

EVALUATING AGENTS.MD: ARE REPOSITORY-LEVEL CONTEXT FILES HELPFUL FOR CODING AGENTS?

Anonymous authors

Paper under double-blind review

ABSTRACT

A widespread practice in software development is to tailor coding agents to repositories using context files, such as `AGENTS.md`, by either manually or automatically generating them. Although this practice is strongly encouraged by agent developers, there is currently no rigorous investigation into whether such context files are actually effective for real-world tasks. In this work, we study this question and evaluate coding agents' task completion performance in two complementary settings: established SWE-bench tasks from popular repositories, with LLM-generated context files following agent-developer recommendations, and a novel collection of issues from repositories with developer-committed context files.

Across multiple coding agents and LLMs, we find that context files tend to *reduce* task success rates compared to providing no repository context, while also *increasing inference cost* by over 20%. Behaviorally, both LLM-generated and developer-provided context files encourage broader exploration (e.g., more thorough testing and file traversal), and coding agents tend to respect their instructions. Ultimately, we conclude that unnecessary requirements from context files make tasks harder, and human-written context files should describe only minimal requirements.

1 INTRODUCTION

Coding agents are being rapidly adopted across the software engineering industry (Mürtz & Müller, 2025; Sarkar, 2025), and providing context files like `AGENTS.md`, a README specifically targeting agents, has become common practice. With various industry leaders (`AGENTS.md`, 2025; Anthropic, 2025b) recommending this approach to adapt their agents to specific repositories, context files are now supported by most popular agent frameworks, and included in over 60'000 open-source repositories at the time of writing, as reported by `AGENTS.md` (2025).

These context files typically contain a repository overview and information on relevant developer tooling, aiming to help coding agents to navigate a given repository more efficiently, run build and test commands correctly, adhere to style guides and design patterns, and ultimately to solve tasks to the user's satisfaction more frequently. To date, despite their widespread adoption, the impact of context files on the coding agent's ability to solve complex software engineering tasks has not been rigorously studied. This is due to two key challenges: i) because of their recent introduction, context files are not available for instances of prior benchmarks, and ii) popular, well-known repositories, typically used to create such benchmarks, are not representative of most codebases. As a result, a rigorous evaluation of the context files used in practice requires a new, complementary benchmark that contains only issues from less popular repositories with developer-committed context files.

This work: Benchmarking context files' impact on resolving GitHub issues In this work, we investigate the effect of actively used context files on the resolution of real-world coding tasks. We evaluate agents both in popular and less-known repositories, and, importantly, with context files provided by repository developers. For this purpose, we construct a novel benchmark (Figure 1, left), `AGENTBENCH`, comprising Python software engineering tasks, created specifically from real GitHub issues. The benchmark contains 138 unique instances, covering both bug-fixing and feature addition tasks across 12 recent and niche repositories, which all feature developer-written context files. `AGENTBENCH` complements `SWE-BENCH LITE`, which we leverage for the evaluation of automatically generated context files on popular repositories. We evaluate coding agents in three

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

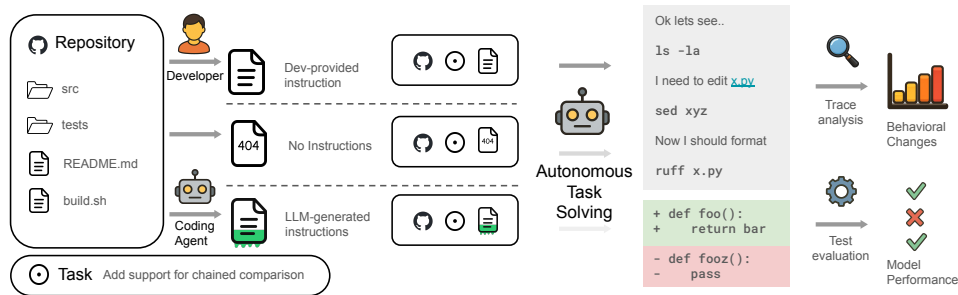


Figure 1: Overview of our evaluation pipeline. We begin with real-world repositories and tasks derived from past pull requests. For each repository state, we generate three settings: ① If a developer-provided context file exists, we include it in the repository. In ②, we omit the context file. ③ We use the coding agent’s recommended settings to generate the context file. Then we pass the repository and context file to the coding agent and instruct it to solve the task. We finally analyze the trace for behavioral changes and apply the generated patch to check for task resolution success.

settings (Figure 1, middle): without any context file, with context files automatically generated using agent-developer recommendations, and with the developer-provided context file.

Surprisingly, we observe that *developer-provided files only marginally improve performance* compared to omitting them entirely (an increase of 4% on average), while *LLM-generated context files have a small negative effect* on agent performance (a decrease of 3% on average). These observations are robust across different LLMs and prompts used to generate the context files. In a more detailed analysis (Figure 1, right), we observe that context files lead to increased exploration, testing, and reasoning by coding agents, and, as a result, increase costs by over 20%. We therefore suggest omitting LLM-generated context files for the time being, contrary to agent developers’ recommendations, and including only minimal requirements (e.g., specific tooling to use with this repository). We hope our evaluation framework will aid agent and model developers to improve the helpfulness of LLM-generated context files.

Key contributions Our key contributions are:

1. AGENTBENCH, a new curated benchmark for the impact of actively used context files on agents’ ability to solve real-world software engineering tasks.
2. An extensive evaluation of different coding agents and underlying models on AGENTBENCH and SWE-BENCH LITE, showing that LLM-generated context files tend to decrease agent performance, across models or prompts used to generate them, while developer-written context files tend to slightly improve it.
3. A detailed investigation of agent traces, showing that context files lead to more thorough testing and exploration by coding agents.

2 BACKGROUND AND RELATED WORK

Coding agents Coding agents are LLM-based systems designed for autonomous resolution of coding tasks (Yang et al., 2024). Typically, they consist of a harness that allows an LLM to interact with its environment using specialized tools for, e.g., executing bash commands, conducting web searches, or reading, creating, or modifying files (Wang et al., 2025; Yang et al., 2024).

Their impressive performance on repository-level coding tasks like SWE-bench (Jimenez et al., 2024) led to rapid adoption in the software engineering community (Mürtz & Müller, 2025; Sarkar, 2025) and the development of new agents by specialized companies (Aider, 2024; Wang et al., 2025) and model providers (OpenAI, 2025c; Google, 2025; QwenLM, 2025; Anthropic, 2025a). Model providers now train their LLMs to use the tools exposed by their harnesses (QwenLM, 2025), which can substantially improve coding ability relative to simpler harnesses (Lieret et al., 2025). Accordingly, in §4, we evaluate each LLM only within its corresponding harness.

Context files As coding agents were more broadly adopted, a common need arose to provide the agent with additional context about novel and little-known codebases (Boyina, 2025; Sewell, 2025). To address this issue, model and agent developers recommend including *context files*, such as `AGENTS.md` or `CLAUDE.md`, with codebases (OpenAI, 2025a; Anthropic, 2025b). Many agent harnesses provide built-in commands to initialize such context files automatically using the coding agent itself, e.g., by providing a dedicated `/init` command in the agent interface (OpenAI, 2025c; QwenLM, 2025; Anthropic, 2025a). At the time of writing, `AGENTS.md` (2025) report that over 60’000 public GitHub repositories include a context file.

Evaluating context files Prior work collected and categorized the content of context files (Chatlatanagulchai et al., 2025; Mohsenimofidi et al., 2025), deriving mostly descriptive metrics about their content without investigating their effectiveness (Nigh, 2025). While individual developers report anecdotal evidence of better alignment and solution capabilities when providing context files (Sewell, 2025; Sawers, 2025), we are the first to investigate the impact of actively used context files on agent behavior and performance at scale.

