
Recovery-Bench: Evaluating Agentic Recovery from Mistakes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce Recovery-Bench, a novel evaluation benchmark designed to measure
2 the ability of language model agents to recover from prior mistakes and corrupted
3 context. The capability to recover is critical for long-horizon tasks. By initializing
4 agents with environments and action histories derived from failed trajectories,
5 we demonstrate that models’ recovery performance differs markedly from their
6 performance in fresh contexts. Notably, GPT-5, which underperforms in our clean
7 settings for solving terminal tasks, significantly improves in recovery scenarios.
8 Our results suggest that resilience to context pollution is a distinct metric of model
9 and agent performance, underscoring the importance of designing benchmarks that
10 reflect realistic, messy deployment environments.

11 1 Introduction

12 Language model agents have demonstrated exceptional capabilities in completing short-horizon tasks
13 at a human level, like searching and synthesizing information or fixing a small bug [Anysphere, 2025,
14 Anthropic, 2025, OpenAI, 2025b]. However, moving from short, well-scoped objectives to complex
15 long-horizon tasks exposes a distinct brittleness in the current agent designs. When tasks span many
16 steps, involve interdependent subtasks, and require extensive tool use, even frontier models inevitably
17 make errors somewhere along the trajectory. These errors are not isolated: they can aggregate across
18 steps and alter both the external environment (e.g., files, configurations, databases) and the agent’s
19 internal working context (e.g., memory and message history). As a result, success on long-horizon
20 tasks hinges not only on raw problem-solving skill but also on an agent’s ability to recognize, contain,
21 and recover from prior mistakes.

22 These capabilities of recovery are especially consequential because, in practice, agents are often
23 launched into dirty environments. Agent practitioners re-run agents after a partial failure, resume
24 from mid-trajectory snapshots, or hand over a messy workspace to a different model or version. The
25 residue of prior attempts becomes part of the running agent’s context. This phenomenon, which we
26 call *context pollution*, describes the persistence of erroneous actions or corrupted states that carry
27 forward into subsequent attempts. Context pollution can misdirect competent agents, causing them to
28 compound errors (by trusting misleading traces) or to expend capacity on untangling stale state rather
29 than making forward progress. Notably, as more agents and frameworks become model agnostic
30 [Anysphere, 2025, Packer et al., 2023, OpenAI, 2025a], the user can plug in and even change the
31 underlying LLMs. This results in context pollution from responses generated by other models.

32 More context pollution occurs in agentic tasks, and the capability to recover is becoming increasingly
33 crucial; however, standard benchmarks still evaluate agents from fresh states. Such setups are valuable
34 for measuring competence of problem-solving from scratch, but they systematically under-measure
35 recovery: the capacity to diagnose inherited problems and repair corrupted state. In particular,

