

NoCode-bench: Evaluating Natural Language-Driven Software Feature Addition

Anonymous ACL submission

Abstract

LLMs have demonstrated remarkable capabilities in supporting software developers, e.g., by automating code generation and code editing. In contrast, their effectiveness and limitations in enabling software users to incrementally improve a piece of software are currently underexplored. A promising paradigm toward this end is *natural language-driven feature addition*, which allows users to specify and modify software functionality purely through natural language (NL) descriptions, sometimes also called “no-code development”. This paper introduces NoCode-bench, a benchmark designed to evaluate LLMs on real-world NL-driven software feature addition tasks. NoCode-bench consists of 634 tasks across 10 popular projects, each of which pairs a user-oriented documentation change and the corresponding code implementation that can be validated against developer-written test cases. To facilitate lightweight and reliable evaluation, we further curate a human-validated subset named NoCode-bench Verified. It covers 114 high-quality tasks across projects, where the task clarity and evaluation validity are manually verified. We use NoCode-bench to assess a range of state-of-the-art LLMs. Experimental results show that despite significant token consumption, the best task success rate remains as low as 37.72% when using the OpenHands scaffold combined with Qwen3-Coder-480B. Our analysis reveals that LLMs face key challenges in performing cross-file edits, understanding existing code modules, and accurately calling tools.

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities in supporting software developers across a wide range of programming tasks (Xu et al., 2022; Hou et al., 2024), including code generation, test generation, and program repair. These advances have led to the widespread adoption of LLM-powered software engineering

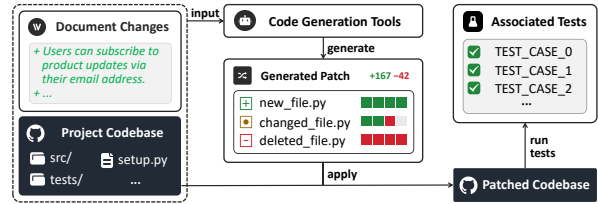


Figure 1: Overview of NL-driven feature addition tasks.

tools (Kumar et al., 2025; Pudari and Ernst, 2023), where models are typically prompted with rich code context and explicit instructions to generate or modify source code (Xia et al., 2025; Wang et al., 2025). As a result, LLMs are increasingly viewed as versatile programming assistants that can automate a wide range of developer-oriented activities.

In contrast, the effectiveness and limitations of LLMs in supporting software users, especially in incrementally improving software, are currently underexplored. Software users typically hope to add or modify specific features in existing software to meet their needs, yet usually lack familiarity with the underlying codebase or development tools (Sahay et al., 2020; Rao et al., 2024). A promising paradigm toward such user-oriented scenarios is natural language-driven feature addition. This paradigm allows users to specify new software functionality or behavioral changes purely through natural language (NL) descriptions with the underlying code automatically generated by the no-code system (Hirzel, 2023), sometimes also called “no-code development”.

To understand and improve the performance of LLMs in NL-driven feature addition, there is a need for high-quality benchmarks. Although several benchmarks have been proposed to evaluate LLMs on software engineering problems (Jiang et al., 2025), they exhibit a substantial gap with respect to no-code scenarios. Specifically, they predominantly use issue reports as task inputs, which are designed for coordination among developers (GitHub,

076 [Inc., 2025a](#)) and are typically authored by active
077 code contributors ([Bissyandé et al., 2013](#)). Prior
078 studies ([Chaparro et al., 2019](#); [Soltani et al., 2020](#))
079 also show that high-quality issue and bug reports
080 often include implementation-level technical de-
081 tails such as reproduction steps, stack traces, and
082 test cases, implicitly assuming familiarity with sys-
083 tem internals. Such developer-centric inputs differ
084 fundamentally from those in no-code development,
085 where requirements are expressed in natural lan-
086 guage without assuming programming expertise.
087 Moreover, existing benchmarks largely focus on
088 bug fixing or issue resolution, while paying lim-
089 ited attention to feature addition, despite feature
090 addition being a central activity and accounting for
091 around 60% of software maintenance efforts ([Rah-
092 man, 2019](#); [Glass, 2001](#)).