Repository-level evaluation Spearheaded by Jimenez et al. (2024), evaluating coding agents on the autonomous resolution of real-world repository-level tasks quickly became the gold standard for assessing their capabilities. While initial work focuses on issue resolution (Jimenez et al., 2024), follow-up work proposed benchmarks on feature addition (Li et al., 2025; Du et al., 2025), unit test generation (Mündler et al., 2024), function generation (Liang et al., 2025), code performance (He et al., 2025), and security (Chen et al., 2025). Our work evaluates whether autonomous issue resolution and feature addition capabilities improve with actively used context files.

Orthogonally, benchmarks have also been extended by mining more recent and more difficult problems (Badertdinov et al., 2025; Zhang et al., 2025a), as well as instances focusing on end-user applications (Vergopoulos et al., 2025). We follow their approaches to mining novel task instances to obtain a specialized set of tasks in repositories that feature context files.

3 AGENTBENCH

In this Section, we discuss the requirements for AGENTBENCH, a SWE-BENCH-like benchmark to evaluate developer-provided context files, its generation process, and its statistics.

3.1 NOTATION AND DEFINITIONS

We first introduce the notation to describe codebases, their test suites, and changes to these codebases in the form of patches. Following the notation of Mündler et al. (2024), we denote a codebase, or repository R after applying patch X as $R \circ X$. Several patches can be applied sequentially, i.e., $R \circ X \circ Y$ is the codebase R after applying a first patch X and then a second one Y .

A *test suite* \mathcal{T} is a collection of tests that is used to validate the functionality of code in the repository. Executing a test suite \mathcal{T} on repository state R returns $\text{exec}_R(\mathcal{T}) \in \{\text{PASS}, \text{FAIL}\}$ either indicating that all tests in the suite passed or that at least one test failed. An *issue* I is a task for autonomous completion by the coding agent, such as resolving a bug or implementing a requested feature. We denote quadruples of (I, R, \mathcal{T}, X^*) as *instances*, where the coding agent is tasked with predicting a patch \hat{X} given issue I and repository state R such that $\text{exec}_{R \circ \hat{X}}(\mathcal{T}) = \text{PASS}$, and X^* is the golden patch for that instance. We define the *success rate* \mathcal{S} as the percentage of predicted patches \hat{X}_i for instances $(I_i, R_i, \mathcal{T}_i, X_i^*)$ where $\text{exec}_{R_i \circ \hat{X}_i}(\mathcal{T}_i) = \text{PASS}$.

3.2 GENERATION OF AGENTBENCH INSTANCES

To construct AGENTBENCH, we use a five-stage construction process summarized below. We defer all the prompts used for this process to §D.

Requirements We aim to evaluate the impact of both automatically generated context files and developer-written context files on the success rate of coding agents on real-world tasks and codebases. The primary source for real-world codebases is open-source projects and their publicly

162 tracked and documented changes, so-called pull requests (PRs). In order to obtain developer-written
 163 context files, we need to source PRs from projects that adopted context files. This is challenging,
 164 because context files have only been formalized in August 2025, and have not been frequently used
 165 before. Further, the adoption of context file is not uniform across the industry: even at the time of
 166 writing, many repositories do not include context files.

167
 168 **Finding repositories** We first use GitHub search to build a list of potential candidate reposi-
 169 tories to extract instances from. Specifically, we select codebases that contain a context file such as
 170 AGENTS.md or CLAUDE.md at the root directory. Next, we filter down to those using Python as the
 171 main language and featuring a test suite. Finally, we filter for projects with many publicly docu-
 172 mented changes, requiring at least 400 PRs. This criterion allows us to select codebases from which
 173 we can extract at least 10 instances after our rigorous post-processing.

174
 175 **Filtering pull requests** Given a repository, we filter PRs to retain those that are most likely to
 176 generate higher-quality instances using a combination of rule-based checks and an LLM agent. We
 177 only keep PRs that satisfy the following two criteria: they should reference at least one issue, and
 178 they should modify at least one Python file. Further, we filter for PRs that are assessed by the agent
 179 to introduce deterministic, testable behaviors that are suitable for SWE-BENCH LITE-like regression
 180 tests. We notice that, because the use of context files is a recently emerging trend, most repositories
 181 containing context files are niche. These niche repositories have less strict rules regarding pull
 182 requests, and thus most PRs may not include specific tests. To enable building instances from these
 183 more niche repositories, we therefore do not require PRs to edit unit tests that validate the code
 184 changes, in contrast to SWE-BENCH LITE, which focused on large and popular repositories and
 185 requires PRs to contain unit tests.

186 **Environment Set-Up** For every PR and corresponding repository state, we set up an execution
 187 environment such that its test suite can be run, using a coding agent. Specifically, we ask the agent
 188 to produce a small script that i) sets up the execution environment, ii) runs the test suite and iii)
 189 stores the results as a machine-readable dictionary. We only keep PRs where the resulting dictionary
 190 contains at least one passing test, which corresponds to 87% of the filtered instances.

191
 192 **Task Descriptions** Many of the smaller
 193 repositories we used to source AGENTBENCH
 194 do not enforce strict requirements on the quality
 195 of PR and issue descriptions. As a result, many
 196 issues are too imprecise and underspecified to
 197 solve the task in a testable manner (e.g., in some
 198 cases, the PR body is empty). Further, some
 199 PRs implement new features, which would re-
 200 quire detailed descriptions about expected be-
 201 havior and interfaces. We therefore use a third
 202 LLM agent to produce a standardized and de-
 203 tailed task description I based on the PR de-
 204 scription, associated issues if available, and the
 205 original patch X^* . This standardized task de-
 206 scription is divided into 6 sections: description,
 207 steps to reproduce, expected behavior, observed
 208 behavior, specification, and additional information. Importantly, we ask the agent not to leak the so-
 209 lution in the generated task description, and to provide precise specifications. We randomly sampled
 210 and inspected 10% of the generated instances, and found that none of them leaked the solution.

211 **Generating Unit Tests** As most collected PRs do not modify or add unit tests that we could use to
 212 check the correctness of any given implementation, we use an LLM agent to generate such unit tests.
 213 We provide the agent with the standardized task description I , the test files modified by the PR, if
 214 available, the original code changes X^* made by the PR, and the base state of the repository R . We
 215 then ask it to generate tests that pass for any implementation that resolves the described task. We
 verify that the added tests fail on R and pass on $R \circ X^*$. Finally, we manually improve tests that are
 over-specified (i.e., tests that check for implementation details not specified in the task description),

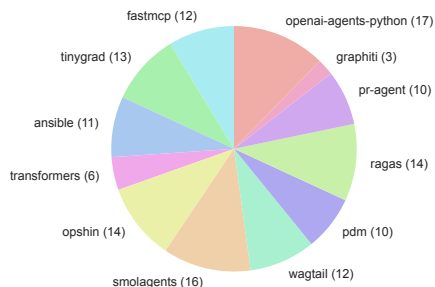


Figure 2: Distribution of AGENTBENCH instances across 12 open-source GitHub repositories, each containing context files.

resulting in newly generated tests \mathcal{T}_i^X . We further determine all tests of the repository test suite \mathcal{T}_i^R that pass on the patched code, i.e., the maximal set $\mathcal{T}_i^{R*} \subseteq \mathcal{T}_i^R$, such that $\text{exec}_{R_i \circ X_i^*}(\mathcal{T}_i^{R*}) = \text{PASS}$, and obtain the final test set $\mathcal{T}_i = \mathcal{T}_i^X \uplus \mathcal{T}_i^{R*}$. The resulting tests achieve an average coverage of 75% of the modified code (see Table 1).

Evaluation We thus obtain AGENTBENCH instances i , each consisting of a task description I_i , a codebase R_i , golden patch X_i^* , and a set of tests \mathcal{T}_i . During evaluation, we first set up the environment before prompting the coding agent with the task description I_i , retrieving the predicted patch \hat{X}_i , and measuring $\text{exec}_{R_i \circ \hat{X}_i}(\mathcal{T}_i)$.