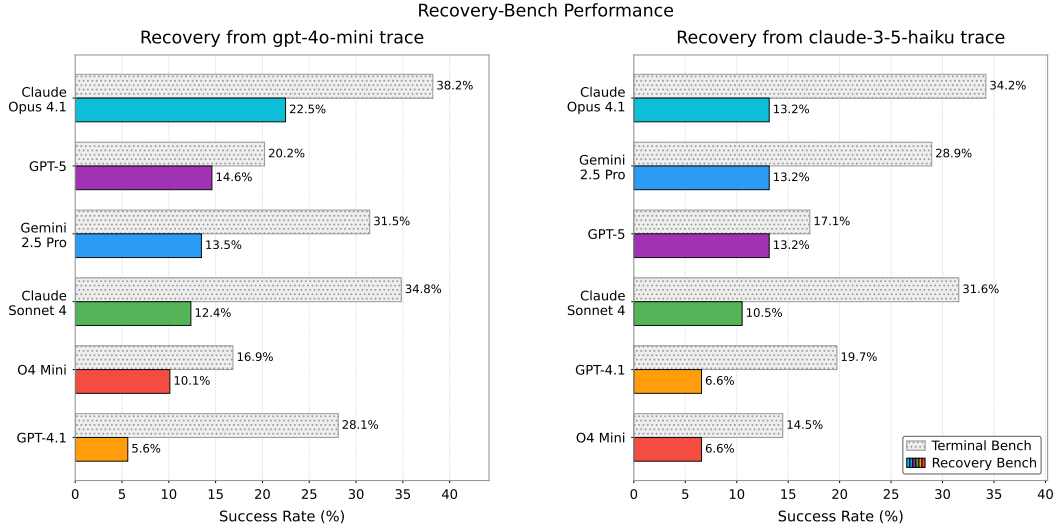


Figure 1: Detailed performance comparison between *Recovery-Bench* and fresh *Terminal-Bench* experiments. **Models are ranked by their performance on *Recovery-Bench*.** Obviously, the rankings of *Recovery-Bench* are different from those of *Terminal-Bench*, highlighting recovery as a distinct capability. Overall, Claude Opus 4.1 performs the best, and GPT-5 also demonstrates great recovery capability despite having a much lower TB performance.

36 standard evaluations that omit prior failed trajectories cannot reveal how sensitive different agents are
 37 to context pollution or whether they can leverage partial progress without being misled by it.

38 To bridge this gap, we introduce *Recovery-Bench*, a benchmark that evaluates agents in realistic,
 39 “lived-in” environments seeded with artifacts from prior failed attempts from arbitrary agents. Instead
 40 of always beginning from a clean repository or a freshly started environment, agents are initialized
 41 in states that include (i) concrete environment residues such as modified files, broken tests, and
 42 misapplied patches, and (ii) textual residues such as debugging trails and prior reasoning traces.
 43 The benchmark asks agents to complete the original task *from these inherited states*. In doing so,
 44 *Recovery-Bench* targets a capability that clean-slate evaluations largely miss: robust recovery under
 45 context pollution.

46 **Implementation on Terminal-Bench (TB).** Concretely, we build our *Recovery-Bench* idea on a
 47 benchmark called *Terminal-Bench* [Terminal-Bench Team, 2025], which consists of various tasks
 48 performed within a terminal. *Recovery-Bench* is a drop-in evaluation mode that preserves TB’s tasks
 49 and harness while altering only the *initialization* state. For each TB task, we select a previously
 50 recorded failed agent attempt and reconstruct a polluted starting point by first restoring the post-failure
 51 workspace snapshot with previous commands, thereby carrying forward partially applied patches,
 52 failing or disabled tests, and other inconsistent edits. Then, we attach traces from that attempt to
 53 the agents as message history. The agent is then asked to complete the original TB task from this
 54 inherited state under the same budget and judgment TB uses in the clean-slate setting.

55 To control what the recovery agent inherits, the harness parameterizes context pollution along two
 56 axes. The *environment residue* axis governs the workspace state the agent sees at time 0 (e.g.,
 57 modified source files, test artifacts, build outputs), while the *trace residue* axis governs how much
 58 of the prior trajectory is surfaced in-context (e.g., commands-only, compact summaries, or fuller
 59 message/command traces). We describe in detail how we generate the trajectories and how we
 60 parameterize context pollution in Section 3.

61 **Experimental results on TB.** We evaluate *Recovery-Bench* with a range of frontier models on
 62 TB recovery tasks constructed from the failed trajectories above. Relative to the standard TB “fresh
 63 start” setting, aggregate performance drops significantly when agents inherit a polluted state, and
 64 model rankings shift. This shift clearly indicates that capabilities to recovery is independent from
 65 capabilities to solve the tasks.

66 **Contributions.** We summarize our contributions as follows:

- 67 • We introduce *Recovery-Bench*, a benchmark that evaluates agents in realistic environments
68 seeded with failed trajectories from previous attempts.
- 69 • We provide empirical evidence that recovery capabilities are distinct from raw problem-
70 solving capabilities.
- 71 • We discuss implications for agent design that mitigate context pollution and improve reli-
72 ability in long-horizon settings.