093 To fill this gap, we present NoCode-bench, a
094 new benchmark designed to evaluate LLMs’ ca-
095 pability in NL-driven feature addition. NoCode-
096 Bench identifies feature addition tasks through re-
097 lease notes and constructs natural language specifi-
098 cations from *changes in user-facing software doc-
099 umentation*, which serves as an authoritative and
100 comprehensive NL specification of software func-
101 tionality ([Aghajani et al., 2019](#)). It consists of 634
102 carefully collected tasks derived from high-quality
103 open-source projects hosted on GitHub. As illus-
104 trated in Figure 1, each instance takes documenta-
105 tion changes as input and expects the model to
106 generate corresponding code changes. The imple-
107 mentation is validated using developer-written test
108 cases. In addition, to provide a lightweight and
109 reliable evaluation option, we manually validate a
110 subset called NoCode-bench Verified. This subset
111 consists of 114 instances with their task clarity and
112 evaluation accuracy manually verified, enabling
113 reliable evaluation under limited resources.

114 We leverage both NoCode-bench Verified and
115 NoCode-bench Full to systematically evaluate a
116 range of state-of-the-art LLMs with the Agentless
117 and OpenHands frameworks. Evaluation results
118 show that, despite incurring high token consump-
119 tion, the best success rate of existing models is
120 only 37.72%, achieved when using Qwen3-Coder-
121 480B with the OpenHands scaffold. In compari-
122 son, when evaluated on SWE-bench, a developer-
123 centric benchmark collected from similar projects,
124 the same model and scaffold achieve a success rate
125 of 69.60%. This stark performance gap highlights
126 the substantial challenges of user-centric scenarios
127 and no-code development. Our analysis further re-

128 veals that LLMs face key challenges in performing
129 cross-file edits, understanding existing code mod-
130 ules, and accurately calling tools when completing
131 NL-driven feature addition tasks. Our code and
132 data are available at [https://anonymous.4open.
133 science/r/NoCode-bench-DFBC](https://anonymous.4open.science/r/NoCode-bench-DFBC).

134 In summary, our main contributions are:

- 135 • We introduce NoCode-bench, a novel user-
136 centric benchmark that evaluates NL-driven soft-
137 ware feature addition, paving the way for no-code
138 software development.
- 139 • We present a systematic five-phase construction
140 pipeline to ensure the quality and fidelity of the
141 benchmark, and further provide a manually vali-
142 dated subset for lightweight and reliable evalua-
143 tion under limited resources.
- 144 • We comprehensively evaluate multiple state-of-
145 the-art LLMs on NoCode-bench, analyze their
146 performance from both quantitative and qualita-
147 tive perspectives, and identify key factors that
148 affect model performance.

149 2 NoCode-bench

150 NoCode-bench focuses on evaluating LLMs in NL-
151 driven feature-addition scenarios. It comprises
152 634 tasks collected from 10 popular open-source
153 projects. Each task corresponds to a real-world fea-
154 ture addition explicitly documented in the project’s
155 official release notes. Each feature addition in-
156 cludes a coordinated change to the documenta-
157 tion, the code change, and the test suite. The
158 goal of each task is to produce a patch that im-
159 plements the feature specified by the documenta-
160 tion change and passes the relevant tests. The in-
161 put includes the documentation changes and the
162 complete codebase, as illustrated in Figure 1. The
163 task design aligns with documentation-first devel-
164 opment practices, such as README-driven devel-
165 opment ([Preston-Werner, 2010](#)), OpenAPI-based
166 workflows ([OpenAPI, 2025](#)), and Spec-driven de-
167 velopment ([GitHub, Inc., 2025b](#)), where documen-
168 tation changes precede and guide code modifica-
169 tions.