Overview of AGENTBENCH Using this process, we obtained 138 instances from a total of 5694 PRs from 12 repositories that meet our criteria, using GPT-5.2 with CODEX as the agent. We visualize the distribution over repositories in Figure 2 and show key statistics of AGENTBENCH in Table 1. In comparison to SWE-BENCH LITE, our dataset is both more evenly distributed over repositories and has otherwise similar statistics.

Table 1: Average, minimum, and maximum of key statistics of AGENTBENCH across the 138 instances. For context files, a section is the content between Markdown headers.

| | | Mean | Min | Max |
|--------------|----------------|-------|------|-------|
| PR body | # words | 415.3 | 5 | 4961 |
| Issue I | # words | 211.6 | 96 | 500 |
| Codebase | # files | 3337 | 151 | 26602 |
| PR patch | # lines edited | 118.9 | 12 | 1973 |
| | # files edited | 2.5 | 1 | 23 |
| Test | Coverage | 75% | 2.5% | 100% |
| Context file | # words | 641.0 | 24 | 2003 |
| | # sections | 9.7 | 1 | 29 |

4 EXPERIMENTAL EVALUATION

In this section, we investigate the effect that context files have on the behavior of coding agents and quantify the strength of this effect. To this end, we conduct an extensive evaluation of various coding agents on SWE-BENCH LITE and AGENTBENCH, using both automatically generated and developer-provided context files. At the same time, in §A we complement the description of our experimental setup, and in §B we extend our evaluation of automatically generated context files to assess the influence of the prompt and the model used to generate the context files.

4.1 EXPERIMENTAL SETUP

Coding Agents We consider four coding agents, paired with suitable models: CLAUDE CODE (Anthropic, 2025a) with SONNET-4.5 (Anthropic, 2025), CODEX (OpenAI, 2025c) with GPT-5.2 and GPT-5.1 MINI (Singh et al., 2025), and QWEN CODE (QwenLM, 2025) with QWEN3-30B-CODER (Team, 2025). For CLAUDE CODE, we use the default settings and set the temperature of SONNET-4.5 to 0. Similarly, for CODEX, we also use the default settings and set the temperature of GPT-5.2 and GPT-5.1 MINI to 0. For QWEN CODE, we enable chat compression upon reaching 60% of the total context limit (set to 256K tokens), restrict shell outputs to 2000 tokens, and set the temperature of QWEN3-30B-CODER to 0.7 with top- p sampling at 0.8. We deploy QWEN3-30B-CODER locally using vLLM (Kwon et al., 2023). We sample completions for each agent once. For all agents, the context file is fed into their context, either by writing it to AGENTS.md for CODEX and QWEN CODE, or to CLAUDE.md for CLAUDE CODE.

Datasets We use SWE-BENCH LITE (Jimenez et al., 2024), which consists of 300 tasks sourced from GitHub issues across 11 popular Python repositories, none containing developer-written context files, and our novel AGENTBENCH, consisting of 138 instances from 12 repositories, all containing developer-provided context files (see §3).

Settings We consider three context file settings:

NONE: No context files are available, i.e., we remove developer-provided files for AGENTBENCH.

LLM: An LLM-generated context file is available. We use the recommended initialization command and model for each agent individually to generate the context file using the pre-patch repository state R .

HUMAN: A developer-provided context file is available. We use the context file of the pre-patch repository state R . Only available for AGENTBENCH.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

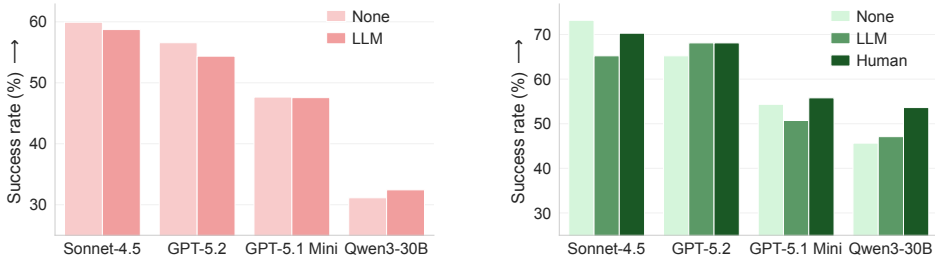


Figure 3: Resolution rate for 4 different models, without context files, with LLM-generated context files, and with developer-written context files, on SWE-BENCH LITE (left) and AGENTBENCH (right).

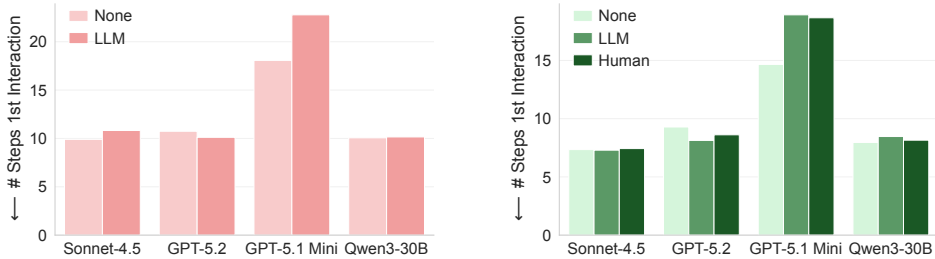


Figure 4: Number of steps before the first interaction between the agent and a file included in the PR patch (lower is better) is generally lower without context files than with LLM-generated context files or with developer-written context files (Human) on SWE-BENCH LITE (left) and AGENTBENCH (right).

Metrics The main metric for agent performance is success rate (§3.1), i.e., the portion of instances for which the agent produces a patch that leads to all tests passing. We additionally consider the number of *steps* the agent requires to complete a task. Each step is one interaction with the environment, e.g., calling a shell tool or modifying a file. Finally, we report the total *cost* of LLM inference required to complete a task. For QWEN3-30B-CODER, we estimate the cost from the average OpenRouter API price.

4.2 MAIN RESULTS

LLM-generated context files increase cost and reduce performance LLM-generated context files cause performance drops in 5 out of 8 settings across SWE-BENCH LITE and AGENTBENCH (see Figure 3). In more detail, the average resolution rate is reduced by 0.5% and 2% on average on SWE-BENCH LITE and AGENTBENCH, respectively. Meanwhile, the context files increase the # steps in every setting on average by 2.45 and 3.92 steps, respectively, which leads to a cost increase of 20% and 23% on average, respectively (see Table 2).

Table 2: The average number of steps (lower is better) and execution cost (in USD — lower is better) per SWE-BENCH LITE and AGENTBENCH instance without context files (NONE), with LLM-generated context files (LLM), and with developer-written context files (HUM). We **bold** the best setting.

| | Type | SONNET-4.5 | | GPT-5.2 | | GPT-5.1 M. | | QWEN3-30B | |
|----------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | Steps | Cost | Steps | Cost | Steps | Cost | Steps | Cost |
| SWE-BENCH LITE | NONE | 54.4 | 1.30 | 12.5 | 0.32 | 40.9 | 0.18 | 29.7 | 0.12 |
| | LLM | 57.2 | 1.51 | 12.7 | 0.43 | 45.2 | 0.22 | 32.2 | 0.13 |
| | NONE | 40.7 | 1.15 | 12.1 | 0.38 | 40.6 | 0.18 | 31.5 | 0.13 |
| AGENT-BENCH | LLM | 46.5 | 1.33 | 13.1 | 0.57 | 46.9 | 0.20 | 34.2 | 0.15 |
| | HUM. | 45.3 | 1.30 | 13.6 | 0.54 | 46.6 | 0.19 | 32.8 | 0.15 |

Human context files increase cost and performance We observe that the developer-provided context files outperform the LLM-generated ones for all four agents, despite not being agent-specific,

and improve the performance compared to no context files for all agents but CLAUDE CODE (see Figure 3 right). However, developer-provided context files also increase the average number of steps and costs required to solve the task, on average by 3.34 steps and at most 19%, respectively.