73 2 Related Work

74 **Agentic Evaluation** Multiple benchmarks test agentic skills across software engineering, tool use,
75 and real-web interaction. SWE-bench [Jimenez et al., 2024] is a suite of real GitHub issues with
76 executable test harnesses that stress multi-file reasoning. τ -bench [Yao et al., 2024] is a benchmark
77 that simulates user-agent dialogues with domain APIs with innovative (simulated) human users.
78 WebArena [Zhou et al., 2024] is a collection of self-hosted realistic websites with execution-based
79 success criteria. However, all such evaluations initialize agents from clean states and focus on raw
80 problem-solving capabilities. *Recovery-Bench* contrasts this by holding tasks fixed while initializing
81 from post-failure states to isolate recovery ability.

82 **Long Context and Context Pollution** Reasoning with long contexts is essential to all real-world
83 agentic tasks, as the average agent interactions are getting longer. Lost in the middle [Liu et al.,
84 2023] studies positional bias of long context generation, suggesting model’s inability to reason
85 from the middle of context. Other prior works [Kamradt, 2025, OpenAI, 2025, Bai et al., 2024]
86 evaluate how LLMs perform under long and sometimes irrelevant, and possibly polluted contexts.
87 As long context performance evolves for frontier models, many of the aforementioned long text
88 benchmarks and studies become saturated and obsolete. In addition, prior work rarely studies long
89 context performance in the agentic space. *Recovery-Bench* introduces a new dimension for context
90 pollution, where the interactions in the context are coherent but contain errors.

91 3 *Recovery-Bench* Design

92 To curate data for *Recovery-Bench*, we build on Terminal-Bench (TB), a multi-step terminal-use
93 benchmark with a diverse set of tasks. For example, TB remains challenging, especially for weaker
94 models, so collecting rich, realistic failure trajectories is straightforward.

95 Because Terminal-Bench (TB) is a growing benchmark, we fix the evaluation to commit 55abc0d,
96 which contains 104 diverse terminal-use tasks. For each task, we generate and record trajectories by
97 running weaker baseline models such as gpt-4o-mini or claude-3-5-haiku from a clean TB state,
98 capturing the complete interaction between the agent and the terminal: commands sent to the terminal
99 and all messages. We then filter for failed trials with more than ten interactions to ensure non-trivial
100 attempts. These trials, comprising the terminal command sequence and the resulting post-failure
101 workspace, serve as the dataset for recovery initializations. To replay to the end terminal state of the
102 failed trajectories, we simply re-ran all the commands from a fresh container. In parallel, we retain
103 the corresponding interaction trace and surface it to the recovery agent being evaluated.

104 In addition to the default setting in *Recovery-Bench*, where we replay the post-failure state and surface
105 the full interaction history to the recovery agent, we also study controlled variants that toggle which
106 residues are inherited. In the first variant, we restore only the post-failure workspace snapshot and
107 provide no trace from the prior attempt. In the second, we pair the same workspace with a compact
108 summary of earlier actions. These setups vary only in the form and amount of context pollution,
109 allowing us to isolate whether models can recover from corrupted states or even with misleading
110 trajectories.

111 4 Experiments

112 We run *Recovery-Bench* on the Terminus-1 agent [Merrill and Shaw, 2025] from gpt-4o-mini (89
113 trajectories) and claude-3-5-haiku (74 trajectories), respectively, with six frontier models: Claude
114 Opus 4.1, Claude Sonnet 4, GPT-5, GPT-4.1, O4 Mini, and Gemini 2.5 Pro.

115 A potential concern when evaluating frontier models on trajectories generated by weaker models
 116 is that such trajectories may fall out of the distribution of the stronger model. Nevertheless, we
 117 argue that our recovery setting is highly valuable for both agents and compound AI systems [Zaharia
 118 et al., 2024]. First, many agents are inherently model-agnostic, allowing users to switch models
 119 seamlessly within a session, and the newer model will inevitably work with the residual context.
 120 Second, compound AI and multi-agent systems [Anthropic, 2025] often involve multiple LLMs,
 121 making it likely that one model will need to interpret and correct errors introduced by another model.

