# CostBench: Evaluating Multi-Turn Cost-Optimal Planning and Adaptation in Dynamic Environments for LLM Tool-Use Agents

**Jiayu Liu**[1]  **Cheng Qian**[2]  **Zhaochen Su**[1]  **Qing Zong**[1]  **Shijue Huang**[1]  **Bingxiang He**[3]
**Yi R. (May) Fung**[1]

[1]Hong Kong University of Science and Technology
[2]University of Illinois Urbana-Champaign  [3]Tsinghua University
{jliufv,yrfung}@ust.hk    chengq9@illinois.edu

## Abstract

Current evaluations of Large Language Model (LLM) agents primarily emphasize task completion, often overlooking resource efficiency and adaptability. This neglects a crucial capability: agents' ability to devise and adjust cost-optimal plans in response to changing environments. To bridge this gap, we introduce **CostBench**, a scalable, cost-centric benchmark designed to evaluate agents' economic reasoning and replanning abilities. Situated in the travel-planning domain, CostBench comprises tasks solvable via multiple sequences of atomic and composite tools with diverse, customizable costs. It also supports four types of dynamic blocking events, such as tool failures and cost changes, to simulate real-world unpredictability and necessitate agents to adapt in real time. Evaluating leading open-sourced and proprietary models on CostBench reveals a substantial gap in cost-aware planning: agents frequently fail to identify cost-optimal solutions in static settings, with even *GPT-5* achieving less than 75% exact match rate on the hardest tasks, and performance further dropping by around 40% under dynamic conditions. By diagnosing these weaknesses, CostBench lays the groundwork for developing future agents that are both economically rational and robust.[1]

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in reasoning (DeepSeek-AI et al., 2025; Wang et al., 2025f), code generation (Mai et al., 2025; Zhoubian et al., 2025), and complex problem-solving (Yu et al., 2025; Dou et al., 2025). When equipped with external tools, they can interact with dynamic environments such as the web (Qiao et al., 2025; He et al., 2025) and interactive systems (Qin et al., 2024a; Su et al., 2025a). This integration of tool use gives LLMs agentic capabilities, enabling them to autonomously execute multi-step tasks.

While tool-augmented reasoning provides impressive autonomy, its effectiveness ultimately hinges on strategic *planning*: deciding not only which tools to use but also how to sequence them efficiently to achieve complex objectives. Existing evaluations, however, focus primarily on task completion (Zhuang et al., 2023; Wang et al., 2024b; Chen et al., 2025), offering limited insight into how different tool-use strategies affect resource consumption. Although prior studies have examined factors such as API fees (Wang et al., 2024a; Chen et al., 2024a), token usage (Han et al., 2025a; Huang et al., 2025b), GPU computation (Kim et al., 2025), and memory overhead (Wu et al., 2024), they often overlook the agent's ability to reason about costs and remain cost-aware in dynamic environments. As a result, it remains unclear whether current agents can effectively plan and adapt under shifting cost conditions. This gap motivates our central question: **how effectively can LLM agents devise and adapt cost-optimal plans for arbitrary cost functions in dynamic environments?** Answering this requires an evaluation environment that *i)* features **diverse and flexible cost structures** to reveal the agent's cost sensitivity, and *ii)* introduces **runtime dynamics** that evolve throughout the interaction, challenging the agent to sustain cost-awareness and replan adaptively.

To this end, we introduce **CostBench**, a scalable and cost-centric benchmark situated in the travel-planning domain. CostBench defines six travel-related tasks, each of which can be completed through sequences of non-separable atomic tools. The costs of these atomic tools are randomly assigned, enabling diverse cost configurations across different runs. Building on atomic tools, we create composite ones whose costs equal the aggregated component costs plus Gaussian noise. This design creates multiple alternative pathways to accomplish

---

[1]Our code is publicly available at https://github.com/JiayuJeff/CostBench.

| Benchmark | Multi-turn Interaction | Cost-optimal Planning | Flexible Cost | Tool Use | Dynamic State | Adjustable Difficulty | Customizable Details | Scalable |
|---|---|---|---|---|---|---|---|---|
| *PlanBench (Valmeekam et al., 2023)* | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ |
| *ToolBench (Qin et al., 2024b)* | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ |
| *TravelPlanner (Xie et al., 2024)* | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ |
| *AucArena (Chen et al., 2023)* | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ |
| *SayCanPay (Hazra et al., 2024)* | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ |
| *ACPBench Hard (Kokel et al., 2025)* | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ |
| *UserBench (Qian et al., 2025b)* | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ |
| *MCP-Bench (Wang et al., 2025g)* | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ |
| *MINT (Wang et al., 2024b)* | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ |
| **CostBench (Ours)** | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: For each benchmark, the table reports whether each trait is fully (✓), partially (✓), or not (✗) addressed. Detailed explanations are provided in Appendix A.

the same goal, enriching the benchmark's dynamics and encouraging agents to maintain cost-awareness when selecting the optimal execution path. Cost-Bench also introduces dynamic blocking events (e.g., tool cost or user preference changes) that simulate four types of real-world disruptions during multi-turn interactions, requiring agents to replan their trajectories while maintaining cost-optimality.

We evaluate ten leading open-source and proprietary LLMs on CostBench and find that all models perform poorly. *GPT-5* yields the best performance yet attains less than 75% exact match rate (i.e., the tool-call sequence exactly matches the ground-truth trajectory) in the most difficult static settings, and its accuracy further declines to approximately 35% under cost-change conditions. Overall, the models are highly sensitive to cost noise and environmental perturbations, especially when tool costs fluctuate or tools become corrupted after use. These results highlight the fragility of current LLM agents in maintaining cost-awareness and the necessity for more robust, dynamically adaptive models. We summarize our main contributions as follows:

- **Scalable Benchmark Framework:** We introduce CostBench, a scalable and cost-centric framework built on a task-generation pipeline. It programmatically generates tasks with multiple atomic and composite tools, featuring diverse and tunable cost assignments to probe agents' economic reasoning.

- **Dynamic Interaction Environment:** We design an interactive, multi-turn environment that simulates real-world unpredictability. It introduces four types of dynamic blocking events (e.g., cost changes, tool failures) to compel agents to adapt and replan in real time. The environment also serves as a reinforcement learning (RL) playground for agent refinement.

- **Comprehensive Analysis:** We conduct a com-

prehensive evaluation of ten leading LLMs, revealing their high sensitivity to cost variations and dynamic disruptions. Our analysis highlights a substantial gap in cost-aware planning, with performance dropping by around 40% under dynamic conditions, underscoring limitations in achieving robust, cost-optimal planning.

Looking ahead, **CostBench** lays the foundation for future research on adaptive and cost-aware LLM agents, encouraging the development of models that can dynamically plan, learn, and optimize under evolving real-world constraints.

## 2 Related Work

**Evaluation of LLM cost-optimal planning.** Recent work emphasizes evaluating LLMs under economically grounded conditions. Existing tool-use benchmarks such as MINT (Wang et al., 2024b), ToolQA (Zhuang et al., 2023), and MTU-Bench (Wang et al., 2025d) focus on task completion but ignore efficiency. Cost-related benchmarks like PlanBench (Valmeekam et al., 2023), SayCanPay (Hazra et al., 2024), and ACPBench Hard (Kokel et al., 2025, 2024) consider cost only in single-turn planning, limiting applicability to dynamic, iterative settings. Multi-turn benchmarks such as UserBench (Qian et al., 2025b) and TravelPlanner (Xie et al., 2024) focus on candidate prices rather than operational costs, while others like R-ConstraintBench (Jain and Wetter, 2025) and Flex-TravelPlanner (Oh et al., 2025) emphasize performance–resource trade-offs under fixed budgets. Most also conflate cost with path functionality and use static environments, obscuring cost sensitivity. In contrast, CostBench decouples cost from functional bias through randomized tool assignments and equivalent solution paths, while introducing dynamic blocking events to assess adaptive planning. A detailed comparison with prior work is provided in Table 1.
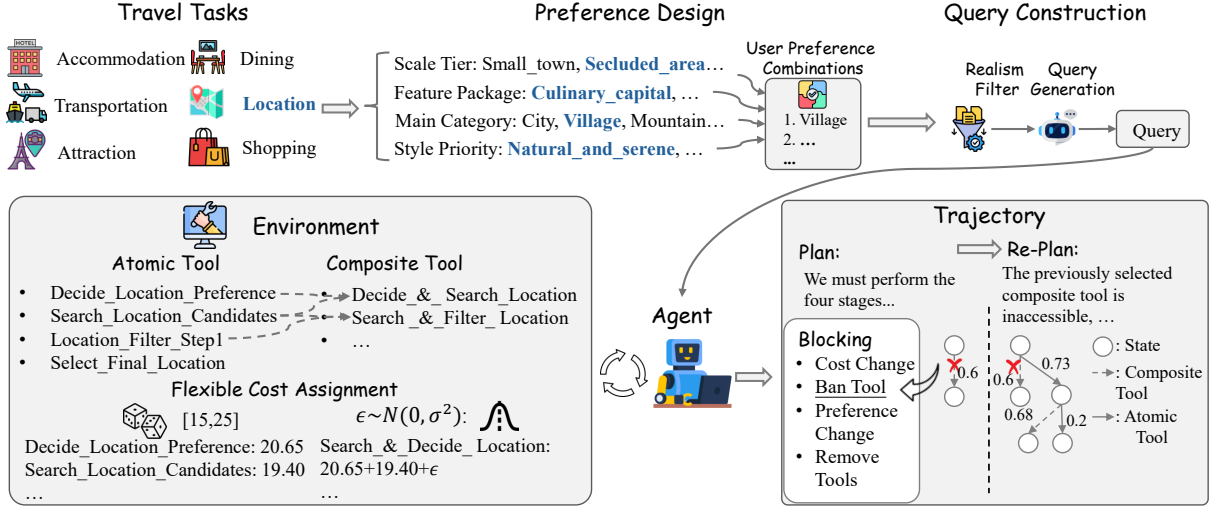
Figure 1: Overview of the **CostBench** pipeline. Starting from high-quality queries generated from combinations of user preferences, the agent constructs its plan, then interacts with an environment set up with atomic and composite tools under flexible cost assignments (atomic tool costs are randomized between 15 and 25 in our experiments), and executes actions along an customizable dynamic blocking module to achieve its goal.

**Cost-centric Agent Design.** Prior work has examined cost-awareness in agents from several angles. Number of tool call represents the simplest form of tool cost which is thoroughly examined and optimized through algorithms (Qian et al., 2025a; Wang et al., 2025c,b). Furthermore, Toolformer (Schick et al., 2023) focuses on when to invoke tools rather than reasoning about fine-grained tool costs or long-term expenditure. Starting from FrugalGPT (Chen et al., 2024b), many studies further explore multi-agent coordination to reduce expenses, achieving cost control through selective and hierarchical model invocation (Zhang et al., 2024; Panda et al., 2025; Gandhi et al., 2025), often in limited or static scenarios. SayCanPay (Hazra et al., 2024) and CATP-LLM (Wu et al., 2024) explicitly optimize for cost but remain limited to one-off planning without dynamic adaptation to changing environments. Other studies treat cost merely as token usage and focus on efficient reasoning under token constraints (Han et al., 2025b; Wen et al., 2025; Huang et al., 2025b), neglecting broader operational costs. Overall, these works target narrow aspects of cost optimization, whereas our benchmark significantly broadens the notion of cost and systematically evaluates agent adaptability under realistic, dynamically changing conditions and diverse practical constraints.

## 3 CostBench

**CostBench** is an interactive, consumption-aware environment designed to evaluate agents' ability for cost-optimal planning in complex reasoning tasks, where the goal is to complete a task while minimizing total cost. It simulates realistic settings in which every action incurs a measurable cost. To enable this, CostBench offers a flexible ecosystem of tools encompassing both atomic operations and higher-level composite functions. Each tool is assigned a randomized cost profile independently resampled for every data instance, preventing memorization and data leakage while encouraging agents to adopt adaptive, cost-sensitive strategies.

Moreover, CostBench can operate in two distinct modes: a **static mode**, where user preferences, tool costs, and tool availability remain constant; and a **dynamic mode**, where these attributes evolve in response to the agent's behavior or environmental changes, capturing the uncertainty inherent in real-world scenarios. Each task unfolds through an iterative interaction loop in which the agent observes the current task state, reasons about its next action, invokes relevant tools, and receives cost-tagged feedback that updates its internal state. This continual feedback enables adaptive decision-making and dynamic replanning until the task goal is achieved. We present the detailed design of this framework in the following.

### 3.1 Environment Setup

In this work, we focus on the travel planning domain. Following UserBench (Qian et al., 2025b) and TravelPlanner (Xie et al., 2024), we categorize travel-related activities into six domains: *Location*,

| Task Name | Pref. Comb. (Unfiltered) | Train Split Filtered | Test Split Filtered |
|---|---|---|---|
| Location | 1552 | 127 (6.68%) | 29 (7.61%) |
| Transportation | 1552 | 296 (15.56%) | 102 (26.77%) |
| Accommodation | 1552 | 417 (21.92%) | 51 (13.39%) |
| Attraction | 1552 | 361 (18.98%) | 56 (14.70%) |
| Dining | 1552 | 282 (14.83%) | 60 (15.75%) |
| Shopping | 1552 | 419 (22.03%) | 83 (21.78%) |
| **Total / Proportion.** | **9312** | **1902 (100%)** | **381 (100%)** |

Table 2: Main statistics of query construction in Cost-Bench. We generate 6,000 training and 1,000 test user preference combinations and filter out those violating common sense. The total number of user preference combinations equal the sum of both splits, with a discussion of the slightly imbalanced distribution provided in Appendix D.1.

| Tools | Value |
|---|---|
| Task Sequence | $N$ |
| Total Atomic Tools | $N$ |
| Total Composite Tools | $N * (N-1)/2$ |
| Total Tools | $N * (N+1)/2$ |

Table 3: Main statistics of tools used in CostBench. The variable $N$ denotes a adjustable task sequence, and the total numbers of atomic and composite tools vary accordingly, as indicated (we use $N = 5$ in most of our experiments).