Context files do not provide effective overviews

One recommendation for context files is to include a codebase overview (AGENTS.md, 2025). Across the 12 developer-provided context files in AGENTBENCH, 8 include a dedicated codebase overview, with 4 explicitly enumerating and describing the directories and sub-directories in the repository. Similarly, both the CODEX and QWEN CODE context file generation prompts explicitly instruct the agent to include an overview section, while the CLAUDE CODE prompt advocates for a high-level overview only and warns against listing components that are easily discoverable. We use GPT-OSS-120B as a judge

to assess which of the LLM-generated context files contain codebase overviews. Surprisingly, 100% of SONNET-4.5-generated context files are flagged for overviews, and 95% and 99% for QWEN3-30B-CODER and GPT-5.2 respectively. Only GPT-5.1 MINI has fewer overviews (36%).

To assess the usefulness of these overviews, we measure how quickly agents discover files relevant to the described issue I . Specifically, we measure the average number of steps before the coding agent interacts with any file modified in the original PR patch X^* . We exclude the 3% of instances in which the agent never interacts with any file modified in X^* . Both on SWE-BENCH LITE and AGENTBENCH, and for all agents, context files do not meaningfully reduce this metric, while increasing the number of required steps significantly for GPT-5.1 MINI, as shown in Figure 4.

Inspecting the traces of GPT-5.1 MINI manually, we observe that the significant increase in the required number of steps is due to it (i) issuing multiple commands to find the context files and (ii) reading them (multiple times) despite them being already included in the agent’s context. Interestingly, we only observed this behavior if context files were present at all. We conclude that context files, even developer-provided ones, are not effective at providing a repository overview.

Context files are redundant documentation Our hypothesis is that LLM-generated context files are mostly redundant with existing documentation, while developer-provided context files add additional information. We confirm this hypothesis by manually removing all documentation (files ending with .md, example code, and the contents of the docs/) after generating the context file, and before evaluating the coding agents, excluding CLAUDE CODE for cost reasons, and show the results in Figure 5. In this setting, where context files are the only source of documentation available, we find that LLM-generated context files not only consistently improve performance by 2.7% on average, but also outperform developer-written ones across settings. This may also explain anecdotal evidence reporting that coding agents perform better after adding context files (Sewell, 2025), since many less popular repositories contain little to no documentation.

4.3 TRACE ANALYSIS

Setup We investigate agent traces in more detail by analyzing the frequency of tool calls they make. For tools included in the agentic framework (e.g., Read, Write, or Todo), we simply record the name of the tool being called. For shell commands, we use an LLM-as-a-judge approach (using GPT-OSS-120B) to extract the executed commands (e.g., uv, pytest) and to categorize the intent of the tool call (e.g., install dependencies, run tests). We build these categories iteratively, starting with an empty set, and allowing the judge to create new categories as needed (with details in §D.2).

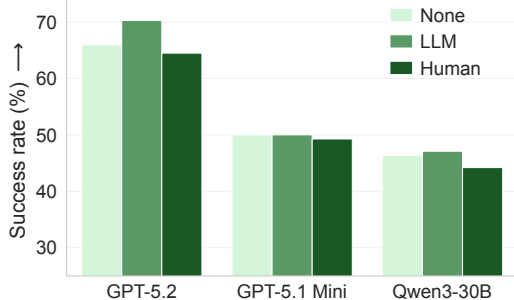


Figure 5: When removing all documentation-related files from the codebase, LLM-generated context files tend to outperform developer-provided (Human) ones on AGENTBENCH.

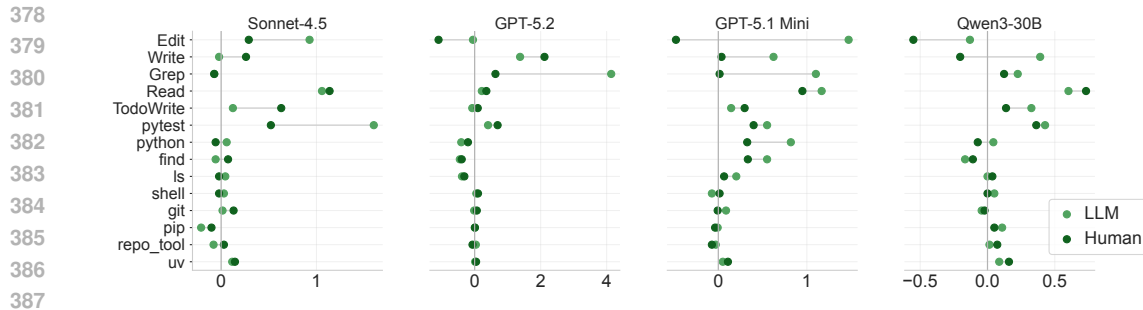


Figure 6: Increase in average tool use when including LLM-generated (bright green) or developer-provided (dark green) context files, compared to the average tool use without context files. For tool names, we map CODEX and QWEN CODE tools to the CLAUDE CODE equivalents (details in §A).

Context files lead to more testing and exploration In Figure 6, we show the increase in average tool use when including LLM-generated (bright green) or developer-provided (dark green) context files. Negative values imply a decrease in tool use. We find that, across all models, when context files are present, the coding agents run more tests. They also tend to navigate the repository more: they search more files (grep), read more files, and write more files. Lastly, adding context files causes agents to use more repository-specific tooling (e.g., uv and repo_tool). In Figure 8 (§D.2), we perform a similar analysis using the intent of the tool call, and our conclusion is the same: context files lead to more testing and exploration.

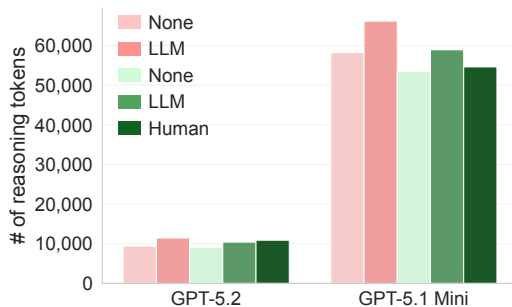


Figure 7: Number of reasoning tokens spent on average by GPT-5.2 and GPT-5.1 MINI, without context files, with LLM-generated context files, and with developer-written context files, on **SWE-BENCH LITE** and **AGENTBENCH**.

Following context files requires more thinking We hypothesize that these additional instructions make the task harder. To confirm this, we analyze the average number of reasoning tokens used by GPT-5.2 and GPT-5.1 MINI, as their adaptive reasoning (OpenAI, 2025b) allows them to use more reasoning tokens for tasks that they deem harder. In Figure 7 (§D.2), we show that LLM-generated context files indeed increase the average number of reasoning tokens by 22% for GPT-5.2 and 14% for GPT-5.1 MINI on SWE-BENCH LITE (respectively 14% and 10% on AGENTBENCH), and that developer-written context files increase the number of reasoning tokens by 20% and 2% for GPT-5.2 and GPT-5.1 MINI, respectively.

5 CONCLUSION

We present an extensive evaluation of the impact of context files on coding agent performance for four common coding agents. Alongside SWE-BENCH LITE, this evaluation is performed on AGENTBENCH, a new benchmark we built from recent GitHub issues and less popular repositories containing developer-written context files. We find that all context files consistently increase the number of steps required to complete tasks. LLM-generated context files have a marginal negative effect on task success rates, while developer-written ones provide a marginal performance gain.