170 2.1 Benchmark Construction

171 We focus on high-quality, real-world software de-
172 velopment data to construct the NoCode-bench
173 through five phases (Figure 2):

174 **Phase 1: Project Selection.** The project selec-
175 tion begins with the 12 projects used in SWE-
176 bench ([Jimenez et al., 2024](#)), which are actively

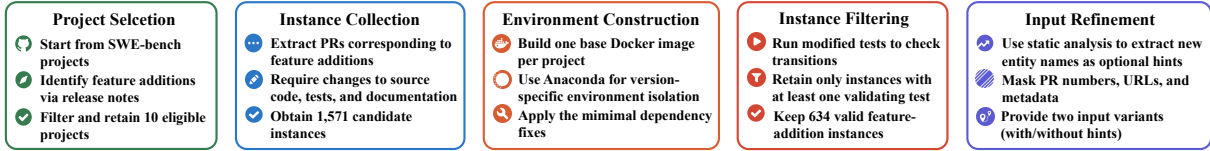


Figure 2: The workflow of building NoCode-bench.

maintained and well-documented open-source projects. Unlike SWE-bench, which constructs instances by linking GitHub issues to corresponding pull requests (PRs), we identify real-world feature addition tasks from release notes, which are written or confirmed by project developers and typically include clearly labeled feature additions, making them a reliable source of accurate task annotations. To ensure that the benchmark is reliable, representative, and suitable for NL-driven feature addition, we apply additional criteria when selecting projects. Eligible projects must: (1) provide release notes that signal sustained maintenance (Bi et al., 2020; Yu, 2009); (2) label each change item with its type (e.g., feature, enhancement), which supports accurate identification of feature-addition instances; and (3) link every change item to at least one GitHub PR, allowing us to trace documented features back to their corresponding code modifications. Following this approach, we evaluate all 12 projects from SWE-bench and find that 10 of them meet all the above conditions. The flask (Pallets, 2025) and sympy (Team, 2025) projects lack feature-labeled annotations and are thus excluded.

Phase 2: Instance Collection. For each project selected in Phase 1, we retrieve all PRs corresponding to the feature additions identified from release notes, resulting in a total of 3,263 PRs. To ensure that every instance reflects a genuine, verifiable feature update, we retain only PRs that simultaneously modify (1) source code, (2) test files, and (3) user-facing documentation. Unlike SWE-bench, which considers only source code and test modifications, we additionally require documentation changes, which serve as the NL specification provided as input to our task. Compared to issue reports, it better aligns with the no-code setting, as they describe the intended feature behavior without exposing implementation details. For each retained PR, we record the base commit, the documentation diff (i.e., the input of the NL-driven feature addition task), and the feature and test patches (i.e., the evaluation oracle). Applying these criteria yields 1,571 candidate instances, which are then processed.

Phase 3: Environment Construction. To evaluate feature implementations reliably, each instance must be executed in a reproducible environment. Unlike SWE-bench, which builds one Docker image per instance, we adopt a more scalable strategy (Yang et al., 2025c): one base Docker image per project, with version-specific isolation handled through lightweight environment management, i.e., Anaconda (Rolon-Mérette et al., 2016). This design leverages the fact that instances within a project usually share similar dependencies, reducing redundancy while preserving reproducibility. For each project, we reconstruct the environment corresponding to the feature’s base commit and install its declared dependencies. When discrepancies arise, such as unpinned or outdated library versions, we apply minimal manual fixes to restore executability. Because such issues often recur across instances of the same project, a small number of fixes typically generalize broadly. This phase ensures that all tasks run in clean and consistent containerized environments, providing a reliable foundation for automated test-based evaluation.