*Transportation*, *Accommodation*, *Attraction*, *Dining*, and *Shopping*. Each task in CostBench follows a structured four-stage workflow comprising *preference identification*, *candidate search*, *candidate filtering*, and *final selection*. Among these, the filtering stage involves multiple steps, while the other three are modeled as single-step operations. We define the task sequence as the total number of steps across all stages, providing a clear metric for task length and complexity. This multi-step design enhances the scalability of the benchmark, enabling the flexible generation of tasks with varying lengths and levels of difficulty.

**Tool Library.** To ensure interpretability and enable precise progress tracking, we introduce a set of self-defined data types that explicitly specify the input–output schema for each tool. This design prevents agents from skipping intermediate reasoning steps through direct guessing and allows the environment to automatically verify whether each action meaningfully contributes to the final solution. In CostBench, each task step maps to an *atomic tool* performing a non-separable operation, while a *composite tool* is a higher-level tool executing a sequence of atomic tools. For example, "Search_and_Decide_Location" combines "Decide_Location_Preference" and "Search_Location_Preference" into an independent

tool equivalent to running them sequentially (See Figure 1). All tools follow the sequential order of task steps, with each tool consuming the outputs of preceding steps as its inputs (See example in Figure 7). This sequential dependency enables any consecutive subset of atomic tools to form a valid composite tool, thereby supporting scalability and the construction of longer or more complex tool chains while maintaining explicit input–output relationships. Details of tool construction and an example of the tool schema is illustrated in Appendix B.1 and Figure 13.

**Tool Cost Assignment.** Each atomic tool is assigned a randomized cost drawn from a customizable range, while the cost of a composite tool is defined as the sum of its constituent atomic tools plus Gaussian noise with zero mean and a customizable standard deviation. Tool costs are independently resampled for every data instance but remain reproducible, ensuring that environment variations are diverse yet comparable. For tasks sharing the same sequence length, cost assignments are kept consistent across static and dynamic blocking settings to ensure fair comparison. Further implementation details are provided in Appendix B.6.1.

**Agent State Definition.** At each time step, the agent observes the current state $s_t = q, \mathcal{T}, \mathcal{C}, \mathcal{D}$, where $q$ denotes the user query, $\mathcal{T}$ represents the available tool library, $\mathcal{C}$ specifies the corresponding tool cost table, and $\mathcal{D}$ contains the set of self-defined data types obtained from previous tool invocations. Once a data type is acquired, it remains accessible throughout the interaction, reflecting realistic scenarios in which previously gathered information can be reused. The agent's objective is to reach the designated goal state specified clearly in its input while minimizing the total accumulated cost over the entire interaction trajectory.

### 3.2 Agent–Environment Interaction

Each task instance in CostBench follows an iterative interaction loop. At each step, the agent observes the current state, including the user query, available tools and their costs, previously obtained data, and feedback from earlier actions, then decides which atomic or composite tool to invoke next. The environment responds with an updated state and cost-tagged feedback. The interaction terminates once the agent produces the final output or reaches the predefined task sequence length.

**Dynamic Blocking.** CostBench includes an optional dynamic blocking module that introduces runtime perturbations to evaluate the robustness of cost-aware agents. Unlike static interventions, blocking events are triggered adaptively during execution in response to the agent's behaviors. Specifically, the agent may encounter four types of runtime disruptions:

- **Explicit Blocks:** disruptions with direct notification (e.g., tool/user message). (1) *Ban Tool*: behavior-dependent failures where specific tools become unavailable due to the agent's prior actions. (2) *Preference Change*: modifications to user preferences that require the agent to replan its strategy mid-execution.

- **Implicit Blocks:** disruptions without explicit prompt. (1) *Cost Change*: global adjustments to tool costs that alter the current cost landscape. (2) *Remove Tools*: removes composite tools consisting of a certain number of component tools, thereby reducing the available action paths.

The mechanism ensures that: *(1) Dynamic and Uniform Triggering*: Blocking events are triggered dynamically based on the estimated length of the cost-optimal path, leading to approximately uniform trigger positions and ensuring that disruptions occur at meaningful stages. *(2) Independence Across Datapoints*: For each run, environments are instantiated independently across datapoints, preventing correlations among blocking events. *(3) Randomized Yet Reproducible Execution*: Even under dynamic conditions, the experiments remain *randomized yet reproducible*, preserving both diversity and fairness in evaluation (see Appendix B.6).

### 3.3 User Query Construction

For each of the six travel-planning tasks, we construct user queries through a multi-stage pipeline. Each task defines four preference dimensions with ten candidate values, where four are used for test query construction and six for the training set, ensuring clear differentiation and non-overlapping splits, resulting in $4^4 + 6^4 = 1,552$ unique user preference combinations per task. All combinations are verified by *GPT-4o* and partially checked manually for clarity (see Appendix E). Natural language queries are generated by *GPT-4o* to capture realistic vagueness while preserving clear distinctions among user preferences within each dimension. An example query is shown in Figure 12, and dataset statistics are summarized in Tables 2 and 3.

## 4 Experiment

### 4.1 Settings

We evaluate both proprietary and open-source models to ensure a balanced and comprehensive assessment of current LLM capabilities. Proprietary models include *GPT* (OpenAI, 2024b), *Claude* (Anthropic, 2025), *Deepseek* (DeepSeek-AI et al., 2025), and *Gemini* (Google DeepMind, 2025); while open-source models include *GLM* (Team et al., 2025), *Qwen3* (Qwen Team, 2025), and *Llama3* (Touvron et al., 2023). All models use a temperature of 0.0 for deterministic decoding and a maximum token length of 16,384 to prevent output truncation.

**Metrics.** We evaluate performance using metrics across four main categories: cost, path, user-intent understanding, and tool use. The majority of these metrics quantify the deviation of an agent's trajectory from a ground-truth trajectory, with a greedy policy serving as a baseline for comparison. Detailed definitions of all metrics are provided in Appendix B.3.

**(1) Cost Gap:** This metric directly measures the gap between gt accumulative costs and agent's accumulative cost.

**(2) Average Edit Distance (AED):** This metric measures the edit distance between the ground-truth and the agent's trajectories to assess the quality of the agent's path.

**(3) Average Normalized Edit Distance (ANED) (%):** This metric normalizes the ED to [0, 1] for fair comparison.

**(4) Exact Match Ratio (EMR) (%):** This metric evaluates the ratio of datapoints that agent implements an identical tool call trajectory as gt.

**(5) User Intent Hit Ratio (UIHR) (%):** This metric reflects whether the agent correctly understands the user intent in the query. Note that there is only one correct answer for each datapoint.

**(6) Invalid Tool-Use Ratio (ITUR) (%)** This metric evaluates the proportion of valid tool use. Invalid behaviour includes calling wrong tool name or calling wrong parameter name, call a tool but doesn't actually have the input type, etc.

**Prompts and Environment Settings** Detailed prompt and environment hyperparameters are in Figure 14 and Table 5. To prevent trivial one-step solutions, tools that complete the entire procedure in a single call are excluded.

| Model Name | Cost | Path | | | User Intent | Tool Call |
| --- | --- | --- | --- | --- | --- | --- |
| | Cost Gap ↓ | AED ↓ | ANED (%) ↓ | EMR (%) ↑ | UIHR (%) ↑ | ITUR (%) ↓ |
| *Greedy Policy* | 0.269 (0.269) | 2.202 | 74.74 | 10.76 | - | - |
| *Qwen3-8B* | 0.477 (0.239) | 2.024 | 71.05 | 17.32 | 82.94 | 0.72 |
| *Qwen3-14B* | 0.670 (0.130) | 1.976 | 51.26 | 33.45 | 88.71 | 0.91 |
| *Qwen3-32B* | <u>0.108</u> (0.108) | 1.748 | 52.51 | 33.60 | 90.55 | 1.73 |
| *Llama-3.1-8B-Instruct* | 0.303 (0.303) | 2.215 | 80.21 | 9.71 | 63.25 | 32.50 |
| *GLM-4.5* | 0.165 (0.066) | 1.076 | 34.79 | 54.33 | 90.55 | 4.22 |
| *Deepseek-V3.1* | 2.221 (0.200) | 2.320 | 65.28 | 17.85 | 87.93 | 21.26 |
| *Gemini-2.5-Pro* | 0.312 (<u>0.015</u>) | <u>0.367</u> | <u>10.98</u> | <u>83.73</u> | **94.49** | **0.00** |
| *Claude-Sonnet-4.0* | 0.249 (0.078) | 0.982 | 31.28 | 59.84 | <u>93.58</u> | 5.04 |
| *GPT-4o* | 1.179 (0.228) | 2.454 | 68.12 | 13.65 | 90.29 | 19.96 |
| *GPT-5* | **0.060 (0.002)** | **0.103** | **3.31** | **95.52** | 92.11 | <u>0.01</u> |

Table 4: **CostBench** main evaluation results for task sequences of length 5 in the static setting. Scores in **bold** indicate the best performance among all models, and <u>underlined</u> scores denote the second-best performance. The metrics in (%) are in percentage forms. Reported values under **Cost Gap** indicate the average cost gap, with numbers in parentheses computed after excluding samples with redundant tool calls (see Appendix D.2 for details). The greedy baseline selects the atomic tool with the lowest average cost at each step (See Appendix B.2).
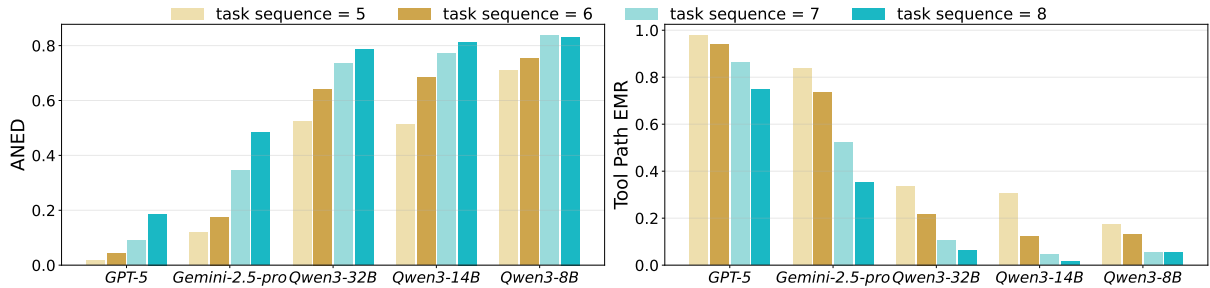


Figure 2: Models' average normalized edit distance and exact match ratio for task sequences of length five to eight. Models show significantly lower performance as task complexity increases. In the most challenging case with task sequence length of eight, even the strongest model, *GPT-5*, achieves less than 75% exact match ratio.

## 4.2 Results

**Most models perform poorly in cost-optimal planning.** Our main results in Table 4 reveal several notable trends. First, only *GPT-5* and *Gemini-2.5-Pro* show strong exact-match performance (95.52% and 83.73%, respectively) at a task length of five, whereas all other models fall below 60%. Performance further declines for longer task sequences (see Section 5.1), where even *GPT-5* drops to below **75%** EMR. This highlights a substantial gap in executing complete cost-aware plans, as even the strongest model, *GPT-5*, falls short of reliably producing optimal trajectories. Second, compared to the greedy baseline, some LLMs fail to show convincing improvements. For example, *Qwen3-8B*, *GPT-4o*, and *Deepseek-V3.1* yield nearly identical NED scores to greedy policy, with *GPT-4o* even exhibiting a higher ED, indicating that its generated plans deviate further from optimal solution than greedy baseline. More concerningly, *Llama-3.1-8B-Instruct* underperforms the greedy baseline across all path-related metrics, suggesting fundamental weaknesses in cost-optimal planning.

**Comparison of Open-Source and Proprietary Models.** Beyond the greedy baseline, we compare open-source and proprietary models. Interestingly, some proprietary models show no clear advantage (Chen et al., 2025): *Qwen3-32B* consistently surpasses *GPT-4o* and *DeepSeek-V3.1* on all path-related metrics, and *GLM-4.5* achieves performance comparable to *Claude-Sonnet-4.0*. While partly due to differing invalid tool-use rates, our metrics decouple tool use failure from cost/path measurements (see Appendix B.3.2), indicating that some proprietary models' shortcomings in cost-sensitive planning are substantive and reflect genuine limitations in their reasoning capabilities.