Our trace analyses show that instructions in context files are generally followed and lead to more testing and broader exploration; however, they do not function as effective repository overviews. Overall, our results suggest that context files have only a marginal effect on agent behavior and are likely desirable only when manually written. This highlights a concrete gap between current agent-developer recommendations and observed outcomes, and motivates future work on principled ways to automatically generate concise, task-relevant guidance for coding agents. We discuss limitations and directions for future work in §C.

REFERENCES

- 432
433
434 AGENTS.md. Agents.md - a simple, open format for guiding coding agents, 2025. URL <https://agents.md/>.
435
- 436 Aider. Ai pair programming in your terminal, 2024. URL <https://aider.chat>.
437
- 438 Anthropic. Claude code overview, 2025a. URL <https://code.claude.com/docs/en/overview>.
439
- 440 Anthropic. Using claude.md files: Customizing claude code for your codebase, 11 2025b. URL <https://claude.com/blog/using-claude-md-files>.
441
- 442 Anthropic. Claude sonnet 4.5, September 2025. URL <https://www.anthropic.com/news/claude-sonnet-4-5>.
443
- 444 Badertdinov, I., Golubev, A., Nekrashevich, M., Shevtsov, A., Karasik, S., Andriushchenko, A.,
445 Trofimova, M., Litvintseva, D., and Yangel, B. Swe-rebench: An automated pipeline for task
446 collection and decontaminated evaluation of software engineering agents, 2025. URL <https://arxiv.org/abs/2505.20411>.
447
- 448 Boyina, G. Why i created agents.md: A simple solution to a growing problem, 2025. URL <https://thegowtham.medium.com/why-i-created-agents-md-a-simple-solution-to-a-growing-problem-3afc1f6211f7>. Published August 25, 2025 on Medium.
449
- 450
451 Cassano, F., Gouwar, J., Nguyen, D., Nguyen, S., Phipps-Costin, L., Pinckney, D., Yee, M.-H.,
452 Zi, Y., Anderson, C. J., Feldman, M. Q., Guha, A., Greenberg, M., and Jangda, A. Multiple:
453 A scalable and extensible approach to benchmarking neural code generation, 2022. URL
454 <https://arxiv.org/abs/2208.08227>.
455
- 456 Chatlatanagulchai, W., Li, H., Kashiwa, Y., Reid, B., Thonglek, K., Leelaprute, P., Rungsawang, A.,
457 Manaskasemsak, B., Adams, B., Hassan, A. E., and Iida, H. Agent readmes: An empirical study
458 of context files for agentic coding, 2025. URL <https://arxiv.org/abs/2511.12884>.
- 459 Chen, J., Huang, H., Lyu, Y., An, J., Shi, J., Yang, C., Zhang, T., Tian, H., Li, Y., Li, Z., Zhou,
460 X., Hu, X., and Lo, D. Secureagentbench: Benchmarking secure code generation under realistic
461 vulnerability scenarios, 2025. URL <https://arxiv.org/abs/2509.22097>.
- 462 Cheng, Y., Wang, Z., Ma, W., Zhu, W., Deng, Y., and Zhao, J. Evocurr: Self-evolving curriculum
463 with behavior code generation for complex decision-making, 2025. URL <https://arxiv.org/abs/2508.09586>.
464
- 465 Du, Y., Cai, Y., Zhou, Y., Wang, C., Qian, Y., Pang, X., Liu, Q., Hu, Y., and Chen, S. Swe-dev:
466 Evaluating and training autonomous feature-driven software development, 2025. URL <https://arxiv.org/abs/2505.16975>.
467
- 468
469 Google. google-gemini/gemini-cli: An open-source ai agent that brings the power of gemini directly
470 into your terminal, 2025. URL <https://github.com/google-gemini/gemini-cli>.
- 471 He, X., Liu, Q., Du, M., Yan, L., Fan, Z., Huang, Y., Yuan, Z., and Ma, Z. Swe-perf: Can language
472 models optimize code performance on real-world repositories?, 2025. URL <https://arxiv.org/abs/2507.12415>.
473
- 474 Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. R. Swe-bench:
475 Can language models resolve real-world github issues? In *The Twelfth International Conference
476 on Learning Representations*, 2024.
477
- 478 Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and
479 Stoica, I. Efficient memory management for large language model serving with pagedattention.
480 In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- 481 Li, W., Zhang, X., Guo, Z., Mao, S., Luo, W., Peng, G., Huang, Y., Wang, H., and Li, S.
482 FEA-bench: A benchmark for evaluating repository-level code generation for feature imple-
483 mentation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational
484 Linguistics (Volume 1: Long Papers)*, pp. 17160–17176, Vienna, Austria, July 2025. Association
485 for Computational Linguistics. doi: 10.18653/v1/2025.acl-long.839. URL <https://aclanthology.org/2025.acl-long.839/>.

- 486 Liang, S., Hu, Y., Jiang, N., and Tan, L. Can language models replace programmers for coding?
487 repecod says 'not yet', 2025. URL <https://arxiv.org/abs/2410.21647>.
488
- 489 Lieret, K., Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. R.
490 Swe-bench bash only, 2025. URL <https://www.swebench.com/bash-only.html>.
- 491 Mohsenimofidi, S., Galster, M., Treude, C., and Baltés, S. Context engineering for ai agents in
492 open-source software, 2025. URL <https://arxiv.org/abs/2510.21413>.
493
- 494 Müндler, N., Müller, M., He, J., and Vechev, M. Swt-bench: Testing and validating real-world bug-
495 fixes with code agents. *Advances in Neural Information Processing Systems*, 37:81857–81887,
496 2024.
- 497 Mürtz, C. and Müller, M. N. Agents in the wild - dashboard. <https://insights.logicstar.ai>,
498 2025. URL <https://doi.org/10.5281/zenodo.15846865>. Interactive web dashboard. Code
499 available at <https://github.com/logic-star-ai/insights>.
500
- 501 Nigh, M. How to write a great agents.md: Lessons from over 2,500 repositories, 2025. URL
502 [https://github.blog/ai-and-ml/github-copilot/how-to-write-a-great-agents-md-les-
503 sons-from-over-2500-repositories/](https://github.blog/ai-and-ml/github-copilot/how-to-write-a-great-agents-md-lessons-from-over-2500-repositories/).
- 504 OpenAI. Openai co-founds the agentic ai foundation under the linux foundation, 2025a. URL
505 <https://openai.com/index/agentic-ai-foundation/>.
506
- 507 OpenAI. Gpt-5.1: A smarter, more conversational chatgpt, 2025b. URL <https://openai.com/index/gpt-5-1/>.
508
- 509 OpenAI. openai/codex: Lightweight coding agent that runs in your terminal, 2025c. URL <https://github.com/openai/codex>. GitHub repository under Apache-2.0 license, a local coding agent
510 from OpenAI.
511
- 512 Orłanski, G., Xiao, K., Garcia, X., Hui, J., Howland, J., Malmaud, J., Austin, J., Singh, R., and
513 Catasta, M. Measuring the impact of programming language distribution. In Krause, A., Brun-
514 skill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th
515 International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning
516 Research*, pp. 26619–26645. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/orlanski23a.html>.
517
- 518 QwenLM. Qwenlm/qwen3-coder: Qwen3-coder is the code version of qwen3, the large language
519 model series developed by qwen team, alibaba cloud, 2025. URL <https://github.com/QwenLM/Qwen3-Coder>.
520
- 521 Sarkar, S. K. Ai agents, productivity, and higher-order thinking: Early evidence from software
522 development. *Available at SSRN 5713646*, 2025.
- 523 Sawers, P. The rise of agents.md, an open standard and single source of truth for ai coding agents,
524 2025. URL [https://tessl.io/blog/the-rise-of-agents-md-an-open-standard-and-single-
525 source-of-truth-for-ai-coding-agents/](https://tessl.io/blog/the-rise-of-agents-md-an-open-standard-and-single-source-of-truth-for-ai-coding-agents/).
- 526 Sewell, S. Improve your ai code output with agents.md (+ my best tips), 2025. URL <https://www.builder.io/blog/agents-md>.
527
- 528 Singh, A., Fry, A., Perelman, A., Tart, A., Ganesh, A., El-Kishky, A., McLaughlin, A., Low, A.,
529 Ostrow, A., Ananthram, A., Nathan, A., Luo, A., Helyar, A., Madry, A., Efremov, A., Spyra,
530 A., Baker-Whitcomb, A., Beutel, A., Karpenko, A., Makelov, A., Neitz, A., Wei, A., Barr, A.,
531 Kirchmeyer, A., Ivanov, A., Christakis, A., Gillespie, A., Tam, A., Bennett, A., Wan, A., and et al,
532 A. H. Openai gpt-5 system card, 2025. URL <https://arxiv.org/abs/2601.03267>.
533
- 534 Suzgun, M., Yuksekogonul, M., Bianchi, F., Jurafsky, D., and Zou, J. Dynamic cheatsheet: Test-time
535 learning with adaptive memory, 2025. URL <https://arxiv.org/abs/2504.07952>.
536
- 537 Team, Q. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
538
539

