates higher-memory accelerators (e.g., H200). The gap is especially relevant for agentic RL, where multi-turn interactions and tool feedback produce long contexts and large key–value (KV) caches during rollouts, amplifying the memory pressure beyond the update phase.

B.6 ESTIMATING TRAINING FLOPS FROM VERL'S MFU

MFU in VERL (what it is). VERL reports a *model FLOPs utilization* (MFU): the fraction of the cluster's "promised" peak compute achieved during *policy updates*. Internally it is computed per

update as

$$ext{MFU} \, = \, rac{f_{
m ach} \, E}{f_{
m peak} \, W} \, , \qquad f_{
m ach} = rac{ ext{FLOPs}_{
m update}}{t_{
m actor}} \, ,$$

where E is the number of GRPO epochs per batch, W is the number of GPUs (world size), $f_{\rm peak}$ is the promised FLOPs rate per GPU used by VERL in its MFU denominator, $t_{\rm actor}$ is the time spent in the parameter-update step, and FLOPs $_{\rm update}$ is the per-step FLOPs consumed by that update. The in-tree FLOPs counter aggregates forward + backward over all layers/tokens.

Reconstruction used in this paper. Solving for ${\rm FLOPs_{update}}$ gives

$$\boxed{\text{FLOPs}_{\text{update}} = \text{MFU} \times f_{\text{peak}} \times W \times \frac{t_{\text{actor}}}{E}}$$

In our runs E=1 and W=16. We set $f_{\rm peak}=9.89\times 10^{14}$ FLOPs/s per GPU—the same constant VERL uses for MFU—so the reconstruction matches its calculation.

Practical proxy for $t_{\rm actor}$. VERL does not log $t_{\rm actor}$ each step, but it logs update_policy_time $(t_{\rm policy})$, which equals the actor update plus brief offload/reload bookkeeping. Because the parameter update dominates, we use

$$t_{\rm actor} \, \approx \, t_{\rm policy} \quad \Rightarrow \quad {\rm FLOPs_{update}} \, \approx \, {\rm MFU} \, \times \, f_{\rm peak} \, \times \, W \, \times \, t_{\rm policy} \, .$$

This yields a *slight upper bound* (since $t_{\rm policy} \geq t_{\rm actor}$); spot checks in our regime found the gap within $\sim 3\%$.

From per-step to cumulative FLOPs. We compute $FLOPs_{update}$ per step from MFU and t_{policy} , then sum over steps:

$$\text{FLOPs}_{1:T} \ = \ \sum_{t=1}^{T} \ \text{MFU}_t \ f_{\text{peak}} \ W \ t_{\text{policy},t} \,.$$

These cumulative totals back the FLOPs-accuracy curves in Sec. 3.2.

Scope: what is counted and what is not. Our accounting *includes only* the parameter-update compute (forward + backward). It *excludes* (i) the rollout engine's generation compute and (ii) the *reference log-prob* pass (both of which VERL does not report MFU/FLOPs for). Under our settings (long responses, multi-sample trajectories), rollout consists of many forward passes, while the update consists of forward+backward over similar tokens; thus their compute is typically of the same *order of magnitude*, but exact ratios depend on response length, batching, and n (number of sampled trajectories). A precise FLOPs tally for rollout and reference log-prob computation is left to future work.

C CASE STUDIES

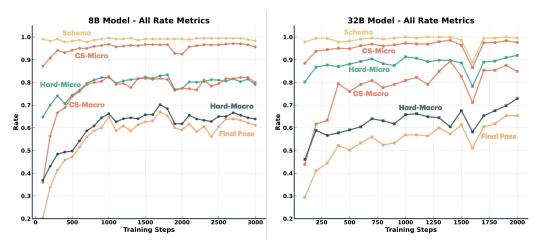


Figure A.1: Progression of six sub rewards for 8B and 32B Planner-R1 during training