Phase 4: Instance Filtering. Then, we verify that each collected PR indeed corresponds to a valid feature addition. For every instance, we execute the test files modified or introduced in the PR both before and after applying the feature patch. Tests that remain passed (PASSED→PASSED) ensure regression safety, while tests that transition from failing to passing (FAILED→PASSED) serve as evidence that the new feature has been implemented. We denote these two categories as F→P tests and P→P tests, respectively. An instance is kept only if it contains at least one such validation test, providing a clear behavioral signal of feature emergence. Unlike SWE-bench, we do not discard instances whose pre-patch executions raise errors (e.g., *ImportError* or *AttributeError*), as such failures naturally arise when a feature introduces previously undefined modules, classes, or functions. Applying this filtering yields 634 validated feature addition instances used in NoCode-bench.

Phase 5: Input Refinement. During benchmark

construction, we observed that some feature additions introduce new program entities (e.g., files, classes, or methods) that are referenced in tests and implementation patches but not mentioned in the documentation change. As the evaluation is test-driven, mismatched or missing entity names can cause models to fail tests even when the generated feature is correct. Simply discarding such cases, as done in prior work (Jimenez et al., 2024), would remove legitimate feature additions and thus reduce benchmark coverage. To preserve these instances while avoiding superficial evaluation failures, we extract the names of newly introduced, test-referenced code entities, such as classes and methods, through lightweight static analysis. These names can be appended as optional hints to the model input to reduce false negatives, or omitted for a stricter evaluation setting. We release both benchmark variants to support different evaluation preferences. Additionally, we automatically mask PR numbers, URLs, and other metadata in documentation diffs to prevent models from leveraging memorized or externally fetchable information.

Together with previous phases, this yields the final set of 634 NL-driven feature addition tasks used in NoCode-bench. A final breakdown of these task instances across projects is presented in Table 1. Technical details about NoCode-bench’s construction pipeline are discussed in Appendix A.1.

2.2 NoCode-bench Verified

To complement the full benchmark of 634 instances (NoCode-bench Full) and support lightweight, reliable evaluation, we construct a manually validated subset, NoCode-bench Verified. This subset is intended for fine-grained debugging, controlled studies of model behavior, and scenarios with limited computational budgets.

Construction Process. We aim to preserve diversity across projects and target a subset of roughly 100 high-quality instances. We sample candidate tasks from all projects and perform double-blind manual verification following the guidelines of SWE-bench Verified (Chowdhury et al., 2024). Five experienced developers independently assess each sampled instance.

Annotation Criteria The annotation adopts two criteria: *task clarity* and *evaluation accuracy*. Both criteria use a coarse 0-3 scale, with lower scores indicating higher quality. The task clarity measures whether a competent engineer, with access to the codebase, can reasonably infer the expected imple-

Table 1: Instance counts of each project in NoCode-Bench Full and NoCode-Bench Verified.

Project	Owner	Full	Verified
seaborn	mwaskom	28	13
scikit-learn	scikit-learn	234	16
xarray	pydata	87	17
astropy	astropy	164	19
django	django	24	8
pylint	pylint-dev	16	6
sphinx	sphinx-doc	51	17
pytest	pytest-dev	4	1
requests	psf	2	1
matplotlib	matplotlib	24	16
Total	-	634	114

mentation from the updated documentation. The evaluation accuracy reflects whether the test patch provides appropriate and targeted functional coverage of the described feature. Annotators may additionally flag instances exhibiting other major issues.

Annotation Results We discard instances whose task clarity or evaluation accuracy receives a score of 2 or higher, as well as any instance flagged for major issues. This process removes 124 candidates and yields 114 verified instances. Inter-annotator agreement, measured by Cohen’s kappa, is 0.42, which is comparable to 0.39 in SWE-bench Verified and indicative of moderate but acceptable reliability.

NoCode-bench Verified thus offers a curated, dependable subset for controlled LLM evaluation while remaining representative of documentation-driven feature addition scenarios. The complete annotation process and results are documented in Appendix A.2.