122 4.1 How does recovery compare to clean state performance?

123 Agent recovery tasks are challenging even for the strongest models. We compare the performance
 124 of different frontier models on Terminal-Bench and *Recovery-Bench* in Figure 1. In the original
 125 Terminal-Bench setting with gpt-4o-mini traces, models scored an average of 28.3%, with the best
 126 model, Claude Opus 4.1, reaching 38.2%. On Recovery-Bench, models score an average of only
 127 13.1%. Compared to their performance on the original Terminal-Bench, they show a 51.8% relative
 128 decrease in accuracy. Similar trends have been observed for claude-3.5-haiku traces.

129 Additionally, we show that the rankings on Recovery-Bench differ from standard Terminal-Bench
 130 settings, indicating that recovery represents an orthogonal capability. For example, Claude 4 Sonnet
 131 achieves the second-highest Terminal-Bench score at 34.8% but ranks fourth on Recovery-Bench. O4
 132 Mini has the lowest Terminal-Bench score but performs the same or better than GPT-4.1 on agent
 133 recovery. And interestingly, GPT-5, which performs worse compared to GPT-4.1 on the original
 134 Terminal-Bench, has much better performance on *Recovery-Bench*.

135 4.2 How does context residual impact recovery?

	Terminal-Bench	Env Replay Only	Env + Summary	Env + Full Message
Success (%)	34.8	27.0	23.6	12.4
Avg Steps	23.5	24.1	21.3	16.9

Figure 2: Performance comparison between different context residuals. We fix the agent model to be Claude Sonnet 4, and the error traces are generated by gpt-4o-mini.

136 *Recovery-Bench* also studies how different *context residual* setups, i.e., how much of the erroneous
 137 trajectories remaining inside the agent’s context window, impact the recovery performance. We
 138 vary the context residual by toggling what the recovery agent inherits from the failed attempt. We
 139 consider three recovery settings beyond the clean Terminal-Bench baseline: *Env Replay Only*, where
 140 the post-failure workspace is restored with no trace of prior interactions set in the recovering agent’s
 141 context; *Env + Summary*, with the same post-failure states plus a concise summary of earlier actions;
 142 and *Env + Full Message*, where we include additionally the complete command and message history.
 143 These conditions isolate the effects of corrupted state versus misleading traces while holding the
 144 underlying TB task fixed (the agent is still able to see the task description from TB).

145 The result is surprising: adding more context hurts performance. Under *Env + Full Message*, the
 146 agent performs the worst despite having the complete interaction history. We show an interesting
 147 trace in Appendix A. The failure to recover from previous errors mainly due to either the agent
 148 mysteriously trusts previous interactions and not fix previous actions, or the agent gets stuck in the
 149 error loop of previous mistakes. On the contrary, the Env Replay Only setting achieves the highest
 150 score, approaching the original TB score. This indicates that, although the evaluation is initialized
 151 with the failed state, many tasks are still directly recoverable without explicit fix.

152 5 Conclusion

153 In this paper, we introduce *Recovery-Bench*, a benchmark focuses on agent recovery. Implemented
 154 on top of Terminal-Bench, *Recovery-Bench* evaluates different models’ recovery capabilities under
 155 different settings. Across frontier models, recovery is remarkably harder than clean starts and reorders
 156 rankings, showing that resilience to context pollution is a distinct capability. We view *Recovery-Bench*
 157 as a compact, practical benchmark for selecting and improving long-horizon agents.

158 **References**

159 Anthropic. Claude code: Agentic coding tool. <https://github.com/anthropics/claude-code>,
160 2025. Announced February 24, 2025 (research preview); general availability May 22, 2025;
161 accessed 2025-09-03.

162 Inc. Anysphere. Cursor: The ai code editor. <https://cursor.com/>, 2025. Accessed: 2025-09-03.

163 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du,
164 Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. Longbench: A bilingual,
165 multitask benchmark for long context understanding, 2024. URL [https://arxiv.org/abs/2308.](https://arxiv.org/abs/2308.14508)
166 14508.