**Model suffers from redundant tool calls and tool call failures.** In Table 4, we report two variants of the cost gap metric: values outside the parentheses are computed over all samples, while those inside exclude two abnormal cases—*repeated tool calls* and *extra tool calls*. After removing these anomalies, the cost metrics align with the path-related metrics in reflecting model performance. Beyond such redundant calls, we also observe fail-
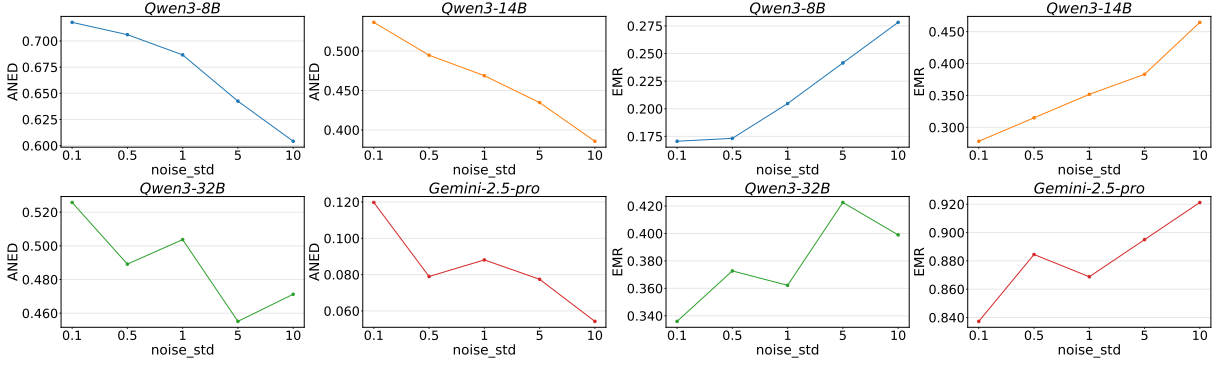
Figure 3: LLMs' performance on CostBench under different standard deviations of composite tool cost noise. Models tend to perform better with higher noise levels, indicating high sensitivity to tool cost variations.
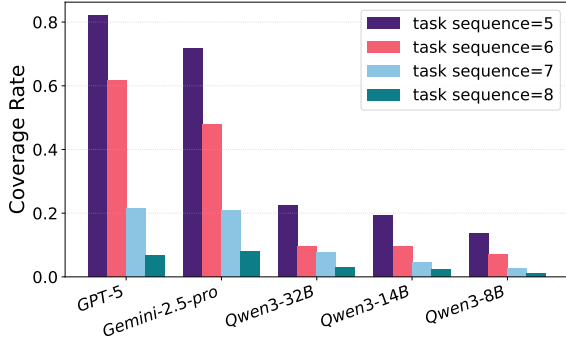


Figure 4: Coverage rates of different LLMs across task sequences of length 5 to 8. Model performance in the static setting shows a strong positive correlation with coverage rate.

ure cases including *invalid parameter inputs* and attempts to invoke *inaccessible tools*, the latter being the most severe. These errors reveal models' limited progress awareness, which in turn hinders their ability to plan in a cost-optimal manner. Further discussion of the failure modes and additional findings in Table 4 can be found in Appendix D.4.

## 5 Analysis

### 5.1 Static Setting

**Model performance declines as task complexity increases.** As shown in Figure 2, performance consistently deteriorates as task sequences grow longer. The best-performing model, *GPT-5*, drops from near-perfect exact match ratio at length 5 to around 75% at length 8, accompanied by a higher normalized edit distance, while weaker models such as *Qwen3-8B* and *Qwen3-14B* collapse entirely (EMR = 0). This overall downward trend highlights the increasing difficulty current models face in sustaining coherent, cost-optimal reasoning as task complexity rises.

**Coverage Rate in planning significantly influences performance.** In the unblocking setting,

the model's task is to enumerate all possible paths and select the optimal one. We define coverage as the proportion of paths explicitly planned, and hypothesize that higher coverage leads to better performance. To identify distinct paths in each model's output, we use *GPT-4o-mini* (OpenAI, 2024a) (prompt in Figure 11). As shown in Table 4, performance follows *GPT-5 > Gemini-2.5-pro >* open-source *Qwen* models, which aligns with the coverage trend in Figure 4. Performance also declines as task sequences grow longer (Figures 2), explained by the corresponding drop in coverage. Importantly, even when given abstract examples that enumerate all paths (Figure 15), models often fail to generalize these planning strategies to real scenarios.

**Models exhibit pronounced cost sensitivity.** To examine models' sensitivity to cost variations, we vary the noise standard deviation when constructing composite tools, which controls how much a composite tool's cost deviates from the sum of its component tools. Larger noise introduces greater diversity in cost structures and reshapes the cost landscape. As shown in Figure 3, model performance, measured by both ANED and EMR, consistently improves as the noise increases. This result suggests that models are highly responsive to cost signals, as greater cost differences between action paths make it easier for them to identify the optimal path without explicitly computing the total cost of each option.

### 5.2 Dynamic Setting

**Models remain highly vulnerable to dynamic blocking.** From Figure 5, we observe that dynamic blocking substantially affects both ANED and EMR across all models. Among the blocking types, *cost change* and *ban tool* consistently induce the most significant degradation. For *cost change*,
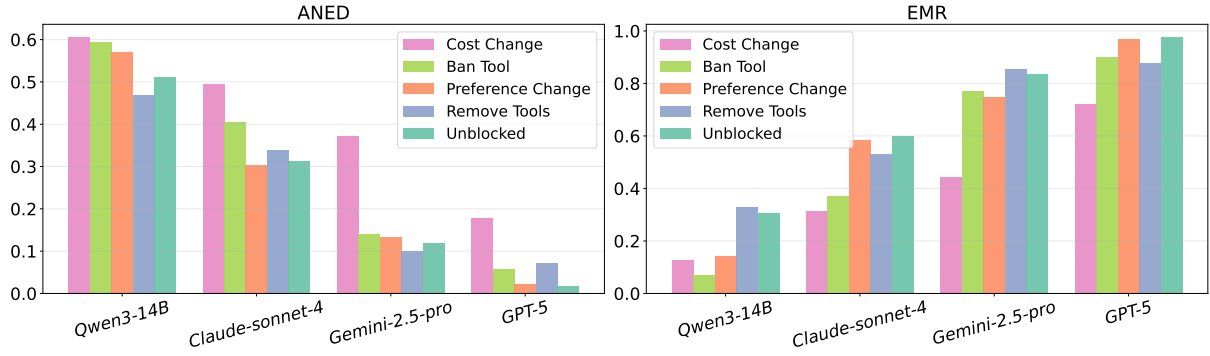
Figure 5: LLMs' performance in CostBench's dynamic blocking setting. All models show consistent EMR drops under ban tool, cost change, and preference change conditions, indicating weak replanning and adaptation abilities in dynamic environments.
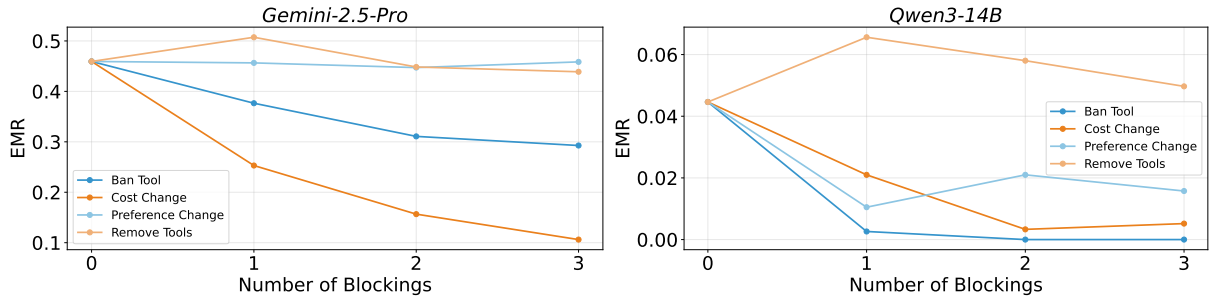


Figure 6: Performance of *Gemini-2.5-Pro* and *Qwen3-14B* under increasing numbers of blocking events (task sequence length = 7). Each curve represents a different blocking type. Both models degrade with more blockings, especially under frequent cost changes or tool bans, with *Qwen3-14B* showing near-total failure.

the EMR of *GPT-5*, which originally achieved nearly perfect EMR, drops by around 20 percentage points, while *Gemini-2.5-Pro* suffers an even larger decline of almost 40 points. To further examine model robustness in the most challenging cases, we evaluate *GPT-5* under the hardest *cost-change* setting, where its Exact Match falls to merely **34.91%**. Under the *ban tool* condition, both *GPT-5* and *Gemini-2.5-Pro* still experience noticeable drops (approximately 5–10%), yet they exhibit stronger resilience compared to the other two models. In addition, *preference change* also causes a stable decrease in Exact Match across all evaluated models. In contrast, the effect of *remove tools* is less consistent, which we attribute to the reduced search space allowing models to achieve higher coverage and thus more stable performance (see Section 5.1 for the coverage–performance relationship).

To better understand performance degradation across blocking types, we estimate the relative adaptation difficulty by measuring how much ground-truth trajectories change under each condition, using the average normalized edit distance (ANED) between unblocked and blocked settings. Higher ANED indicates greater deviation in the optimal path and thus higher adaptation difficulty. Results show that *cost change* (mean$_{\text{GT-NED}}$ =

0.378) and *ban tool* (0.538) cause larger deviations than *preference change* (0.295) and *remove tools* (0.214), aligning with the stronger performance drops in Figure 5. Notably, although *ban tool* produces greater trajectory deviation, *cost change* yields the steepest performance decline, suggesting that implicit blocking (e.g., cost perturbation) poses extra challenges, as models must first detect the change before replanning—revealing their limited adaptability.

**Multiple blockings severely undermines model robustness.** We further analyze the results under multiple blocking events (Figure 6), focusing on how repeated disruptions impact model robustness and cost-optimal planning. We find that increasing the frequency of both *ban tool* and *cost change* blocks leads to a substantial drop in EMR across all models. For instance, the EMR of *Gemini-2.5-Pro*, which initially achieved around 45% accuracy, drops to nearly 0.1 after three consecutive *cost-change* events, while *Qwen3-14B* completely fails (EMR = 0) after only two such events. This indicates that even state-of-the-art commercial and open-source models struggle to maintain reasoning consistency and cost-awareness under dynamically perturbed environments, highlighting their limited robustness and adaptability to repeated disruptions.

## 6    Conclusion

We present CostBench, a scalable, cost-focused benchmark for evaluating LLM agents' ability to plan and adapt cost-optimally in dynamic environments. Centered on travel planning, it includes atomic and composite tools with randomized costs and four types of dynamic blocking events that force adaptive replanning. Evaluating ten leading LLMs, we find substantial weaknesses: even the best-performing *GPT-5* achieves under 75% exact match on the hardest static tasks, with performance dropping to around 35% under cost-change conditions and showing similar vulnerability when tools are banned. These results reveal limitations in current models' cost-aware reasoning, path enumeration, and adaptability. Beyond diagnosing these weaknesses, CostBench provides a principled framework to drive the next generation of LLM agents toward economically rational, resource-efficient, and resilient decision-making in complex, real-world scenarios.

## Limitations

While CostBench provides a principled and scalable framework for evaluating cost-aware planning, several limitations remain. First, our current tasks are restricted to the travel-planning domain, which, though diverse, may not fully capture the breadth of real-world cost trade-offs. Nevertheless, Cost-Bench can be easily extended to other domains by adding configuration files. Second, our simulation abstracts away true API latency, resource consumption, and stochastic failures, which may further influence agent behavior in deployment. Third, the blocking events we design, while representative, are still manually parameterized; future work could explore learning-based or user-driven dynamics to better approximate natural environmental shifts. Finally, incorporating multimodal agents capable of processing and reasoning over diverse data types such as text and images represents another critical research direction (Su et al., 2025b).

## Ethics Statement

**Offensive Content Elimination.**    Our benchmark focuses on the travel-planning domain, and its curation does not rely on content generated by LLMs. Instead, all data were carefully sampled and validated to ensure the dataset is free of offensive material. Consequently, we are confident that it poses no risk of negative societal impact.

**Licenses.**    Our code will be released under the MIT license to allow unrestricted research use. The CostBench will be distributed under a Creative Commons (CC) license, providing free access for the academic community. We take full responsibility for any potential rights violations or licensing issues, and all resources comply with their respective terms of use while supporting research purposes.

**Models.**    All open-source models were hosted and executed locally using the vLLM library (Kwon et al., 2023), while all closed-source models were accessed through their respective official APIs. For reproducibility, the experimental settings are detailed in Section 4.1.

**Data Annotations.**    All data annotation was performed by the paper's co-authors, who are qualified researchers with relevant expertise, ensuring that the process was conducted responsibly and in accordance with ethical standards.

## References

Anthropic. 2025. Introducing claude 4. `https://www.anthropic.com/news/claude-4`.

Arturs Backurs and Piotr Indyk. 2015. Edit distance cannot be computed in strongly subquadratic time (unless SETH is false). In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 51–58. ACM.

Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Yuefeng Huang, Wulong Liu, Xinzhi Wang, Defu Lian, Baoqun Yin, Yasheng Wang, and Wu Liu. 2025. Acebench: Who wins the match point in tool usage? *Preprint*, arXiv:2501.12851.

Jiangjie Chen, Siyu Yuan, Rong Ye, Bodhisattwa Prasad Majumder, and Kyle Richardson. 2023. Put your money where your mouth is: Evaluating strategic planning and execution of LLM agents in an auction arena. *CoRR*, abs/2310.05746.

Lingjiao Chen, Matei Zaharia, and James Zou. 2024a. FrugalGPT: How to use large language models while reducing cost and improving performance. *Transactions on Machine Learning Research*.

Lingjiao Chen, Matei Zaharia, and James Zou. 2024b. FrugalGPT: How to use large language models while reducing cost and improving performance.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong

Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *Preprint*, arXiv:2501.12948.

Edsger W Dijkstra. 1959. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1):269–271.

Zhihao Dou, Qinjian Zhao, Zhongwei Wan, Dinggen Zhang, Weida Wang, Towsif Raiyan, Benteng Chen, Qingtao Pan, Yang Ouyang, Zhiqiang Gao, Shufei Zhang, and Sumon Biswas. 2025. Plan then action:high-level planning guidance reinforcement learning for llm reasoning. *Preprint*, arXiv:2510.01833.

Shubham Gandhi, Manasi Patwardhan, Lovekesh Vig, and Gautam Shroff. 2025. Budgetmlagent: A cost-effective llm multi-agent system for automating machine learning tasks. *Preprint*, arXiv:2411.07464.

Xinyu Geng, Peng Xia, Zhen Zhang, Xinyu Wang, Qiuchen Wang, Ruixue Ding, Chenxi Wang, Jialong Wu, Yida Zhao, Kuan Li, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025. Webwatcher: Breaking new frontier of vision-language deep research agent. *Preprint*, arXiv:2508.05748.

Google DeepMind. 2025. Gemini 2.5 pro overview. https://deepmind.google/technologies/gemini/pro/.