3 Benchmark Characteristics

To better understand the unique characteristics and challenges posed by NoCode-bench, we analyze the distribution of instances in NoCode-bench and SWE-bench across dimensions related to input complexity, localization difficulty, and editing difficulty, as shown in Figure 3. For better visualization, we remove the top 1% of the data in Figure 3, but we keep them during the following analysis.

Input complexity. Unlike SWE-bench, which uses issue descriptions as input, NoCode-bench takes document changes as input and requires the LLM to understand the semantic differences before and after the change. Long document changes pose extra challenges for the LLM to comprehend the

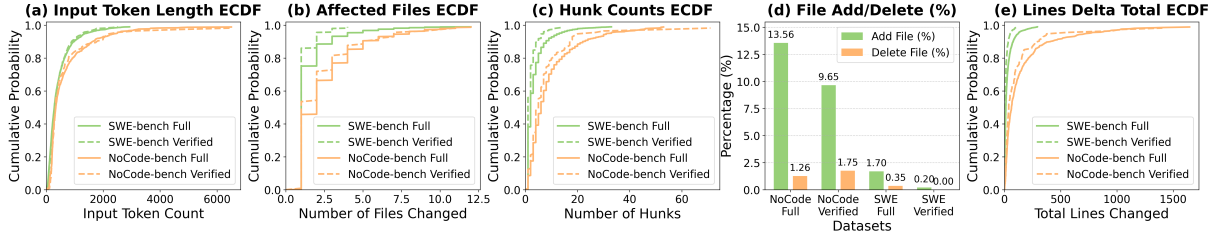


Figure 3: Comparison of golden patch distributions between NoCode-bench and SWE-bench along dimensions related to localization and editing difficulty. ECDF stands for “empirical cumulative distribution function”.

task and generate appropriate outputs. As shown in Figure 3 (a), the average length of issue descriptions in SWE-bench Full and Verified is 480.37 and 447.72, respectively, while the average length of document changes in NoCode-bench Full and Verified is 739.06 and 820.65, which is nearly twice as long as that in SWE-bench. This indicates that tasks in NoCode-bench require LLMs to capture information from longer and more complex task inputs, posing greater challenges to their abilities in long-text comprehension, identifying key information, and attending to fine-grained details.

Localization Difficulty. Accurate localization in the codebase is a necessary condition for generating the correct patches. Localization is particularly challenging when multiple files are involved or when the changes are scattered across many code regions. As shown in Figure 3 (b,c), on average, golden patches in NoCode-bench Full involve edits to 2.65 files and 10.32 hunks, whereas in SWE-bench Full, they involve only 1.66 files and 4.12 hunks. In the Verified subsets, the number of files involved is 2.39 for NoCode-bench and 1.24 for SWE-bench, with corresponding hunk counts of 8.96 and 2.43. Moreover, Figure 3 (d) reveals that NoCode-bench includes a much higher proportion of patches that involve adding or deleting files. 13.56% of instances in NoCode-bench Full require new files, while the percentage is only 1.70% in SWE-bench Full. These results indicate that our benchmark requires LLMs to locate multiple code blocks spread across different files within large codebases, which goes far beyond the demands of issue-solving tasks in SWE-bench, leading to localization complexity.

Editing Difficulty. Once the correct edit location is identified, the correct code segments must be generated to realize the intended functionality. Larger edits typically increase the chance of introducing syntactic or semantic errors. We measure the total number of lines changed in golden patches, i.e., the

ground truth edits excluding changes to tests and documentation files, as shown in Figure 3 (e). In NoCode-bench Full, golden patches involve an average of 179.12 lines of edits, and 179.22 lines in NoCode-bench Verified. In contrast, golden patches in SWE-bench Full average 37.71 lines of edits, and only 14.32 lines in SWE-bench Verified. Furthermore, we find that SWE-bench patches are small, with nearly all patches under 200 lines. In contrast, NoCode-bench includes a substantial portion of larger edits, with nearly 20% of patches exceeding 200 lines. These larger modifications may require stronger code generation capabilities and greater coherence across multi-line edits.