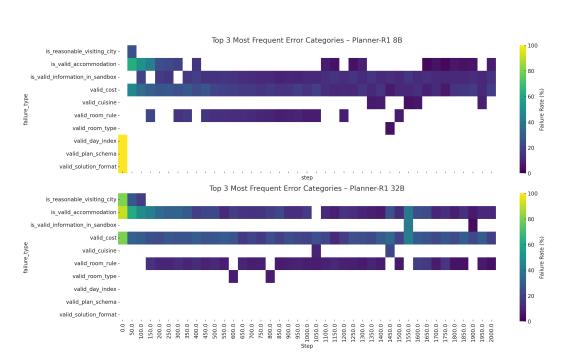


Figure A.2: Top 3 Most Frequent Error Categories for the Planner-R1 models. The upper heatmap shows the 8B model and the lower heatmap shows the 32B model. Both plots visualize, at each training step across five runs (up to 2000), the three most frequent failure categories and their relative rates. Rows are aligned so categories match across models; blank cells indicate that a category did not appear in the top three for that model at that step.

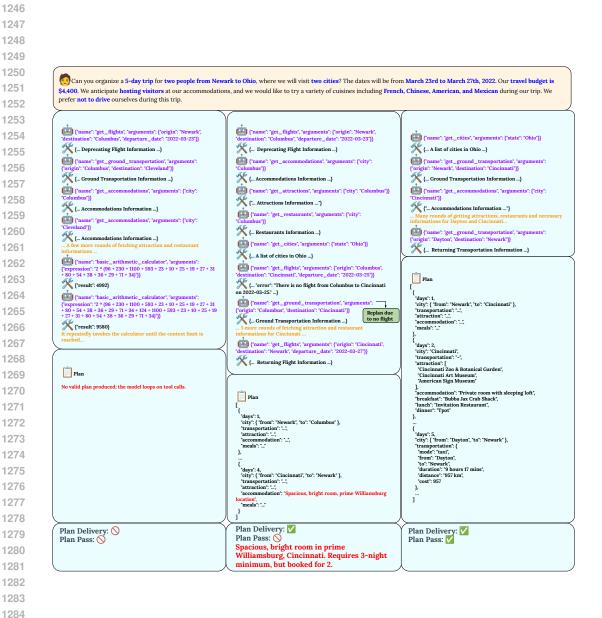


Figure A.3: Model tool-call trajectories across checkpoints. The base model (left) loops on repetitive tool calls once the context is saturated. After 100 training steps (middle), the model produces a coherent travel plan but fails to satisfy all constraints. By 500 steps with the 32B Planner-R1 checkpoint (right), the model successfully generates a valid plan that meets all requirements.

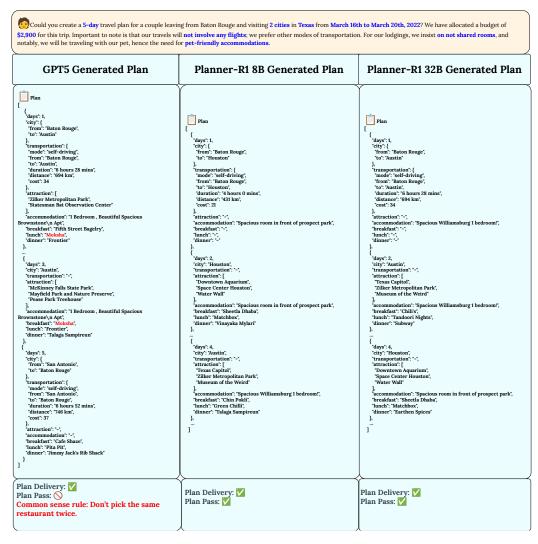


Figure A.4: The GPT-5 model (left) failed to avoid selecting the same restaurant twice, thus violating the common-sense rule. In contrast, the Planner-R1 8B model (middle) and Planner-R1 32B model (right) both generated plans that satisfied all requirements.