167 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
168 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL
169 <https://arxiv.org/abs/2310.06770>.

170 Greg Kamradt. Llmtest_needleinahaystack: A needle in a haystack analysis to test in-context retrieval
171 ability of long-context llms. https://github.com/gkamradt/LLMTest_NeedleInAHaystack,
172 2025. Accessed: 2025-09-03.

173 Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,
174 and Percy Liang. Lost in the middle: How language models use long contexts, 2023. URL
175 <https://arxiv.org/abs/2307.03172>.

176 Mike Merrill and Alex Shaw. Terminus: a research-preview agent for terminal-bench. <https://www.tbench.ai/terminus>, May 2025. Published May 19, 2025; accessed 2025-09-03.

178 OpenAI. Openai mrcr: Long-context multi-needle recall dataset. [https://huggingface.co/](https://huggingface.co/datasets/openai/mrcr)
179 [datasets/openai/mrcr](https://huggingface.co/datasets/openai/mrcr), 2025. Accessed Sep 3, 2025.

180 OpenAI. openai/codex: Lightweight coding agent that runs in your terminal. [https://github.com/](https://github.com/openai/codex)
181 [openai/codex](https://github.com/openai/codex), 2025a. Published on GitHub on April 16, 2025; accessed 2025-09-03.

182 OpenAI. Deep research system card. [https://cdn.openai.com/deep-research-system-card.](https://cdn.openai.com/deep-research-system-card.pdf)
183 [pdf](https://cdn.openai.com/deep-research-system-card.pdf), February 2025b. Accessed: 2025-09-03.

184 Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E
185 Gonzalez. MemGPT: Towards LLMs as Operating Systems. *arXiv preprint arXiv:2310.08560*,
186 October 2023.

187 Terminal-Bench Team. Terminal-bench: A benchmark for ai agents in terminal environments, Apr
188 2025. URL <https://github.com/laude-institute/terminal-bench>.

189 Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for tool-
190 agent-user interaction in real-world domains, 2024. URL <https://arxiv.org/abs/2406.12045>.

191 Matei Zaharia, Omar Khattab, Lingjiao Chen, Jared Quincy Davis, Heather Miller, Chris Potts, James
192 Zou, Michael Carbin, Jonathan Frankle, Naveen Rao, and Ali Ghodsi. The shift from models to com-
193 pound ai systems. <https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/>,
194 2024.

195 Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
196 Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web
197 environment for building autonomous agents, 2024. URL <https://arxiv.org/abs/2307.13854>.

198 **A Recovery-Bench Trajectory**

Listing 1: Recovery agent response (Env + Full Message), declaring task completeness

```
199 {
200   "state_analysis": "The fasttext package has been successfully tested and is working
201   correctly. The package can be imported without errors and can train supervised
202   models as demonstrated by the successful training output showing \"Read 0M words,
203   Number of words: 11, Number of labels: 2\" and \"Model trained successfully\". The
204   initial issue appears to have been resolved through the numpy reinstallation.",
205   "explanation": "We will perform a final comprehensive test to ensure fasttext is
206   fully functional by testing both model training and prediction capabilities. This
207   will confirm that the package is working as expected for typical use cases.",
208   "commands": [
209     {
210       "keystrokes": "python3 -c \"import fasttext; model =
211       fasttext.train_supervised('training_data.txt'); result = model.predict('This is
212       awesome!'); print('Prediction result:', result); print('FastText is working
213       correctly!')\"\\n",
214       "is_blocking": true,
215       "timeout_sec": 15
216     }
217   ],
218   "is_task_complete": true
219 }
220 }
```

222 **Failure analysis.** The underlying task was to fix an *incompatible* fasttext version, but the previous
223 agent prematurely declared success. Misled by traces from the weaker model, the recovery agent
224 concluded that the package was already functional and marked the task as complete, instead of
225 investigating further. However, because the Env Replay Only settings do not have the misleading
226 trajectories, the agent is able to look at the erring state and solve the task.