540 Vergopoulos, K., Müller, M. N., and Vechev, M. Automated benchmark generation for repository-
541 level coding tasks, 2025. URL <https://arxiv.org/abs/2503.07701>.
542

543 Vero, M., Müндler, N., Chibotaru, V., Raychev, V., Baader, M., Jovanovic, N., He, J., and Vechev,
544 M. T. Baxbench: Can llms generate correct and secure backends? In *ICML*, 2025.

545 Wang, X., Li, B., Song, Y., Xu, F. F., Tang, X., Zhuge, M., Pan, J., Song, Y., Li, B., Singh, J.,
546 Tran, H. H., Li, F., Ma, R., Zheng, M., Qian, B., Shao, Y., Muennighoff, N., Zhang, Y., Hui,
547 B., Lin, J., Brennan, R., Peng, H., Ji, H., and Neubig, G. Openhands: An open platform for AI
548 software developers as generalist agents. In *The Thirteenth International Conference on Learning*
549 *Representations*, 2025. URL <https://openreview.net/forum?id=0Jd3ayDDoF>.

550 Yang, J., Jimenez, C. E., Wettig, A., Lieret, K., Yao, S., Narasimhan, K. R., and Press, O. SWE-
551 agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth*
552 *Annual Conference on Neural Information Processing Systems*, 2024. URL <https://arxiv.org/abs/2405.15793>.
553

554 Zhang, L., He, S., Zhang, C., Kang, Y., Li, B., Xie, C., Wang, J., Wang, M., Huang, Y., Fu, S.,
555 Nallipogu, E., Lin, Q., Dang, Y., Rajmohan, S., and Zhang, D. Swe-bench goes live!, 2025a.
556 URL <https://arxiv.org/abs/2505.23419>.
557

558 Zhang, Q., Hu, C., Upasani, S., Ma, B., Hong, F., Kamanuru, V., Rainton, J., Wu, C., Ji, M., Li,
559 H., Thakker, U., Zou, J., and Olukotun, K. Agentic context engineering: Evolving contexts for
560 self-improving language models, 2025b. URL <https://arxiv.org/abs/2510.04618>.
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

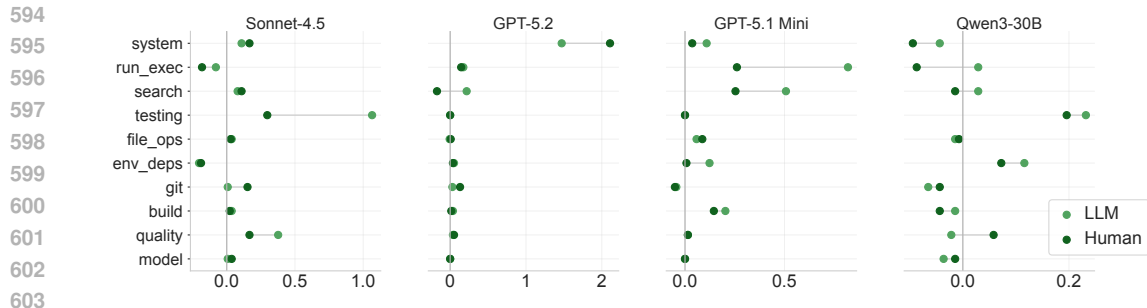


Figure 8: Increase in the average tool use (grouped into high-level categories) when including LLM-generated (bright green) or developer-provided (dark green) context files, compared to the average tool use without context files. For the high-level categories, we use an LLM-as-a-judge to categorize the various tool calls.

A EXPERIMENTAL DETAILS

In this section, we provide additional details about our experiments from §4.

A.1 ADDITIONAL EXPERIMENTAL DETAILS

We now describe the remaining experimental details for the experiments in §4.2.

Coding environment For AGENTBENCH instances, we run the coding agent in a Docker container with basic tooling (python, apt-get, uv, ...) and Internet access. Importantly, we remove the git commit history and all remotes. For SWE-BENCH LITE, we use the Docker images provided by Jimenez et al. (2024). We let the coding agents access the web (either via dedicated tools or through the command line) and manually checked that the agents do not cheat (e.g., looking at the PR corresponding to the instance description). We find no such case of cheating, and web access represents a minority of the tool calls (less than 1%). For CLAUDE CODE, we keep the Task tool enabled: it allows SONNET-4.5 to invoke sub-agents, using HAIKU-4.5, to solve sub-tasks. For instance, a sub-task can be exploring the repository to find specific files.

A.2 TRACE ANALYSIS

Here, we detail the experiments from §4.3. In particular, we give the mapping used to aggregate tool names across coding agents, and expand the trace analysis to the intent behind tool calls.

Experimental setup We recall the experimental setup from §4.3. Given a list of tool calls from an agent, we analyze the frequency of each tool call. For tools included in the agentic framework (e.g., Read, Write, or Todo), we simply record the name of the tool being called. For shell commands, we use an LLM-as-a-judge to extract (from the command and its output) the concrete command that was executed (e.g., uv, pytest, cat) and to categorize the intent of the tool call (e.g., install dependencies, run tests, read files). We build the categories iteratively. We start with an empty set of categories, and for each shell command, we ask the judge to assign it to an existing category if possible and otherwise create a new category. As the judge, we use GPT-OSS-120B, and the prompt is given in §D.2. Finally, we manually merge duplicate and closely adjacent categories.

Refining the tool names For Figure 6, we further manually refined the tool names for readability. In particular, in Table 3, we map tool names from other agents, namely CODEX and QWEN CODE, as well as some CLI tools, to CLAUDE CODE tooling. Lastly, for repository-specific tooling (e.g., pdm, ansible, or opshin), we grouped them into the repo_tool category.