Dadi Guo, Jiayu Liu, Zhiyuan Fan, Zhitao He, Haoran Li, Yumeng Wang, and Yi R. Fung. 2025. Mathematical proof as a litmus test: Revealing failure modes of advanced large reasoning models. *Preprint*, arXiv:2506.17114.

Rishin Haldar and Debajyoti Mukhopadhyay. 2011. Levenshtein distance technique in dictionary lookup methods: An improved approach. *Preprint*, arXiv:1101.1232.

Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. 2025a. Token-budget-aware LLM reasoning. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 24842–24855, Vienna, Austria. Association for Computational Linguistics.

Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. 2025b. Token-budget-aware LLM reasoning. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 24842–24855, Vienna, Austria. Association for Computational Linguistics.

Rishi Hazra, Pedro Zuidberg Dos Martires, and Luc De Raedt. 2024. Saycanpay: Heuristic planning with large language models using learnable domain knowledge. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 20123–20133. AAAI Press.

Zhitao He, Zijun Liu, Peng Li, Yi R. Fung, Ming Yan, Ji Zhang, Fei Huang, and Yang Liu. 2025. Advancing language multi-agent learning with credit re-assignment for interactive environment generalization. In *Second Conference on Language Modeling*.

Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun Chen, Chiyuan Zhang, and Mengdi Wang. 2025a. Mathperturb: Benchmarking llms' math reasoning abilities against hard perturbations. *CoRR*, abs/2502.06453.

Shijue Huang, Hongru Wang, Wanjun Zhong, Zhaochen Su, Jiazhan Feng, Bowen Cao, and Yi R Fung. 2025b.

Adactrl: Towards adaptive and controllable reasoning via difficulty-aware budgeting. *arXiv preprint arXiv:2505.18822*.

Raj Jain and Marc Wetter. 2025. R-constraintbench: Evaluating llms on np-complete scheduling. *Preprint*, arXiv:2508.15204.

Kyoungmin Kim, Kijae Hong, Caglar Gulcehre, and Anastasia Ailamaki. 2025. Optimizing llm inference for database systems: Cost-aware scheduling for concurrent requests. *Preprint*, arXiv:2411.07447.

Harsha Kokel, Michael Katz, Kavitha Srinivas, and Shirin Sohrabi. 2024. Acpbench: Reasoning about action, change, and planning. *Preprint*, arXiv:2410.05669.

Harsha Kokel, Michael Katz, Kavitha Srinivas, and Shirin Sohrabi. 2025. Acpbench hard: Unrestrained reasoning about action, change, and planning. *CoRR*, abs/2503.24378.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. *Preprint*, arXiv:2309.06180.

Pengxiang Li, Zechen Hu, Zirui Shang, Jingrong Wu, Yang Liu, Hui Liu, Zhi Gao, Chenrui Shi, Bofei Zhang, Zihao Zhang, Xiaochuan Shi, Zedong YU, Yuwei Wu, Xinxiao Wu, Yunde Jia, Liuyu Xiang, Zhaofeng He, and Qing Li. 2025. Efficient multi-turn RL for GUI agents via decoupled training and adaptive data curation. *CoRR*, abs/2509.23866.

Yujian Li and Bi Liu. 2007. A normalized levenshtein distance metric. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(6):1091–1095.

Jiayu Liu, Qing Zong, Weiqi Wang, and Yangqiu Song. 2025. Revisiting epistemic markers in confidence estimation: Can markers accurately reflect large language models' uncertainty? In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 206–221, Vienna, Austria. Association for Computational Linguistics.

Jiarui Lu, Thomas Holleis, Yizhe Zhang, Bernhard Aumayer, Feng Nan, Haoping Bai, Shuang Ma, Shen Ma, Mengyu Li, Guoli Yin, Zirui Wang, and Ruoming Pang. 2025. ToolSandbox: A stateful, conversational, interactive evaluation benchmark for LLM tool use capabilities. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 1160–1183, Albuquerque, New Mexico. Association for Computational Linguistics.

Xinji Mai, Haotian Xu, Zhong-Zhi Li, Xing W, Weinong Wang, Jian Hu, Yingying Zhang, and Wenqiang Zhang. 2025. Agent rl scaling law: Agent rl with spontaneous code execution for mathematical problem solving. *Preprint*, arXiv:2505.07773.

National Institute of Standards and Technology. 2015. Secure hash standard (shs). Technical Report FIPS PUB 180-4, U.S. Department of Commerce.

Juhyun Oh, Eunsu Kim, and Alice Oh. 2025. Flextravelplanner: A benchmark for flexible planning with language agents. *Preprint*, arXiv:2506.04649.

OpenAI. 2024a. Gpt-4o mini: advancing cost-efficient intelligence. *OpenAI*.

OpenAI. 2024b. Hello gpt-4o. *OpenAI*.

Pranoy Panda, Raghav Magazine, Chaitanya Devaguptapu, Sho Takemori, and Vishal Sharma. 2025. Adaptive llm routing under budget constraints. *Preprint*, arXiv:2508.21141.

Cheng Qian, Emre Can Acikgoz, Hongru Wang, Xiusi Chen, Avirup Sil, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. 2025a. Smart: Self-aware agent for tool overuse mitigation. *arXiv preprint arXiv:2502.11435*.

Cheng Qian, Zuxin Liu, Akshara Prabhakar, Zhiwei Liu, Jianguo Zhang, Haolin Chen, Heng Ji, Weiran Yao, Shelby Heinecke, Silvio Savarese, Caiming Xiong, and Huan Wang. 2025b. Userbench: An interactive gym environment for user-centric agents. *Preprint*, arXiv:2507.22034.

Zile Qiao, Guoxin Chen, Xuanzhong Chen, Donglei Yu, Wenbiao Yin, Xinyu Wang, Zhen Zhang, Baixuan Li, Huifeng Yin, Kuan Li, Rui Min, Minpeng Liao, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025. Webresearcher: Unleashing unbounded reasoning capability in long-horizon agents. *Preprint*, arXiv:2509.13309.

Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2024a. Tool learning with foundation models. *Preprint*, arXiv:2304.08354.

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. 2024b. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. In *The Twelfth International Conference on Learning Representations*.

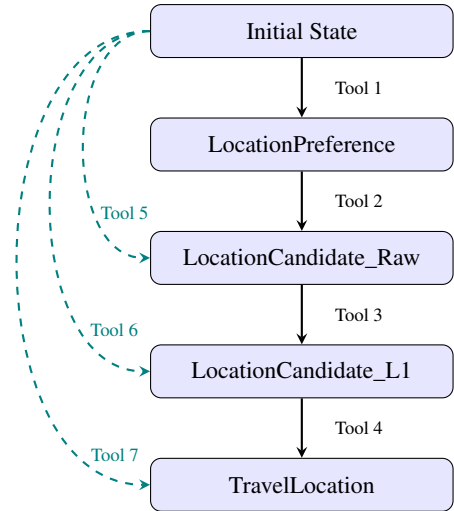Qwen Team. 2025. Qwen3: Think Deeper, Act Faster. https://qwenlm.github.io/blog/qwen3/. Blog post.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.*

Zhaochen Su, Linjie Li, Mingyang Song, Yunzhuo Hao, Zhengyuan Yang, Jun Zhang, Guanjie Chen, Jiawei Gu, Juntao Li, Xiaoye Qu, et al. 2025a. Openthinkimg: Learning to think with images via visual tool reinforcement learning. *arXiv preprint arXiv:2505.08617.*

Zhaochen Su, Peng Xia, Hangyu Guo, Zhenhua Liu, Yan Ma, Xiaoye Qu, Jiaqi Liu, Yanshu Li, Kaide Zeng, Zhengyuan Yang, et al. 2025b. Thinking with images for multimodal reasoning: Foundations, methods, and future frontiers. *arXiv preprint arXiv:2506.23918.*

5 Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang, Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei, Yean Cheng, Yifan An, Yilin Niu, Yuanhao Wen, Yushi Bai, Zhengxiao Du, Zihan Wang, Zilin Zhu, Bohan Zhang, Bosi Wen, Bowen Wu, Bowen Xu, Can Huang, Casey Zhao, Changpeng Cai, Chao Yu, Chen Li, Chendi Ge, Chenghua Huang, Chenhui Zhang, Chenxi Xu, Chenzheng Zhu, Chuang Li, Congfeng Yin, Daoyan Lin, Dayong Yang, Dazhi Jiang, Ding Ai, Erle Zhu, Fei Wang, Gengzheng Pan, Guo Wang, Hailong Sun, Haitao Li, Haiyang Li, Haiyi Hu, Hanyu Zhang, Hao Peng, Hao Tai, Haoke Zhang, Haoran Wang, Haoyu Yang, He Liu, He Zhao, Hongwei Liu, Hongxi Yan, Huan Liu, Huilong Chen, Ji Li, Jiajing Zhao, Jiamin Ren, Jian Jiao, Jiani Zhao, Jianyang Yan, Jiaqi Wang, Jiayi Gui, Jiayue Zhao, Jie Liu, Jijie Li, Jing Li, Jing Lu, Jingsen Wang, Jingwei Yuan, Jingxuan Li, Jingzhao Du, Jinhua Du, Jinxin Liu, Junkai Zhi, Junli Gao, Ke Wang, Lekang Yang, Liang Xu, Lin Fan, Lindong Wu, Lintao Ding, Lu Wang, Man Zhang, Minghao Li, Minghuan Xu, Mingming Zhao, Mingshu Zhai, Pengfan Du, Qian Dong, Shangde Lei, Shangqing Tu, Shangtong Yang, Shaoyou Lu, Shijie Li, Shuang Li, Shuang-Li, Shuxun Yang, Sibo Yi, Tianshu Yu, Wei Tian, Weihan Wang, Wenbo Yu, Weng Lam Tam, Wenjie Liang, Wentao Liu, Xiao Wang, Xiaohan Jia, Xiaotao Gu, Xiaoying Ling, Xin Wang, Xing Fan, Xingru Pan, Xinyuan Zhang, Xinze Zhang, Xiuqing Fu, Xunkai Zhang, Yabo Xu, Yandong Wu, Yida Lu, Yidong Wang, Yilin Zhou, Yiming Pan, Ying Zhang, Yingli Wang, Yingru Li, Yinpei Su, Yipeng Geng, Yitong Zhu, Yongkun Yang, Yuhang Li, Yuhao Wu, Yujiang Li, Yunan Liu, Yunqing Wang, Yuntao Li, Yuxuan Zhang, Zezhen Liu, Zhen Yang, Zhengda Zhou, Zhongpei Qiao, Zhuoer Feng, Zhuorui Liu, Zichen Zhang, Zihan Wang, Zijun Yao, Zikang Wang, Ziqiang Liu, Ziwei Chai, Zixuan Li,

Zuodong Zhao, Wenguang Chen, Jidong Zhai, Bin Xu, Minlie Huang, Hongning Wang, Juanzi Li, Yuxiao Dong, and Jie Tang. 2025. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models. *Preprint*, arXiv:2508.06471.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.

Karthik Valmeekam, Matthew Marquez, Alberto Olmo Hernandez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.*

Haoming Wang, Haoyang Zou, Huatong Song, Jiazhan Feng, Junjie Fang, Junting Lu, Longxiang Liu, Qinyu Luo, Shihao Liang, Shijue Huang, Wanjun Zhong, Yining Ye, Yujia Qin, Yuwen Xiong, Yuxin Song, Zhiyong Wu, Aoyan Li, Bo Li, Chen Dun, Chong Liu, Daoguang Zan, Fuxing Leng, Hanbin Wang, Hao Yu, Haobin Chen, Hongyi Guo, Jing Su, Jingjia Huang, Kai Shen, Kaiyu Shi, Lin Yan, Peiyao Zhao, Pengfei Liu, Qinghao Ye, Renjie Zheng, Shulin Xin, Wayne Xin Zhao, Wen Heng, Wenhao Huang, Wenqian Wang, Xiaobo Qin, Yi Lin, Youbin Wu, Zehui Chen, Zihao Wang, Baoquan Zhong, Xinchun Zhang, Xujing Li, Yuanfan Li, Zhongkai Zhao, Chengquan Jiang, Faming Wu, Haotian Zhou, Jinlin Pang, Li Han, Qi Liu, Qianli Ma, Siyao Liu, Songhua Cai, Wenqi Fu, Xin Liu, Yaohui Wang, Zhi Zhang, Bo Zhou, Guoliang Li, Jiajun Shi, Jiale Yang, Jie Tang, Li Li, Qihua Han, Taoran Lu, Woyu Lin, Xiaokang Tong, Xinyao Li, Yichi Zhang, Yu Miao, Zhengxuan Jiang, Zili Li, Ziyuan Zhao, Chenxin Li, Dehua Ma, Feng Lin, Ge Zhang, Haihua Yang, Hangyu Guo, Hongda Zhu, Jiaheng Liu, Junda Du, Kai Cai, Kuanye Li, Lichen Yuan, Meilan Han, Minchao Wang, Shuyue Guo, Tianhao Cheng, Xiaobo Ma, Xiaojun Xiao, Xiaolong Huang, Xinjie Chen, Yidi Du, Yilin Chen, Yiwen Wang, Zhaojian Li, Zhenzhu Yang, Zhiyuan Zeng, Chaolin Jin, Chen Li, Hao Chen, Haoli Chen, Jian Chen, Qinghao Zhao, and Guang Shi. 2025a. Ui-tars-2 technical report: Advancing gui agent with multi-turn reinforcement learning. *Preprint*, arXiv:2509.02544.

Hongru Wang, Cheng Qian, Manling Li, Jiahao Qiu, Boyang Xue, Mengdi Wang, Heng Ji, and Kam-Fai Wong. 2025b. Toward a theory of agents as tool-use decision-makers. *arXiv preprint arXiv:2506.00886.*

Hongru Wang, Cheng Qian, Wanjun Zhong, Xiusi Chen, Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang, Kam-Fai Wong, and Heng Ji. 2025c. Acting less is

reasoning more! teaching model to act efficiently. *arXiv preprint arXiv:2504.14870.*

Junlin Wang, Siddhartha Jain, Dejiao Zhang, Baishakhi Ray, Varun Kumar, and Ben Athiwaratkun. 2024a. Reasoning in token economies: Budget-aware evaluation of LLM reasoning strategies. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19916–19939, Miami, Florida, USA. Association for Computational Linguistics.