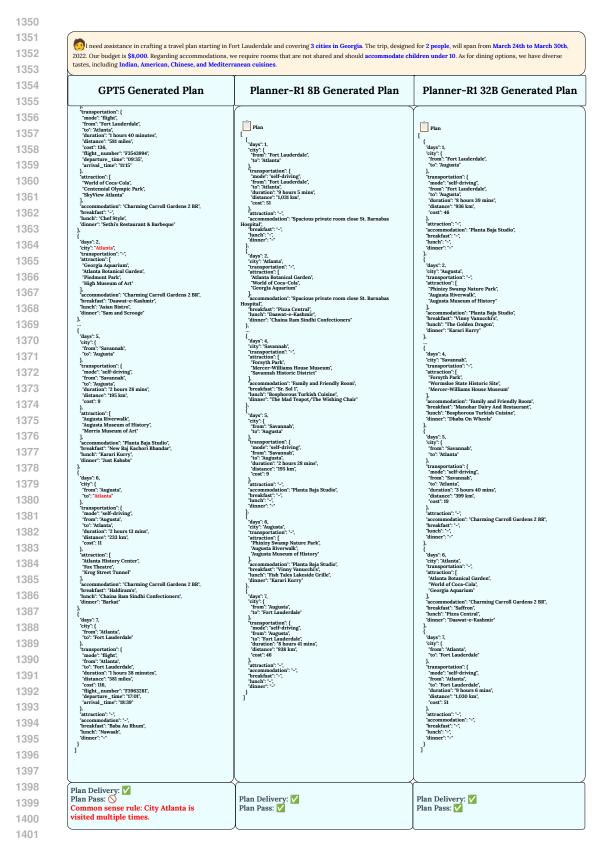


Figure A.5: The GPT-5 model (left) selected the same restaurant twice on different dates, whereas the Planner-R1 8B (middle) and Planner-R1 32B (right) models produced plans that satisfied all requirements.

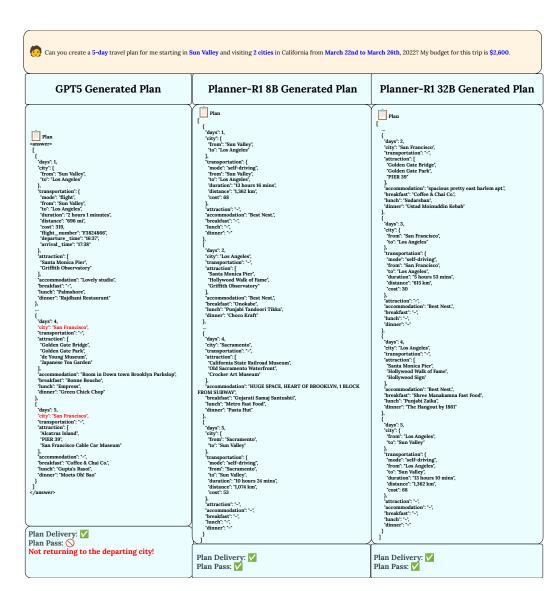


Figure A.6: The GPT-5 model (left) failed to return to the departing city, whereas the Planner-R1 8B (middle) and Planner-R1 32B (right) models produced plans that satisfied all requirements.