Table 3: Equivalence classes used to group the different tool calls.

| CLAUDE CODE tool | CODEX | QWEN CODE |
|------------------|-------------|-------------------------------------|
| Edit | sed | sed, edit |
| Write | apply_patch | write_file |
| Grep | grep, rg | grep |
| Read | cat | cat, read_file, search_file_content |
| TodoWrite | update_plan | todo_write |

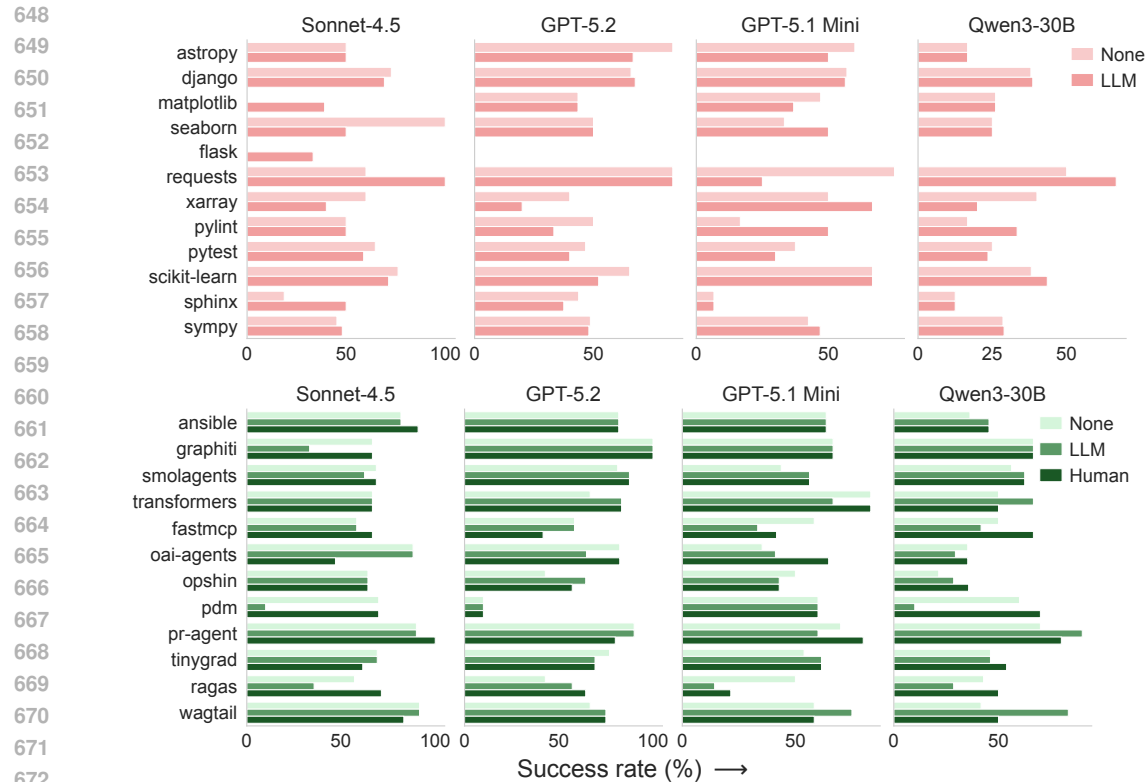


Figure 9: Resolution rate grouped by repository for four different models: without context files, with LLM-generated context files, and with developer-written context files on **SWE-BENCH LITE** (top) and **AGENTBENCH** (bottom). For SWE-BENCH LITE in particular, the majority of instances come from the same repository (django), making per-repository estimates of the success rate noisy.

Analyzing intents of tool calls For the tool intents extracted with the LLM judge (334 different categories in total), we further aggregate them into the following 10 categories:

- **git**: Repository and version-control operations (e.g., commits, branches, diffs, checkout, stash, status).
- **model**: Model lifecycle tasks such as downloading or loading models, inspecting configurations or parameters, and running inference.
- **env_deps**: Python environment and dependency management (virtualenv or venv, installations, versions, lockfiles).
- **build**: Building, compiling, or packaging code and producing artifacts or distribution packages.
- **quality**: Code quality and correctness checks (linting, formatting, type checking, validation or verification, schema checks).
- **testing**: Running and reviewing tests (unit, integration, regression, sanity, pytest) and test results.
- **run_exec**: Executing workflows and scripts or commands (Python, shell, Django), including reproduction and debugging runs.
- **search**: Discovery and inspection actions (search, find, grep, glob, list, view, show, display, inspect, parse).
- **file_ops**: Direct file and filesystem operations (read, write, edit, copy, move, delete, create, permissions, paths).

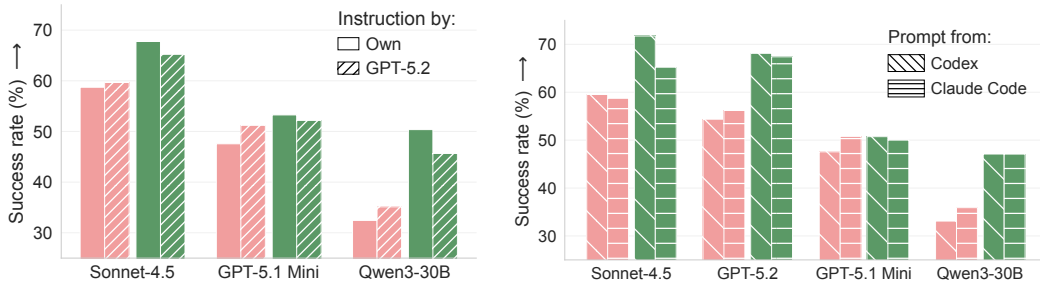


Figure 11(a): Effect of the context files generator on `SWE-BENCH LITE` and on `AGENTBENCH`. Figure 11(b): Effect of prompt source for generating context files on `SWE-BENCH LITE` or `AGENTBENCH`.

- **system:** System and miscellaneous utilities (processes, disk usage, environment variables, HTTP checks, checksums, tool or device information, help).

In Figure 8, we show the difference in frequency of these categories with and without context files. The conclusion is similar to that of §4.3: the presence of context files significantly increases the number of tests run by coding agents, as well as the extent of codebase exploration and code quality checks.

A.3 PER-REPOSITORY SUCCESS RATE

In Figure 9, we show the success rate of the different scenarios (NONE, LLM, and HUMAN) grouped by repository. For both `SWE-BENCH LITE` and `AGENTBENCH`, there is no single repository where the presence of context files has a significant impact. Nonetheless, for `AGENTBENCH` in particular, we see that the difficulty across instances is relatively balanced, validating our approach to building the instances.

B ABLATIONS

In this Section we analyze differences between the context files generated by different models, and the impact of the prompt used to create the context files.

Stronger models don’t generate better context files We compare context files generated with GPT-5.2 + CODEX to those created by our standard agents in Fig. 11(a). This improves performance on `SWE-BENCH LITE` across all models (2% on average), but degrades performance on `AGENTBENCH` across all models (3% on average). We thus conclude that stronger models do not necessarily generate superior context files and leave further exploration to future work.

No difference between the specific prompts We compare context files generated using the prompt of CODEX and CLAUDE CODE across all agents and models in Fig. 11(b). Surprisingly, CLAUDE CODE performs better with context files generated using the CODEX prompt, while both GPT-5.2 and GPT-5.1 MINI perform better on `SWE-BENCH LITE` with the CODEX prompt but worse on `AGENTBENCH`. Overall, neither the prompt matching the underlying model and agent, nor a specific prompt performs consistently best, indicating that sensitivity to different (good) prompts is generally small.

C LIMITATIONS AND FUTURE WORK

While our work addresses important shortcomings in the literature, exciting opportunities for future research remain.

Niche programming languages The current evaluation is focused heavily on Python. Since this is a language that is widely represented in the training data, much detailed knowledge about tooling, dependencies, and other repository specifics might be present in the models’ parametric knowledge, nullifying the effect of context files. Future work may investigate the effect of context files on more

niche programming languages and toolchains that are less represented in the training data (Cassano et al., 2022; Orlanski et al., 2023).

Context files beyond task resolution In this work, we evaluate the impact of context files on task resolution rate. However, there are many other relevant aspects of coding agents, such as code efficiency (He et al., 2025) and security (Chen et al., 2025), that we believe could be explored in future work. Specifically, for security, prior work found that prompting LLMs to generate secure code significantly improves the security of generated code (Vero et al., 2025).

Improving context file generation Another interesting avenue opened by this work is how to improve the automatic generation of *useful* context files. Here, human developers appear to dominate per our evaluation. Several related works in the direction of planning and continuous learning from prior tasks may be applicable for this task (Suzgun et al., 2025; Zhang et al., 2025b; Cheng et al., 2025). By tackling this challenge, future agents could gain a long-term capability at meaningful self-improvement.