Pei Wang, Yanan Wu, Noah Wang, Jiaheng Liu, Xiaoshuai Song, Z. Y. Peng, Ken Deng, Chenchen Zhang, Jiakai Wang, Junran Peng, Ge Zhang, Hangyu Guo, Zhaoxiang Zhang, Wenbo Su, and Bo Zheng. 2025d. Mtu-bench: A multi-granularity tool-use benchmark for large language models. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net.

Rui Wang, Qihan Lin, Jiayu Liu, Qing Zong, Tianshi Zheng, Weiqi Wang, and Yangqiu Song. 2025e. Prospect theory fails for llms: Revealing instability of decision-making under epistemic uncertainty. *Preprint*, arXiv:2508.08992.

Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2024b. MINT: evaluating llms in multi-turn interaction with tools and language feedback. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.

Yumeng Wang, Zhiyuan Fan, Jiayu Liu, Jen tse Huang, and Yi R. Fung. 2025f. Diversity-enhanced reasoning for subjective questions. *Preprint*, arXiv:2507.20187.

Zhenting Wang, Qi Chang, Hemani Patel, Shashank Biju, Cheng-En Wu, Quan Liu, Aolin Ding, Alireza Rezazadeh, Ankit Shah, Yujia Bao, and Eugene Siow. 2025g. Mcp-bench: Benchmarking tool-using llm agents with complex real-world tasks via mcp servers. *Preprint*, arXiv:2508.20453.

Hao Wen, Xinrui Wu, Yi Sun, Feifei Zhang, Liye Chen, Jie Wang, Yunxin Liu, Yunhao Liu, Ya-Qin Zhang, and Yuanchun Li. 2025. Budgetthinker: Empowering budget-aware llm reasoning with control tokens. *Preprint*, arXiv:2508.17196.

Duo Wu, Jinghe Wang, Yuan Meng, Yanning Zhang, Le Sun, and Zhi Wang. 2024. CATP-LLM: empowering large language models for cost-aware tool planning. *CoRR*, abs/2411.16313.

Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025a. Webdancer: Towards autonomous information seeking agency. *Preprint*, arXiv:2505.22648.

Mingqi Wu, Zhihao Zhang, Qiaole Dong, Zhiheng Xi, Jun Zhao, Senjie Jin, Xiaoran Fan, Yuhao Zhou, Huijie Lv, Ming Zhang, Yanwei Fu, Qin Liu, Songyang Zhang, and Qi Zhang. 2025b. Reasoning or memorization? unreliable results of reinforcement learning due to data contamination. *Preprint*, arXiv:2507.10532.

Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024. Travelplanner: A benchmark for real-world planning with language agents. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024.* OpenReview.net.

Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R. Narasimhan. 2025. $\{\tau\}$-bench: A benchmark for \underline{T}ool-\underline{A}gent-\underline{U}ser interaction in real-world domains. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net.

Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. 2025. Dapo: An open-source llm reinforcement learning system at scale. *Preprint*, arXiv:2503.14476.

Jieyu Zhang, Ranjay Krishna, Ahmed Hassan Awadallah, and Chi Wang. 2024. Ecoassistant: Using LLM assistant more affordably and accurately.

Sining Zhoubian, Dan Zhang, and Jie Tang. 2025. Rest-rl: Achieving accurate code reasoning of llms with optimized self-training and decoding. *Preprint*, arXiv:2508.19576.

Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. Toolqa: A dataset for LLM question answering with external tools. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.*

Qing Zong, Jiayu Liu, Tianshi Zheng, Chunyang Li, Baixuan Xu, Haochen Shi, Weiqi Wang, Zhaowei Wang, Chunkit Chan, and Yangqiu Song. 2025a. Critical: Can critique help llm uncertainty or confidence calibration? *Preprint*, arXiv:2510.24505.

Qing Zong, Zhaowei Wang, Tianshi Zheng, Xiyu Ren, and Yangqiu Song. 2025b. ComparisonQA: Evaluating factuality robustness of LLMs through knowledge frequency control and uncertainty. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 4101–4117, Vienna, Austria. Association for Computational Linguistics.

## A Comparison Traits Details

We distill eight core traits that capture the essential characteristics of recent benchmarks designed for cost-optimal evaluation. These traits concern both the agent's ability to conduct realistic, multi-step, and efficiency-oriented interactions, as well as the infrastructure's capacity to enable scalable and extensible evaluation.

- **Multi-turn Interaction**: The benchmark supports and requires multi-turn dialogue or decision-making. This trait is essential, as recent agentic paradigms increasingly emphasize agent performance in multi-turn interactions. From web search (Geng et al., 2025; Wu et al., 2025a) to GUI-based scenarios (Li et al., 2025; Wang et al., 2025a), agentic performance can only be thoroughly evaluated in dynamic, multi-turn environments.

- **Cost-optimal Planning**: Agent performance is explicitly assessed by cost-optimality. We distinguish *budget awareness*, or maximizing performance under a fixed budget, from *cost-optimality*, which requires reaching the goal with the minimum possible cost.

- **Flexible Cost**: The environment allows user-defined costs or provides sufficiently diverse cost combinations. As discussed in Section 1, a more flexible and diverse cost assignment is crucial for accurately evaluating an agent's genuine cost awareness, such as its sensitivity to varying costs.

- **Tool Use**: The benchmark incorporates or necessitates the use of external tools. This trait is particularly important, as tool use has emerged as a defining capability of modern LLM agents—central to their ability to interact with and act upon the external world (Wang et al., 2024b; Qin et al., 2024b). Evaluating tool use is thus indispensable for understanding an agent's true reasoning and decision-making competence.

- **Dynamic State**: The environment introduces runtime-controllable obstacles that may block the agent and enforce replanning. This design is important because traditional benchmarks often suffer from data leakage risks (Guo et al., 2025; Wu et al., 2025b). In contrast, a dynamic environment can mitigate such risks



**Atomic Tools:**
**Tool 1:** `Decide_Location_Preference`;
**Tool 2:** `Search_Location_Candidates`;
**Tool 3:** `Location_Refinement_Step1`;
**Tool 4:** `Select_Final_Location`

**Composite Tools:**
**Tool 5:** `Location_Preference_and_Search` (= Tool 1+2);
**Tool 6:** `Location_Full_Planning_to_Step1` (= Tool 1+2+3);
**Tool 7:** `Location_All_pipeline` (= Tool 1+2+3+4)
······

Figure 7: Example of sequential tool execution flow in CostBench (task sequence = 4, requiring only one refinement step). Each atomic tool transforms one data type to another, ensuring explicit intermediate states and preventing shortcuts. In this figure, we illustrate only the composite tool starting from the initial state. In practice, a composite tool may commence from any state, encapsulating a sequence of atomic tool operations.

by introducing runtime factors during evaluation, ensuring a fairer and more reliable assessment.

- **Adjustable Difficulty**: The difficulty level can be systematically adjusted, for example by scaling task parameters. A naturally adjustable difficulty hierarchy allows systematic evaluation of model adaptability under varying levels of challenge (Qian et al., 2025b; Zong et al., 2025b; Huang et al., 2025a), facilitating more fine-grained analysis.

- **Customized Details**: The environment offers fine-grained configurability of task details. Enabling this feature makes the benchmark more adaptable to other domains and facilitates broader applicability across diverse evaluation settings (Lu et al., 2025; Wang et al.,

2024b).

- **Scalability**: The benchmark can be automatically expanded in size without additional human annotation and easily generalized to other tasks. This property is crucial, as scalability and expandability are essential for large-scale evaluation (Qian et al., 2025b; Yao et al., 2025; Zong et al., 2025a; Wang et al., 2025e).

| Parameter | Value |
|---|---|
| Task length | 5 |
| Maximum tool steps | 20 |
| Global seed ($S$) | 42 |
| Min atomic cost ($c_{\min}$) | 15 |
| Max atomic cost ($c_{\max}$) | 25 |
| Noise standard deviation ($\sigma$) | 0.1 |

Table 5: Environment details for Table 4. $c_{\min}-c_{\max}$: range for randomizing atomic tool costs (see Figure 1). noise std: standard deviation of Gaussian noise added to composite tool costs. The global seed $S$ is used to control the random seed for all datapoints in a given run, ensuring full reproducibility.

## B  Experiment Details

In this section, we provide a comprehensive overview of the experimental setup and evaluation methodology. We begin by detailing the construction of tools, including their names and descriptions. Next, we explain the simulation-based derivation of greedy and ground-truth trajectories using a tool graph framework, covering both non-blocking and blocking scenarios. We then present the evaluation metrics, including cost gaps, path dissimilarities, user intent alignment, and tool call validity, along with their formulas and boundary conditions. Finally, we describe the mechanisms for implementing blocking events, prove the solvability of blocking scenarios, and outline the randomized yet reproducible environment generation process to ensure fairness and reproducibility in evaluations. An illustration of our environment is shown in Figure 7, while the primary hyperparameter configurations are presented in Table 5.

### B.1  Tool Construction Details

To better evaluate the model's genuine planning ability, we mitigate potential overfitting to existing data (Wu et al., 2025b; Liu et al., 2025). Since biases acquired during training may distort the model's assessment of tool cost-effectiveness, we construct a new set of tools for evaluation. In this section, we mainly introduce two aspects of the tool schema: name construction and description construction.

**Tool Name Construction.**  The tool names in CostBench are constructed systematically to reflect their function and position within the task hierarchy. This structured naming convention is designed to be both human-readable and machine-parsable, facilitating a clear understanding of each tool's purpose.

For **atomic tools**, names are generated based on a `Function_Task_Level` template. The *Function* specifies the core action (e.g., `Decide`, `Search`, `Refine`, `Select`). The *Task* denotes the domain-specific task (e.g., `Location`, `Transportation`). The *Level* indicates the stage in the refinement process (e.g., `Preference`, `Step1`). An example is `Location_Refinement_Step1`, which clearly communicates its role.

For **composite tools**, which are sequences of atomic tools, names are generated algorithmically to summarize the encapsulated workflow. The naming logic depends on the composition of the toolchain. For instance, a two-step tool combining preference decision and initial search is named `[Task]_Preference_and_Search`. A tool that covers the entire workflow from initial search to a specific refinement level is named `[Task]_Search_Planning_to_[Level]`. A complete pipeline from decision to final selection is named `[Task]_Complete_[N]Steps_Pipeline`, where `N` is the number of atomic tools involved. This programmatic approach ensures that every possible contiguous sub-sequence of tools in a task plan has a unique and descriptive name.

**Tool Description Construction.**  The descriptions for tools are crucial for the language model to understand their functionality and parameters. We employ a hybrid strategy for generating these descriptions, combining templated generation with LLM-based refinement to achieve both consistency and naturalness.

For **atomic tools**, descriptions are generated from a template that explains the tool's specific purpose. For example, the description for a refinement tool explicitly states the filtering dimension it applies, such as "Filtered by availability" for `Step1` or "Filtered by location" for `Step2`. This ensures that the model is aware of the distinct value each tool provides.

For **composite tools**, descriptions are generated using *Gemini-2.5-pro* (Google DeepMind, 2025). We provide the LLM with the names and descriptions of the constituent atomic tools and prompt it to synthesize a concise, high-level summary of the overall workflow. For example, instead of merely concatenating the descriptions of a `Search` tool and a `Refine` tool, the LLM might generate a more intuitive description such as, "Searches for options based on user preferences and then filters the results for availability". This process is automated in our framework to generate descriptions for all possible composite tools.

It is worth noting that various aspects of tool information, such as cost, atomic or composite type, and the component tools of composite ones, are incorporated into the tool descriptions.

## B.2   Greedy and Ground-truth Calculation

We emphasize that both the greedy and the ground-truth trajectories are obtained through *simulation* rather than actual tool executions. To support this, we construct a **tool graph** that encodes the possible transitions between agent states.

Formally, we define the tool graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ as follows. Each vertex $v \in \mathcal{V}$ corresponds to an *agent state*, represented as a set of datatypes already available to the agent:

$$v \equiv S = \{d_1, d_2, \ldots, d_n\}, \quad S \subseteq \mathcal{D}, \qquad (1)$$

where $\mathcal{D}$ is the universe of all possible datatypes. Importantly, states evolve *monotonically*: invoking a tool never consumes existing datatypes but only augments the state with new ones, since the model always has access to the full conversation history.

The initial state $S_0$ is determined by the task specification:

$$S_0 = \{\texttt{TimeInfo}\} \cup \{\texttt{LocationPreference} \text{ (if needed)}\} \cup \{\texttt{UserPreference}\}. \qquad (2)$$

For each task, we define a unique target datatype $\tau_{\text{task}}$, representing the final result required by the task, which is also the output of the last atomic tool. Accordingly, any state $S$ that contains $\tau_{\text{task}}$ is regarded as a goal state:

$$S \in \mathcal{V}, \quad \text{goal}(S) \iff \tau_{\text{task}} \in S. \qquad (3)$$

Each edge $e \in \mathcal{E}$ corresponds to the invocation of a tool. A tool $T$ is defined as a typed mapping

$$T : \mathcal{I}(T) \to \mathcal{O}(T), \qquad (4)$$

where $\mathcal{I}(T) \subseteq \mathcal{D}$ denotes the required input datatypes and $\mathcal{O}(T) \subseteq \mathcal{D}$ denotes the output datatypes. An edge exists from state $S$ to state $S'$ if and only if $\mathcal{I}(T) \subseteq S$, in which case

$$S' = S \cup \mathcal{O}(T). \qquad (5)$$

Each edge is further assigned a *cost* $c(T)$, representing the computational or interaction cost of invoking tool $T$. Thus, the transition can be written as

$$(S \xrightarrow{T, c(T)} S') \in \mathcal{E}. \qquad (6)$$

This construction allows us to simulate greedy and ground-truth trajectories over the tool graph without actually executing tools. While such stimulation offers a convenient and controllable proxy, it fundamentally differs from real tool use: being algorithm-driven, it bypasses the core challenge of user-intent understanding. Consequently, although useful for isolating certain planning aspects, it fails to deal with the vagueness, ambiguity, and negotiation with human preferences inherent in real scenarios. Our benchmark therefore becomes essential: cost-aware planning cannot be addressed by traditional algorithmic stimulation, but calls for evaluation of large language models, which uniquely possess the capacity to reason over ambiguity and align with human intent (Yao et al., 2025).