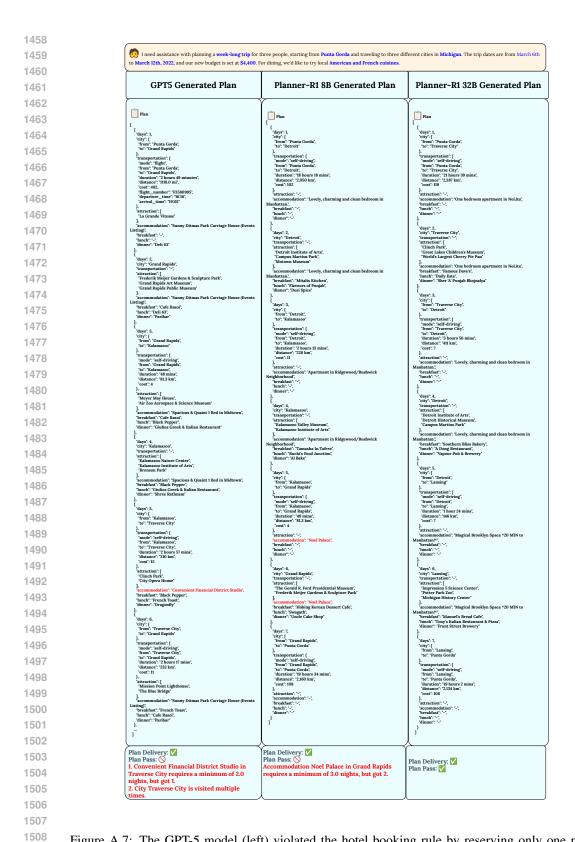


Figure A.7: The GPT-5 model (left) violated the hotel booking rule by reserving only one night instead of the required minimum of two, and also erred by visiting Traverse City multiple times, failing the common-sense requirement. The Planner-R1 8B model (middle) likewise failed the accommodation requirement, while only the Planner-R1 32B model (right) satisfied all requirements.

1512 🧑 Could you create a travel plan for a party of 8 departing from Fayetteville and heading to New York 1513 for 5 days, which will cover 2 cities between March 25th and March 29th, 2022? We have a maximum budget of \$6,900. We are particularly interested in experiencing American and Mexican cuisines during 1515 1516 1517 **GPT5 Generated Plan** Planner-R1 8B Generated Plan Planner-R1 32B Generated Plan 1518 1519 Plan 1520 1521 Plar Plan 1522 1523 1524 1525 'duration": "11 hours 11 mins", 'distance": "1,169 km", 1526 1527 n": "Modern Brooklyn oasis (PRIVATE ROOM)". on": "Williamsburg Gem: Sleep up to 5". 1528 days': 2, city: "Buffalo", transportation": '--'; attraction": ["The Buffalo Zoo", "Buffalo and Erie County Botanical Gardens", "Buffalo AKG Art Museum" nch": "Gurgaon Hights" nner": "G Dot" 1530 1531 mmodation': "Williamsburg Gem: Sleep up to 5', kfast': "Red Mango', h': "Tibbys New Orleans Kitchen', er': "Shokitini" 1532 attraction": ["Statue of Liberty", "9/11 Memorial & ! ommodation": "Modern Brooklyn oasis akfast": "Gurgaon Hights", ch": "G Dot", ner": "Aryan's Rajasthani Pyaz Ki Kach 1533 Brooklyn Bridge" The High Line 1534 1535 inch": "Amchur", inner": "Gurgaon Hights" 1536 "transportation": {
"mode": "self-driving",
"from": "New York",
"to": "Boston",
"duration": "3 hours 41:
"distance": "346 km",
"cost": 17 1537 1538 1539 ion": "A Contemporary Homelike Stay in the Best of , ion": "Sunlight + Space on Eastern Parkway" 1540 reakfast": "-". 1541 1542 1543 attraction": [
"Top of The Rock",
"One World Observatory",
"SUMMIT One Vanderbilt" ommodation": "Sunlight + akfast": "Irish Democrat", ommodation": "A Contemporary Homelike Stay in the Best of 1545 ", "breakfast": "Gurgaon Hights", "lunch": "G Dot", "dinner": "Rambhog" ransportation": "-", ttraction": [Empire State Building", Rockefeller Center", Radio City Music Hall" 1547 city": { "from": "Boston", "to": "Fayetteville" 1548 1549 "breakfast": "-", "lunch": "G Dot", 1550 1551 days": 5, city": "New York", 'transportation": "-", attraction": ["One World Observatory", "The Battery", "Flatiron Building" 1552 1553 1554 1555 1556 1557 Plan Delivery: 🔽 Plan Delivery: ✓ Plan Pass: ○ Plan Delivery: V Plan Pass: V 1558 Plan Pass: 🔽 1. Missing meals on day 2.
2. Accommodation Magical rowhouse and 1559

Figure A.8: In this query, the user plans a trip for a party of 8. The GPT-5 model (left) missed meals on Day 2 and hallucinated non-existent hotel names. The Planner-R1 8B model (middle) generated a plan that exceeded the \$6,900 budget, while only the Planner-R1 32B model (right) satisfied all requirements.