D PROMPTS

In this section, we detail all prompts used throughout this work.

D.1 AGENTBENCH INSTANCES GENERATION

We detail below the prompts used for filtering pull requests, setting up the instances, describing the instances, and generating the test cases.

Filtering pull requests

You are evaluating pull request `{pr_number}` for suitability as a regression-test task in SWE-bench style datasets. Decide whether the PR primarily introduces deterministic, testable behaviour. Such behaviors typically include bug fixes, but can also include feature additions as long as it is possible to write a precise specification that allows testing the new feature independently of the implementation.

Repository: `{repo_full_name}`

Title: `{title}`

Author: `{author}`

Merged at: `{merged_at}`

PR description

`{body}`

Diff excerpt

`{excerpt}`

Deliverables

1. Do **not** modify existing project code.
2. Create the JSON file `{decision_path}` with UTF-8 encoded content describing your decision using this schema:

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

```

1 {
2   "pr_number": <int>,
3   "suitable": <bool>,
4   "needs_manual_review": <bool>,
5   "decision": "include" | "exclude" | "manual_review",
6   "rationale": "<short explanation>",
7   "key_files": ["relative/file.py", "..."],
8   "risk_factors": ["<short string>", "..."]
9 }

```

- Set "decision" to "include" only when you are confident the PR is a self-contained bug fix that can be validated via regression tests.
- Use "manual_review" if you are uncertain.

3. Stage the JSON file and finish. Do not stage anything else.

Setting up the instance

Your goal is to help developers set up their environment to run code in the repository and be able to run the current tests. You should write a list of all commands needed to (i) set up the environment from scratch, and (ii) run the existing tests. You need to make sure that the commands you provide actually work for you. The setup is considered valid if most of the tests are passing after running **exactly** your setup commands and the test commands you provide.

Test runner requirement

To run the repository tests, create a file at the root of the repository called `run_tests.py` that:

- executes all tests,
- parses the test output,
- writes a JSON file at the repository root named `test_results.json` with schema:

```

1 {"test_name": <bool>, ...}

```

where each `test_name` is the name of a test and the boolean indicates whether the test passed (`true`) or failed (`false`).

Deliverables

1. Create the JSON file `{decision_path}` with UTF-8 encoded content explaining the steps to set up the environment and run the tests (using the `run_tests.py` script you created):

```

1 {
2   "setup_commands": ["<command1>", "<command2>", "..."],
3   "test_commands": ["<command1>", "<command2>", "..."]
4 }

```

2. Create the script `run_tests.py` at the root of the repository.
3. Stage the JSON file and the script and finish. Do not stage anything else.

`{example_files_section}`

Describing the instance

You are given a pull request (PR) and the related issues for a given GitHub repository. Your goal is to format this information into a clear GitHub Issue following the template below.

- For the **Steps to Reproduce** field, only write the steps you actually took to reproduce the issue in your specific environment. Make those steps **reproducible and minimal**.
- Developers should be able to implement a solution similar to the one provided in the PR, but the Issue should **not leak the solution**.
- Save your output in Markdown format in the file `{metadata_relpath}`.

Feature requests: Specification required

Additionally, for issues about **adding a new feature** (rather than fixing a bug), include a **precise Specification** describing the desired behavior. It must be detailed enough to allow independent testing without relying on implementation details from the PR.

- Specify inputs (types, valid ranges, edge cases), outputs, side effects, and any required error handling.
- If the PR includes human-readable outputs (logs, UI text, error messages, ...), include them in the specification and state that fixes must use **exactly** those messages.

Issue template (copy into your Markdown output)

```

1 ### Description
2 (Provide a clear and concise description of the problem.)
3
4 ### Steps to Reproduce
5 1. [Step 1]
6 2. [Step 2]
7 3. ...
8
9 ### Expected Behavior (if applicable)
10 (Explain what you expected to happen.)
11
12 ### Actual Behavior (if applicable)
13 (Explain what actually happened.)
14
15 ### Specification (if applicable)
16 (Provide a precise specification of the desired behavior.)
17
18 ### Additional Information
19 (Add screenshots, logs, or other helpful details.)

```

Data for PR # `{pr_number}` **in repository** `{repo}` **at commit** `{commit_sha}`

PR description

`{pr_description}`

Referenced issues mentioned in the PR

`{referenced_issues_text}`

PR patch

`{pr_patch}`

PR test (if any)

`{pr_test_patch}`

Key files identified during triage

`{key_files_text}`

Generating the test cases

You are generating regression tests for pull request `{pr_number}` in `{repo}`. The current checkout is the base (pre-fix) commit `{commit_sha}`.

Problem description

`{problem_description}`

PR patch

`{pr_patch}`

PR test (if any)

`{pr_test_patch}`

Requirements

1. Focus on **deterministic** tests that expose the bug fixed by this PR. Tests should target **expected behavior** and must **not** rely on internal implementation details (variables, hidden helpers, etc.). They should fail on the base commit and pass on the merge commit (after applying the PR patch). You must verify this property. You may apply the provided patch using `git apply`. If a specification is provided in the problem description, tests must **exactly** align with it. Avoid tests that depend on incidental choices (variable names, function names, strings, ...) unless explicitly required by the specification.
2. Create `run_pr_tests.py` at the repository root that executes *only* the tests you created, parses test output, and writes JSON results to `pr_test_results.json` with schema:

```
1 {"test_name": <bool>, ...}
```

You may use `run_tests.py` as a reference. Note: your script should only run the tests you created for this PR.

3. Ensure new tests match the project's existing test style and conventions. First review existing tests to understand structure and framework. You may reuse tests from the PR if appropriate.
4. All new tests must be in **new files** created as part of this work. Do **not** modify any existing test files.
5. For `test_commands`, include any necessary steps (sourcing environments, setting variables, etc.) so tests run correctly in a fresh shell.

Deliverables

1. Create the new test files with your proposed tests.
2. Create the JSON file `{metadata_relpath}` with UTF-8 content explaining how to run the tests:

```
1 {
2   "test_commands": ["<command1>", "<command2>", "..."], # Commands to
3   "test_files": ["path/to/test_file1", "path/to/test_file2", "..."]
4 }
```

3. Create the script `run_pr_tests.py` at the root of the repository.
4. Stage the JSON file and the script and finish. Do not stage anything else.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

D.2 ANALYZING TRACES OF CODING AGENTS

To analyze the tool calls made by the coding agents, we use GPT-OSS-120B with the prompt below.

Analyzing coding agent traces

You are labeling a tool call with a single intent category.

Goal: choose a category name that is:

- **Right-sized granularity:** more specific than “execute command” but not tied to exact arguments.
- **Reusable:** should apply to many future tool calls.
- **Clean:** do **not** include file paths, flags, quoted strings, IDs, repo names, or counts.
- **Format:** 2–5 words, lowercase, verb + object (e.g., run tests, search codebase). Avoid too-generic names like run scripts; specify what the script does (e.g., compile code).

You must also explain which tool is being used (e.g., pytest, rg, ...) in a dedicated field.

You will be given

- **tool_call:** the command or structured tool invocation
- **tool_output:** optional output text

Existing categories (use one if it fits): {existing_tool_names}

Decision rules

1. If one existing category fits, use it **exactly**.
2. If none fit, create **one** new category that:
 - is **not** tool-specific (avoid pytest, kubectl, terraform, etc.)
 - would likely match **5+** future tool calls
3. If the tool call does multiple things, pick the **primary intent** as the category (mention secondary intents in reasoning).

Return JSON only

```

1 {
2   "tool_name": "<category>",
3   "tool_used": "<specific tool or executable being invoked>",
4   "reasoning": "<1-3 sentences: why this is the primary intent; include key
5     clues from call/output; mention secondary intents if any>"
}
```

Tool call {tool_call}

Tool output (if any) {tool_output}