**Scenarios without Blocking.**   For scenarios without blocking, the construction of greedy and ground-truth trajectories is relatively straightforward. In the greedy policy, at each step given the current state $S$, we first compute the set of feasible tools:

$$\mathcal{A}(S) = \{T \mid \mathcal{I}(T) \subseteq S \ \wedge \ \mathcal{O}_{\text{prev}} \cap \mathcal{I}(T) \neq \varnothing\}, \qquad (7)$$

where $\mathcal{I}(T)$ denotes the input types of tool $T$, and $\mathcal{O}_{\text{prev}}$ denotes the output type(s) of the previously invoked tool. For each candidate tool $T \in \mathcal{A}(S)$, we define a utility score that measures cost-efficiency:

$$u(T) = \frac{c(T)}{comp(T)}, \qquad (8)$$

where $c(T)$ is the cost of invoking $T$, and $comp(T)$ is the number of atomic tools contained in $T$. The greedy trajectory is then constructed by repeatedly selecting the tool with the smallest $u(T)$ until reaching a goal state. We explicitly restrict greedy to select only tools in $\mathcal{A}_{\text{chain}}(S)$ because any globally cost-optimal path must satisfy the chain

property (i.e., each tool in the optimal path consumes the previous tool's output); by searching only within this subset of chain-consistent paths, the greedy policy explores exactly those paths that could possibly be globally optimal.

In contrast, the ground-truth trajectory is obtained by directly running *Dijkstra's algorithm* (Dijkstra, 1959) over the tool graph $\mathcal{G}$. Formally, let $\pi(S_0, S)$ denote the minimum-cost path from the initial state $S_0$ to a state $S$. The ground-truth trajectory is defined as

$$\pi(S_0, S^*) = \arg \min_{\substack{S \in \mathcal{V} \\ \tau_{\text{task}} \in S}} \text{Cost}(S_0 \rightsquigarrow S), \qquad (9)$$

where $\text{Cost}(S_0 \rightsquigarrow S)$ is the cumulative cost of the path from $S_0$ to $S$. Unlike the greedy policy, this construction imposes no additional constraint that each tool must consume the immediate predecessor's output, but instead corresponds to the standard shortest-path setting.

**Scenarios with Blocking.** For scenarios with blocking, trajectory construction is more involved due to unpredictable environment changes. Let $S_{t_i}$ denote the state immediately after the $i$-th blocking event, and $\mathcal{G}_{t_i} = (\mathcal{V}_{t_i}, \mathcal{E}_{t_i})$ the updated tool graph at that time.

For the greedy policy, the feasible tool set at state $S_{t_i}$ is

$$\mathcal{A}^{(i)}(S_{t_i}) = \begin{cases} \{T \mid \mathcal{I}(T) \subseteq S_{t_i}\}, & \text{if this is the first step after blocking } i, \\ \{T \mid \mathcal{I}(T) \subseteq S_{t_i} \wedge \\ \quad \mathcal{O}_{\text{prev}} \cap \mathcal{I}(T) \neq \varnothing \}, & \text{otherwise.} \end{cases}$$

$$(10)$$

Here, the first-step relaxation occurs because the environment has changed, which effectively resets the greedy policy with a new initial state $S_{t_i}$. Greedy then selects

$$T^* = \arg \min_{T \in \mathcal{A}^{(i)}(S_{t_i})} u(T), \quad u(T) = \frac{c(T)}{comp(T)}, \quad (11)$$

and continues until reaching a goal state or the next blocking event. The tool graph $\mathcal{G}_{t_i}$ is updated according to the type of blocking: an edge is removed for a *ban tool*, edge weights updated for a *cost change*, unchanged for a *preference change*, and accessible edges updated for a *step length change*.

For the ground-truth trajectory, we construct it in a segmented fashion. Let $S_{t_0} = S_0$ be the initial state. For each blocking interval $[t_i, t_{i+1})$, we

compute the shortest path $\pi_i$ on the updated graph $\mathcal{G}_{t_i}$:

$$\pi_i = \arg \min_{\substack{S \in \mathcal{V}_{t_i} \\ \tau_{\text{task}} \in S}} \text{Cost}_{\mathcal{G}_{t_i}}(S_{t_i} \rightsquigarrow S), \qquad (12)$$

where $\text{Cost}_{\mathcal{G}_{t_i}}(\cdot)$ denotes the cumulative cost in the updated graph.

The full ground-truth trajectory is then obtained by concatenating all segments:

$$\text{GT\_traj} = \pi_0 \oplus \pi_1 \oplus \cdots \oplus \pi_N, \qquad (13)$$

where $\oplus$ denotes path concatenation. This construction ensures that each segment is globally optimal under the current environment, while the full trajectory adapts to blocking events by replanning from the updated state each time.

**Ground-Truth Trajectory under Blocking:**

Initial GT path:

$S_0 \xrightarrow{A} S_1 \xrightarrow{B} S_2 \xrightarrow{C} S_3 \xrightarrow{D} S_4$

Blocking 1 after $B$, replan from $S_2$ :

$\pi_1 : S_2 \xrightarrow{E} S_5 \xrightarrow{F} S_6 \xrightarrow{G} S_7$

Blocking 2 after $F$, replan from $S_6$ :

$\pi_2 : S_6 \xrightarrow{H} S_8 \xrightarrow{I} S_9$

Full GT trajectory:

$\text{GT\_traj} = \pi_0 \oplus \pi_1 \oplus \pi_2$

Above is an example of a ground-truth trajectory under blocking. Each blocking event triggers replanning from the current state. Segments $\pi_i$ correspond to shortest paths computed on the updated tool graph after the $i$-th blocking event. The $\oplus$ operator denotes concatenation of path segments.

### B.3 Metric Details

#### B.3.1 Detailed Implementation and Formulas

**Cost Metrics.** The design intuition of the Cost Gap metric is to directly capture how much more (or less) cost an agent incurs compared to the ground-truth trajectory. Unlike absolute cost values, which depend on the specific scale of the tool graph and are not directly comparable across different settings, the gap formulation provides a simple, interpretable measure of deviation from the optimal reference.

Formally, let $\tau_{\text{agent}} = (a_1, a_2, \ldots, a_m)$ and $\tau_{\text{gt}} = (g_1, g_2, \ldots, g_n)$ denote the sequences of actions in the agent and ground-truth trajectories, respectively. Each action $x$ is associated with a cost $\text{cost}(x)$. The cost gap is defined as:

$$\text{Cost Gap} = \sum_{a \in \tau_{\text{agent}}} \text{cost}(a) - \sum_{g \in \tau_{\text{gt}}} \text{cost}(g). \qquad (14)$$

By construction, CostGap $\geq 0$ since the agent is less cost-efficient than the ground-truth in the static setting, and CostGap $= 0$ indicates perfect alignment with the ground-truth trajectory.

In Table 4, we present two variants of the cost gap metric. The values outside the parentheses denote the average cost gap computed over all samples, whereas the values in parentheses are obtained after excluding two types of progress-awareness failures: (i) extra tool calls, where the agent continues invoking tools after reaching the goal state, and (ii) repeated tool calls, where the agent redundantly invokes the same tool multiple times even after the initial invocation has already succeeded. Our intuition is to decouple progress-awareness from cost-awareness, since extra or repeated calls primarily indicate a lack of progress-awareness rather than deficiencies in cost reasoning. Detailed statistics of such progress-aware errors are reported in Table 6.

Note that we do not compare the aggregated cost optimality over the entire trajectory under blocking scenarios. This is because, at the initial state or immediately after each blocking event, the planning procedure assumes the current goal state is the true final goal. As a result, the cumulative trajectory cost under blocking does not reliably capture cost optimality, since a lower total cost does not necessarily correspond to a better trajectory.

**Path Metrics.** For **Average Edit Distance (AED)**, we use an extended definition of the *Levenshtein Distance* (Backurs and Indyk, 2015; Haldar and Mukhopadhyay, 2011), defined as the minimum number of operations required to transform one path into another. Here, each path is represented as a sequence of tool calls, and the allowed operations are insertion, deletion, and substitution of tool calls. This metric directly measures the structural dissimilarity between the agent's trajectory and the ground-truth trajectory.

$$\text{ED}(\tau_{\text{agent}}, \tau_{\text{gt}}) = \min_{\Gamma \in \mathcal{O}(\tau_{\text{agent}} \to \tau_{\text{gt}})} |\Gamma|, \quad (15)$$

where $\mathcal{O}(\tau_{\text{agent}} \to \tau_{\text{gt}})$ denotes the set of valid edit-operation sequences transforming $\tau_{\text{agent}}$ into $\tau_{\text{gt}}$, and $|\Gamma|$ is the number of operations in $\Gamma$. We compute the Average Edit Distance across all queries in the test set as:

$$\text{AED} = \frac{1}{N} \sum_{i=1}^{N} \text{ED}(\tau_{\text{agent}}^{(i)}, \tau_{\text{gt}}^{(i)}), \quad (16)$$

where $N$ is the total number of queries in the test set.

For **Average Normalized Edit Distance (ANED)**, we follow Li and Liu (2007) and normalize the edit distance by the maximum of the two path lengths. For each query, the Normalized Edit Distance (NED) is defined as:

$$\text{NED}(\tau_{\text{agent}}, \tau_{\text{gt}}) = \frac{\text{ED}(\tau_{\text{agent}}, \tau_{\text{gt}})}{\max\left(|\tau_{\text{agent}}|, |\tau_{\text{gt}}|\right)}, \quad (17)$$

where $|\tau|$ denotes the length of a trajectory. By construction, this metric ranges from 0 to 1, where 0 indicates identical paths and 1 indicates maximal dissimilarity. We then average across all queries:

$$\text{ANED} = \frac{1}{N} \sum_{i=1}^{N} \text{NED}(\tau_{\text{agent}}^{(i)}, \tau_{\text{gt}}^{(i)}). \quad (18)$$

For **Exact Match Ratio (EMR)**, we use a binary indicator that checks whether the agent's trajectory exactly matches the ground-truth trajectory. For each query, the exact match is defined as:

$$\text{EM}(\tau_{\text{agent}}, \tau_{\text{gt}}) = \mathbb{1}\left[\tau_{\text{agent}} = \tau_{\text{gt}}\right], \quad (19)$$

where $\mathbb{1}[\cdot]$ is the indicator function that equals 1 if the two trajectories are identical and 0 otherwise. The Exact Match Ratio is then computed as the proportion of exact matches across the test set:

$$\text{EMR} = \frac{1}{N} \sum_{i=1}^{N} \text{EM}(\tau_{\text{agent}}^{(i)}, \tau_{\text{gt}}^{(i)}). \quad (20)$$

**User Intent Metric.** We report **User Intent Hit Ratio** as the proportion of samples where the agent's predicted final candidate exactly matches the ground-truth answer. This captures agent's ability to understand and intepret human intent in natural language.

**Tool Call Metrics.** For the **Invalid Tool-Use Ratio**, we report the proportion of invalid tool calls over all tool calls.

### B.3.2 Boundary Conditions

**Cost Metrics.** For cost-related metrics, we exclude all datapoints where the agent fails to reach the goal state within the maximum of 20 turns, since comparing a trajectory that stops at an intermediate state with one that reaches the goal is not meaningful. When computing the total cost of each datapoint, we also disregard invalid search attempts. This design ensures that the reported cost reflects the agent's genuine cost awareness, while disentangling it from its tool-use capability.

**Path Metrics.** For path-related metrics, we focus on evaluating genuine cost-optimal planning ability. Trajectories that fail to reach the goal state are excluded from metric computation to ensure that cost awareness is measured only on successful completions. To avoid being influenced by agents' tool calling ability and reflect true cost-optimal planning capability, we also exclude failed tool calls. We thereby disentangle agents' tool calling ability and cost-optimal planning ability.

**User Intent Metrics.** For **User Intent Hit Ratio**, we exclude all samples in which the agent fails to reach the goal state, since such predictions amount to mere guesses rather than answers derived from a complete tool-call trajectory.

**Tool Call Metrics.** For the **Invalid Tool-Use Ratio**, we include all datapoints regardless of outcome. This metric reflects the model's zero-shot tool-use ability on our benchmark, providing an intuition of how such ability impacts cost-optimal planning performance.

Notably, in blocking scenarios, we exclude all samples that do not experience the designated number of blocking events. Including such samples would make the metrics unfair; therefore, we report results based on the intersection of valid sets that meet the blocking criteria.

### B.4 Blocking Details

**Block Trigger Time.** To ensure blocking events occur at strategically meaningful moments throughout the agent's planning process, we employ a dynamic trigger mechanism that adapts to the evolving optimal path length after each environment change. Specifically, let $L_i$ denote the optimal path length after the $i$-th blocking, $t_i$ denote the current execution step, and $B$ denote the total number of planned blockings. For the next $(i + 1)$-th blocking, we calculate the relative trigger position as $\Delta = \lfloor L_i/(B - i) \rfloor$, and set the trigger step to $t_{i+1} = t_i + \max(1, \Delta)$. This strategy distributes blocking events approximately evenly across the remaining optimal trajectory, preventing clustering at the beginning or end of execution. Critically, the trigger step for each blocking is recalculated dynamically after every environment change, as each blocking (especially cost changes or tool removals) may alter the optimal path length $L_i$.