garden is not presented in database

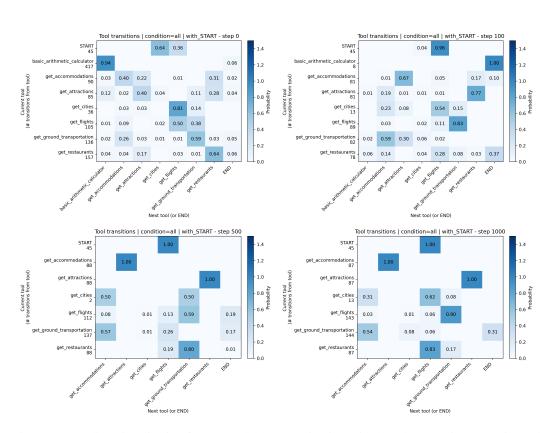


Figure A.9: Policy visualization for 8B model across 45 trajectories based on previous (y-axis) and next (x-axis) tool calls across various steps of learning: $\{0, 100, 500, 1000\}$. As learning progresses, the policy becomes more deterministic.

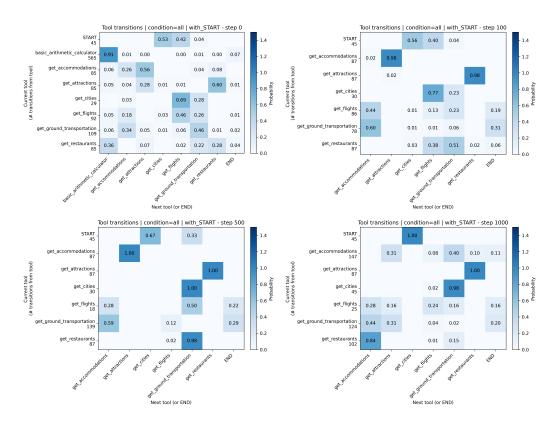


Figure A.10: Policy visualization 32B model across 45 trajectories based on previous (y-axis) and next (x-axis) tool calls across various steps of learning: $\{0, 100, 500, 1000\}$. As learning progresses, the policy becomes more deterministic.



Figure A.11: Tool call sequence behavior as 8B training progresses. The base model (leftmost) repeatedly invoked the calculator and restaurant tools until reaching the rollout cap (30 turns), exhibiting poor tool-use behavior as context grew. With longer training, the model developed more consistent and structured patterns for tool calls.

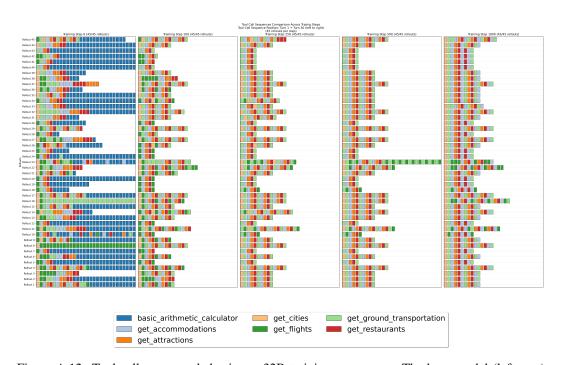


Figure A.12: Tool call sequence behavior as 32B training progresses. The base model (leftmost) repeatedly invoked the calculator until reaching the rollout cap (30 turns), exhibiting poor tool-use behavior as context grew, similar to the 8B model. With longer training, the model developed more consistent and structured patterns for tool calls. In particular, it learned to invoke get cities early to check available cities within states before searching for tickets and attractions. We also observed that the model made fewer get flights calls across queries, instead preferring to select more grounded transportation options.

LLM USAGE DISCLOSURE

We used large language models (LLMs) *only* as general-purpose writing assistants for surface-level language edits (grammar, phrasing, and clarity), LaTeX formatting suggestions (e.g., table/figure spacing, caption wording), and copyediting (e.g., consistent terminology, acronym expansion). LLMs did *not* contribute to research ideation, dataset or method design, experimental planning, implementation, analysis, or conclusions. All technical content, experiment setup, results, and interpretations were created and verified by the authors, who take full responsibility for the paper's substance and correctness.