### B.5 Proof of Solvability of Blocking Scenarios

For four block types: ban tool, preference change, cost change and remove tools, we both ensure that the agent could reach the goal state in theory, that is to say, the tool graph remains solvable.

- **Ban Tool** Intuitively, the worst case occurs when all blockings prevent the same state from transitioning to other states. As long as the agent can still move to alternative states, it can reach the goal. Therefore, when the tool that completes the entire task in a single step is unavailable, the maximum number of blockings that can be tolerated is $task\_length - 2$. A rigorous proof is provided later in this section.

- **Preference Change** Since the tool graph does not change, the solvability is unaffected.

- **Cost Change** It is trivial that the tool graph could still be solved, since no edges are removed from the graph.

- **Remove Tools** In this case, we always keep the atomic tools unremoved to maintain solvability.

Rigorous proof for the "Ban Tool" scenario: Consider a linear task of length $n$. Suppose every contiguous interval $[i, j]$ with $1 \le i \le j \le n$ corresponds to a tool, except that the full-length tool $[1, n]$ is banned by default. Then the minimum number of additional tools whose removal destroys *all* partitions of $[1, n]$ is

$$h = n - 1.$$

Consequently, the robustness bound is

$$K_{\text{robust}} = h - 1 = n - 2,$$

meaning that as long as fewer than $n - 1$ tools are banned, the tool graph always remains solvable.

For each $k \in \{1, \ldots, n - 1\}$, consider the partition

$$P_k = \{[1, k], [k + 1, n]\}.$$

Since $[1, n]$ is not allowed, each $P_k$ is a valid partition. Distinct $k$ yield disjoint sets of tools, so to block all partitions one must remove at least one tool from each $P_k$, i.e. at least $n - 1$ tools. Thus $h \ge n - 1$.

Conversely, let

$$H = \{[1, 1], [1, 2], \ldots, [1, n - 1]\}.$$

Any partition of $[1, n]$ must begin with some $[1, k]$ ($1 \leq k \leq n-1$), because $[1, n]$ is banned. Hence every partition uses a tool from $H$. Therefore $H$ is a hitting set of size $n-1$, so $h \leq n-1$. Combining gives $h = n-1$, and thus the robustness bound is $K_{\text{robust}} = n-2$, which completes the proof.

## B.6 Randomized yet Reproducible Environment

To enable systematic evaluation while maintaining experimental variability, we design a fully deterministic environment generation framework controlled by a global seed $S$. This ensures that all random aspects (tool costs, blocking parameters, and trigger timings) are reproducible across runs while remaining sufficiently diverse across different queries.

### B.6.1 Deterministic Cost Assignment

For each query with identifier $q$, and for each tool (atomic or composite) with name $\text{name}(\cdot)$, we deterministically derive random seeds by applying SHA-256 hashing (National Institute of Standards and Technology, 2015) to the tuple $(S, q, \text{name})$:

$$c_{q,a} \sim \text{Uniform}(c_{\min}, c_{\max}; \text{seed} = h(S, q, a)) \tag{1}$$

$$\epsilon_{q,T} \sim \mathcal{N}(0, \sigma^2 k; \text{seed} = h(S, q, T)) \tag{2}$$

$$C_{q,T} = \max\left(1.00, \text{round}\left(\sum_{a \in T} c_{q,a} + \epsilon_{q,T}, 2\right)\right) \tag{3}$$

Here $a$ denotes an atomic tool, $T$ a composite tool with $k$ atomic components, and $h(\cdot)$ the hash-based seed derivation function. The $\sqrt{k}$ scaling in equation (2) reflects empirical observations that execution variance grows sublinearly with pipeline length.

### B.6.2 Deterministic Blocking Parameters

When evaluating agents under a specific blocking type (ban_tool, preference_change, remove_tools, or cost_change), all blocking events are pre-planned using a hierarchical seeding strategy. Given a base seed $S_q$ for query $q$, we derive the seed for the $i$-th blocking as $S_{q,i} = S_q + i \cdot \Delta_s$, where $\Delta_s = 100$ is a fixed interval ensuring seed independence. This seed $S_{q,i}$ deterministically controls the parameters for blocking $i$:

- **Cost Change:** A new global seed $S'$ is sampled from a predefined range using $S_{q,i}$, which then regenerates all tool costs via equations (1)–(3).

- **Preference Change:** $S_{q,i}$ selects an alternative query from the same task category to extract new user requirements and preferences.

- **Remove Tools:** Using $S_q$, we deterministically sample $n$ non-overlapping length intervals $[l_i, r_i]$ from the feasible range $[2, L_{\max}]$, where $n$ is the total count of *remove tools* events and $L_{\max}$ depends on the refinement level. As shorter tools are more prevalent, we impose a minimum width constraint $(r_i - l_i) \geq L_{\max}/2$ to ensure the removal of a sufficient number of tools, thereby increasing the likelihood of modifying the ground-truth plan. At blocking $i$, only composite tools with length exactly equal to $l_i$ (the left endpoint) are removed from the visible tool set, while all atomic tools remain accessible.

- **Ban Tool:** $S_{q,i}$ produces deterministic failure messages, where the banned tool is determined at runtime according to the agent's actual tool call, thereby enabling targeted blocking of the agent's actions.

Among the four types of blocking events, all trigger times are determined based on the optimal path. Therefore, given identical cost assignments, the trigger times remain fixed. The blocking parameters are consistent across all blocking types except for *ban tool*, as they are uniquely controlled by the seed $S_{q,i}$. In contrast, the *ban tool* event is model-dependent—it is triggered at a specified time, and whichever tool the model invokes at that timestep will result in a failure.

Although this design means the environment is not perfectly identical across models, we argue that it does not compromise fairness. If a model follows the correct trajectory, both the ground truth and the agent are equally affected after the ban event, requiring replanning and thus preserving fairness in evaluation. Conversely, if a model already deviates from the optimal path, the model-dependent blocking has no impact, as the datapoint's EM score would already be zero. Consequently, the **EMR** metric remains unaffected. Moreover, as shown in Table 4, **EMR** and **ANED/AED** exhibit opposite trends, suggesting that the latter two effectively capture performance differences as well.

Thus, this hierarchical seeding ensures that: (i) given $(S, q)$ and the blocking type, the entire blocking sequence and all associated parameters are fully determined; (ii) different queries exhibit different blocking patterns even with the same $S$; and (iii)

| Steps | Redundant Tool calls | | Failure Tool Calls | |
|---|---|---|---|---|
| | Repeated Calls ↓ | Extra Calls ↓ | Wrong Parameters ↓ | Unaccessible Calls ↓ |
| *Qwen3-8B* | 1 | 0 | 0 | 6 |
| *Qwen3-14B* | 1 | 7 | 1 | 11 |
| *Qwen3-32B* | 0 | 0 | 19 | 0 |
| *Llama-3.1-8B-Instruct* | 0 | 0 | 19 | 347 |
| *GLM-4.5* | 2 | 0 | 3 | 43 |
| *Deepseek-V3.1* | 16 | 0 | 6 | 314 |
| *Gemini-2.5-pro* | 0 | 2 | 0 | 0 |
| *Claude-sonnet-4.0* | 3 | 0 | 0 | 53 |
| *GPT-4o* | 10 | 0 | 1 | 318 |
| *GPT-5* | 11 | 0 | 0 | 1 |

Table 6: Detailed statistics of failure modes for the tested models in the main experiment (Table 4). Numbers in parentheses represent the ratio (%) to the total number of tool calls.

experimental results are perfectly reproducible, facilitating controlled comparisons across models and further analysis. At the same time, users can increase the benchmark's **diversity** by adjusting random seeds, allowing robustness evaluation under varied conditions and mitigating potential data contamination.

By controlling a randomized yet reproducible environment, we ensure **fair comparison** across different runs and model configurations.

## C Additional Experiment Results

### C.1 The Effect of Random Seeds

To ensure the robustness of our conclusions regarding the cost assignment, which is controlled by a random seed, we conducted three additional experiments using seeds 1000, 2000, and 3000 to evaluate model performance. As illustrated in Figure 8, the LLMs exhibit no significant performance variation across different seeds, with differences in ANED and EMR remaining within 5%, confirming low seed sensitivity. However, Figure 3 reveals that LLMs are highly sensitive to variations in cost assignment. We attribute the observed stability across seeds to the diversity and size of our test set, which mitigates the impact of seed-dependent cost assignments and ensures robust evaluation.

## D Analysis and Discussion

### D.1 Data Distribution

We note a slight imbalance in our final filtered dataset, where the "Location" task contains the fewest instances in both the training and test splits. This skew originates from our generation pipeline: we initially sampled an equal number of raw preference combinations for each task, but a subsequent
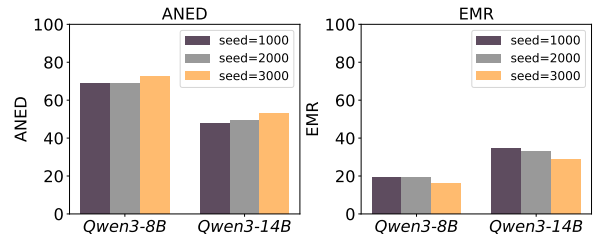


Figure 8: Performance of Qwen3-8B and Qwen3-14B on ANED and EMR **(%)** across different seeds. Results demonstrate low variations with different seeds, with variations in ANED and EMR not exceeding 5% across seeds.

commonsense filter retained a different proportion of combinations for each task. The "Location" task had a lower filter pass rate due to its unique dimensional features, resulting in the final distribution.

We argue that this imbalance is non-critical for our primary objective of evaluating cost-optimal planning. The user preference features are only utilized in the initial step to understand user intent and select the first tool. The subsequent agent trajectory, which comprises the core planning and tool-calling sequence, is independent of these initial preferences. Therefore, while the skew may slightly influence the evaluation of intent understanding, it does not affect the integrity of the cost-optimal planning assessment. For future applications, users of Cost-Bench requiring a more balanced distribution or augmented 'Location' data can readily achieve this by adjusting parameters in our provided codebase.

### D.2 Error Analysis

Building upon the main discussion in Section 5, we provide a detailed examination of all observed error types, including both redundant and failure tool calls. These analyses shed light on how progress-awareness limitations hinder the agents' ability to

| Steps | Redundant Tool calls | | Failure Tool Calls | |
|---|---|---|---|---|
| | Repeated Calls | Extra Calls | Wrong Parameters | Unaccessible Calls |
| Step 1 | Decide_Attraction_Preference | Decide_Shopping_Preference | Attraction_Full_Planning_to_Step1 | Decide_Location_Preference |
| Step 2 | Attraction_Preference_and_Search | Search_Shopping_Candidates | Attraction_Finalize_from_Step1_3Steps | Search_Location_Candidates |
| Step 3 | Attraction_Refine_to_Step2 | Shopping_Finish_from_Step1_3Steps | - | Location_Refinement_Step2 |
| Step 4 | Select_Final_Attraction | Shopping_Refinement_Step1 | - | - |
| Step 5 | - | Shopping_Refinement_Step2 | - | - |
| Step 6 | - | Select_Final_Shopping | - | - |

Table 7: Representative examples of the four failure modes, taken from *Qwen3-14B*. Tool calls highlighted in green indicate the step where the goal state is reached, while those in red mark the erroneous calls. Each model in Table 4 exhibits at least one of these error types.

plan and act cost-efficiently.

### D.2.1 Redundant Calls

As defined in the Section 5, redundant calls consist of two subtypes: (1) *Repeated calls*, where the model invokes the same tool multiple times even after a successful call, and (2) *Extra calls*, where the model continues invoking tools after reaching the goal state.

Such behaviors commonly occur when models fail to maintain an internal notion of task completion or overlook that an equivalent operation has already been performed. In Table 6, we report the frequency of these redundant patterns across models. The overall ratio of redundant calls correlates with the noise observed in the raw cost metrics, as discussed in Section 5.

Table 7 illustrates representative examples:

- **Repeated Calls:** The agent redundantly invokes both "Decide_Attraction_Preference" and "Attraction_Preference_and_Search," even though the latter already subsumes the former's function. This results in unnecessary tool usage and inflated total cost.

- **Extra Calls:** The model successfully reaches the goal state by Step 3 but continues to call additional tools instead of terminating the process. This behavior suggests a lack of awareness regarding task completion.

### D.2.2 Failure Calls

Apart from redundant behavior, we observe two major categories of *failure calls*, which directly lead to invalid tool invocations:

**(1) Wrong Parameters.** These errors arise when the model specifies incorrect tool names or malformed parameter formats. For instance, a model may call a non-existent tool like "Attraction_Finish" instead of the correct "Attraction_Finish_from_Step1_3Steps." This reflects inadequate grounding in the available tool schema

and incomplete adherence to the prompt's enumerated tool list (see Figure 14).

**(2) Unaccessible Calls.** These represent the most severe failure mode. An unaccessible call occurs when the model invokes a tool whose input dependencies have not yet been satisfied. For example, after executing "Search_Location_Candidates," a model may directly call "Location_Refinement_Step2" without completing "Step 1", which provides the required structured input. Such cases reveal a fundamental failure to reason about task dependencies, despite the explicit prompt instructions outlining the correct execution sequence.

As shown in Table 6, unaccessible calls dominate among all failure cases across models. This suggests that models struggle to maintain a consistent internal state of intermediate results or to map current progress to the corresponding valid action space.

### D.3 Summary and Implications

Across all analyzed cases, the dominant failure modes, repeated, extra, wrong-parameter, and unaccessible calls, point to a shared underlying limitation: insufficient progress awareness. Models often fail to track which subgoals have been achieved and which inputs are currently available, leading to redundant or logically inconsistent tool invocations. This lack of situational grounding undermines cost sensitivity and prevents the models from performing truly cost-optimal planning, even when they possess the necessary cost-related reasoning capabilities.

### D.4 More Discussions

**Limited gains from parameter scaling in Qwen models.** Scaling up model size leads to performance improvements, but the gains are limited. Within the Qwen family, moving from *Qwen3-8B* to *Qwen3-14B* yields a clear boost, whereas the

improvement from *Qwen3-14B* to *Qwen3-32B* is marginal. This indicates diminishing returns from parameter scaling alone, pointing to the need for algorithmic or architectural advances beyond sheer model size.

**Models demonstrate strong performance in human intent understanding.** As shown in Table 4, most models exhibit excellent capabilities in understanding user intent. However, as pointed out by UserBench (Qian et al., 2025b), there remains a considerable gap before models can dynamically adapt to evolving user needs. This observation highlights the importance of dynamic benchmarks that can capture the evolution of conversational states across multi-turn interactions.

## E  Data Annotation

**Query Validation.** We enlisted three PhD-level researchers, all co-authors with expertise in NLP, to annotate a sample of 200 travel queries across six scenarios: accommodation, transportation, attractions, location, dining, and shopping. A screenshot of the annotation interface is shown in Figure 9. The annotators achieved individual accuracy rates of 97%, 95%, and 98% against the reference answers, confirming the high quality of the annotations. They also verified that the dataset contains no offensive language or personal information.

**Coverage Rate Calculation.** For the calculation of the coverage rate described in Section 5.1, we adopt the LLM-as-a-judge approach. To verify its reliability, we manually inspected 10 sampled cases for each setting and found the automatically extracted results to be fully consistent with the human annotations. Therefore, we are confident in the credibility of the reported results.

**Query ID: <Query00001>**

Task: location

Requirement: I've been craving an adventure in a bustling urban environment, not just any town but a sprawling metropolis teeming with life and diversity. My heart is set on a place with a rich and storied past, where history seeps from its every corner and tradition is woven into the fabric of daily life. I'm particularly drawn to destinations renowned for their breathtaking architecture, where each building tells a story and serves as a testament to human creativity and achievement. This is the kind of cityscape I long to explore, one that offers a deep dive into the past while standing tall in the present.

**Question 1: 🗒️ category**

⊖ city

⊖ seaside

⊖ mountain

⊖ village

**Question 2: 🏆 tier**

⊖ major_metropolis

⊖ mid_sized_city

⊖ small_town

⊖ secluded_area

**Question 3: 🎨 style**

⊖ historic_and_traditional

⊖ modern_and_cosmopolitan

⊖ natural_and_serene

⊖ entertainment_and_vibrant

**Question 4: 📦 feature_package**

⊖ architectural_marvel

⊖ religious_center

⊖ signature_theme_city

⊖ culinary_capital

Figure 9: Annotation Screenshot

**Prompts used in query construction stage**

**User preference validation**
You are a helpful assistant for validating commonsense conflicts in user queries. Given a set of user requirements, determine whether there are any commonsense conflicts among them. Apply strict checking: even minor inconsistencies should be marked as conflicts. Your response must be either **conflict** or **no conflict**, nothing else.

Example: User prompt: Task: Location search. User requirements: " 1. I want the Location category to be 'city'.
2. I want the Location tier to be 'secluded_nature'.
3. I want the Location style to be 'adventure'.
4. I want the Location features to include 'nightlife_central'.
Generated response: **conflict**

User prompt:

**Prompt construction**
You are a helpful assistant for generating queries.
Please generate a search query for a [task] task based on detailed user requirements. The user requirements will be comprised of four dimensions (Category requirement, Tier requirement, Style requirement, Feature package requirement). The query should be written as a long, self-contained user statement that clearly describes the user's needs and intentions. You should follow these rules:

1. The query clearly discribe the user requirements without any possibilities of misunderstanding. For each requirement dimension, you should clearly distinguish the user required one from any other possible candidates in the generated query. Possible candidates are listed below:
Category requirement: [category_candidates]
Tier requirement: [tier_candidates]
Style requirement: [style_candidates]
Feature package requirement: [features_candidates]
2. You should use human-like language to express the user requirements. That is to say, you shouldn't use the exact word to describe the user requirements. Instead, you should paraphrase and rephrase the requirements to imply the user needs in a natural way. For example: For if the user has a 'luxury' requirement, then you could say something like 'money is totally not a concern, and I want a extravagant experience'. For some special cases (proper nouns), you can use the exact wording.
3. The query should be concise and to the point, avoiding unnecessary details or overly complex sentences.
4. All the information you could use is from the user preferences. The location and time information are just meaningless placeholders.

Please **DO NOT GENERATE ANYTHING OTHER THAN THE QUERY**.

Figure 10: The prompts used in our query construction stage. All the word surrounded with "[ ]" would be replaced with real parameters in construction time.

Figure 11: The prompts used in our path extraction stage. All the word surrounded with "[ ]" would be replaced with real parameters in construction time.

**An example user query**

**User Query**:
I've been craving an adventure in a bustling urban environment, not just any town but a sprawling metropolis teeming with life and diversity. My heart is set on a place with a rich and storied past, where history seeps from its every corner and tradition is woven into the fabric of daily life. I'm particularly drawn to destinations renowned for their breathtaking architecture, where each building tells a story and serves as a testament to human creativity and achievement. This is the kind of cityscape I long to explore, one that offers a deep dive into the past while standing tall in the present.

**Corresponding user preference**:
Location main category: city
Location scale tier: major_metropolis
Location style priority: historical_and_traditional
Location feature package: architectural_marvel

Figure 12: An example user query in CostBench.

| Dimension | Values |
|---|---|
| Category | flight, train, bus, car rental |
| Tier | luxury class, business class, standard class, budget class |
| Style | speed priority, comfort priority, scenic route, schedule flexibility priority |
| Feature Package | onboard connectivity and power, full meal and beverage service, special luggage allowance, lie flat or sleeper facility |

Table 8: The test set feature list for task transportation.

```
Name: Decide_Transportation_Preference
Description: With this atomic tool, you describe a transportation by choosing
    its category, tier, style, and feature package according to the user
    requirements. The tool then gives you a unique label that represents your
    exact combination of transportation preferences. This tool has a cost of
    20.06 units. The output type of this tool is TransportationPreference.
Parameters:
    - Type: object
    - Properties:
        - LocationPreference:
            - Type: String
            - Description: A LocationPreference object identifier representing
                user preferences for location selection
        - TransportationCategory:
            - Type: String
            - Enum: [
                flight,
                train,
                bus,
                car_rental
            ],
            - Description: A transportation category enumeration value
                specifying the main type/category for transportation selection.
                Available options: flight, train, bus, car_rental
        - TransportationTier:
            - Type: string
            - Enum: [luxury_class, business_class, standard_class, budget_class
                ],
            - Description: A transportation tier enumeration value specifying
                the quality/price level for transportation selection. Available
                options: luxury_class, business_class, standard_class,
                budget_class.
        - TransportationStyle:
            - Type: string
            - Enum: [speed_priority, comfort_priority, scenic_route,
                schedule_flexibility_priority],
            - Description: A transportation style enumeration value specifying
                the preferred style/approach for transportation selection.
                Available options: speed_priority, comfort_priority,
                scenic_route, schedule_flexibility_priority.
        - TransportationFeaturePackage: {
            - Type: String,
            - Enum: [onboard_connectivity_and_power,
                full_meal_and_beverage_service, special_luggage_allowance,
                lie_flat_or_sleeper_facility.
            ],
            - Description: A transportation feature package enumeration value
                specifying additional features/services for transportation
                selection. Available options: onboard_connectivity_and_power,
                full_meal_and_beverage_service, special_luggage_allowance,
                lie_flat_or_sleeper_facility.
    - Required: [LocationPreference, TransportationCategory, TransportationTier,
        TransportationStyle, TransportationFeaturePackage]
```

Figure 13: An example tool schema for the tools we used in the CostBench.

**Inference Prompt Part 1**

You are an AI assistant for planning task-related schedules.

<Task description>
Your only objective is to obtain the required information (goal type: 'TravelLocation', represent by a unique ID '<LocationCandidate{Candidate_ID}>') by following the tool path with the **LOWEST TOTAL COST**. The task consists of 4 parts: Deciding Preference, Searching Candidates, Refining Options, Final Recommendation. In the refinement stage, you should take charge of filtering the Location candidates. You should refine the possible candidate set from these 1 dimensions: availability and seasonal suitability. Note that the order of the refinement steps is fixed as specified above, and using other order will result in incorrect behavior.
</Task description>

<Tool description>
1. **Tool Cost**. Each tool call has a predefined cost listed in the tool description.
2. **Tool Input and Output Types**. Each tool defines its input types through its parameters (the parameter name indicates the data type) and its output type in its description.
3. **Tool Dependencies**. Some tools depend on others through their input/output types. Carefully read each tool's input/output fields and description before calling the tool.
4. **Data types**. Each Tool has a list of input data types and a output data type. You should infer LocationCategory, LocationTier, LocationStyle, LocationFeaturePackage, TimeInfo from the user query. For other data types, you only obtain them when a certain tool explicitly returns them. The data types are specially designed, and using them incorrectly will result in incorrect behavior.
5. **Atomic vs Composite Tools**. The tools available could categorized into atomic tools and composite tools, which is specified in the tool description. An atomic tool performs a single and unseparable operation. A composite tool chains multiple atomic tools in sequence and lists its component atomic tools in its description. The cost of a composite tool is specified in its description. Inputs/outputs of a composite tool follow the component chain. Despite being multi-step internally, it still counts as ONE tool call and must obey the one-tool-per-step rule. The cost of a composite tool might be higher or lower than the sum of its component atomic tools.
6. **Sample Atomic Tool Sequence**. For this task, the basic atomic tool calling sequence is: Decide_Location_Preference, Search_Location_Candidates, Location_Refinement_Step1, Select_Final_Location. You should replace some atomic tools with composite tools if that reduces cost. You must then compare all possible equivalent tool-calling paths and pick the one with the lowest total cost.
</Tool description>

Figure 14: The first part of the prompt used to benchmark agents during runtime. The example shown corresponds to a task sequence of length 4. The filtering steps and demonstration examples are dynamically generated according to the task sequence length. This part is concatenated with Part 2 in Figure 15 to form the complete prompt.

<Expected workflow>
1. **Explain your reasoning.** Write out your plan clearly, showing how you'll minimize cost. To ensure the optimality of your plan, you should list out all possible tool-calling paths, sum up the cost of each path, and then select the path with the lowest cost.
2. **Execute your plan.** Right after the explanation, invoke the required tool. Do not describe or print the tool call in text, just make the call directly.
3. **Adapt and continue.** You should always keep an eye on the environment. On every step of execution, you should always check if anything about the tool changes (e.g. cost, availability, etc.). If something goes wrong or changes, adapt and continue along the most cost-optimal path.
</Expected workflow>

<Important rules>
- **Cost is the most important.** Your performance is evaluated solely based on the total cost of tool calls upon reaching the goal state. Always pick the cost-minimal tool path. If there are two path with the same cost, you should pick the one with the least number of tool calls.
- **One tool per step.** You may only call one tool at a time and SHOULD NOT call multiple tools in one request. If you try to call multiple, only the first will count.
- **Exact parameters.** Use the provided values exactly as given (e.g., if "<TimeInfo00000>" is given, the "TimeInfo" parameter must be "<TimeInfo00000>", if "<LocationPreference00000>" is given, the "LocationPreference" parameter must be "<LocationPreference00000>"). - **Final answer format.** Once you obtain the "Candidate_ID" representing your goal type, stop calling tools immediately and return the answer in this exact format: "<answer> <LocationCandidate{Candidate_ID}> </answer>". Only incorporate the "<answer>", "</answer>" tag when you want to provide the final answer. If you output the format, your conversation would be terminated.
</Important rules>

<example>
Here is an example of how to plan your tool call paths in a cost-optimal way for your reference. You should adapt to the task and available tools instead of memorizing this example.

Given that: 1. The basic atomic tool calling sequence is: A(Cost: Cost_A), B(Cost: Cost_B), C(Cost: Cost_C), D(Cost: Cost_D). 2. The available tools are: A(Cost: Cost_A), B(Cost: Cost_B), C(Cost: Cost_C), D(Cost: Cost_D), AB(Cost: Cost_AB), BC(Cost: Cost_BC), CD(Cost: Cost_CD), ABC(Cost: Cost_ABC), BCD(Cost: Cost_BCD). Composite tools are those whose names contain at least two letters; each letter represents an atomic tool included within the composite, while their costs are not necessarily the sum of their component atomic tools (e.g., 'AB' is equivalent in effect to performing A then B, but Cost_AB may differ from Cost_A + Cost_B).

Then you should list out all possible tool calling paths first:
<path> 1. A(Cost: Cost_A) -> B(Cost: Cost_B) -> C(Cost: Cost_C) -> D(Cost: Cost_D). Total Cost: Cost_A + Cost_B + Cost_C + Cost_D.</path> <path> 2. AB(Cost: Cost_AB) -> C(Cost: Cost_C) -> D(Cost: Cost_D). Total Cost: Cost_AB + Cost_C + Cost_D.</path> <path> 3. A(Cost: Cost_A) -> BC(Cost: Cost_BC) -> D(Cost: Cost_D). Total Cost: Cost_A + Cost_BC + Cost_D.</path> <path> 4. A(Cost: Cost_A) -> B(Cost: Cost_B) -> CD(Cost: Cost_CD). Total Cost: Cost_A + Cost_B + Cost_CD.</path> <path> 5. AB(Cost: Cost_AB) -> CD(Cost: Cost_CD). Total Cost: Cost_AB + Cost_CD.</path> <path> 6. ABC(Cost: Cost_ABC) -> D(Cost: Cost_D). Total Cost: Cost_ABC + Cost_D.</path> <path> 7. A(Cost: Cost_A) -> BCD(Cost: Cost_BCD). Total Cost: Cost_A + Cost_BCD.</path>
At last, you should select and execute the path with the lowest total cost. </example>

Figure 15: The second part for agent runtime prompt.