4 Experimental Setup

Model Selection. Completing a task in NoCode-bench typically requires a deep understanding of the target project, demanding strong reasoning ability and the capacity to handle long contexts. Considering these requirements, we selected 8 state-of-the-art LLMs for evaluation on NoCode-bench Verified. Four of them are open-source models, i.e., DeepSeek-v3-0324 (Liu et al., 2024), DeepSeek-R1-0528 (Guo et al., 2025a), Qwen3-235B-A22B (Yang et al., 2025a), and Qwne3-Coder-480B-A35B, where DeepSeek-R1 is a reasoning model and Qwen3-235B-A22B is used in thinking mode. The other four are closed-source models, i.e., Claude-4-Sonnet (Anthropic, 2024), GPT-4o-2024-11-20 (OpenAI, 2024), Gemini-2.5-Pro-Exp-03-25 (Comanici et al., 2025), and GPT-5-mini (OpenAI, 2025), where Gemini-2.5-Pro is a reasoning model and GPT-5-mini is used in medium thinking. For brevity, version numbers will be omitted when referencing model names throughout the paper. NoCode-bench Full is significantly larger than NoCode-bench Verified, with 5.6 times the number of instances. Due to the high cost of Claude-4-Sonnet, Qwen3-Coder, GPT-5-mini, and Gemini-2.5-Pro (total estimated at

Table 2: Performance comparison across different LLMs on NoCode-bench Verified.

Model	Agentless							OpenHands						
	Success%	Applied%	RT%	FV-Micro	FV-Macro	File%	#Token	Success%	Applied%	RT%	FV-Micro	FV-Macro	File%	#Token
Open-Source Model														
Qwen3-235B	13.16%	100.00%	76.32%	8.75%	22.39%	42.98%	0.12M	7.89%	64.91%	47.37%	1.96%	14.03%	45.61%	0.24M
DeepSeek-v3	21.05%	98.25%	78.95%	7.96%	32.80%	57.14%	0.30M	11.40%	65.79%	49.12%	1.68%	18.29%	56.14%	0.59M
DeepSeek-R1	25.44%	93.86%	73.68%	10.87%	35.52%	50.47%	0.50M	7.02%	47.37%	46.49%	0.47%	10.86%	38.60%	0.39M
Qwen3-Coder-480B	28.07%	100.00%	82.46%	9.03%	40.12%	54.39%	0.45M	37.72%	92.11%	84.21%	10.96%	50.45%	63.16%	2.91M
Closed-Source Model														
GPT-4o	13.16%	100.00%	67.54%	7.97%	22.32%	47.37%	0.11M	5.26%	63.16%	57.89%	0.45%	6.76%	34.21%	0.58M
Gemini-2.5-Pro	12.28%	100.00%	74.56%	6.22%	20.55%	48.25%	0.29M	0.00%	54.39%	61.40%	0.01%	0.29%	0.00%	0.47M
Claude-4-Sonnet	28.07%	100.00%	79.82%	8.47%	38.48%	57.89%	0.33M	25.44%	68.42%	69.30%	11.25%	36.48%	67.54%	2.32M
GPT-5-mini	26.32%	98.25%	82.46%	10.63%	37.82%	53.98%	0.21M	36.84%	87.72%	80.70%	12.61%	49.64%	64.60%	1.69M
Average	20.94%	98.80%	76.97%	8.74%	31.25%	51.56%	0.29M	16.45%	67.98%	62.06%	4.92%	23.35%	46.23%	1.15M

Table 3: Performance comparison across different LLMs on NoCode-bench Full with Agentless.

Model	Success%	Applied%	RT%	FV-Micro	FV-Macro	File%	#Token
Qwen3-235B	6.62%	99.37%	75.24%	10.78%	15.77%	39.46%	0.19M
DeepSeek-v3	11.83%	99.84%	77.44%	16.02%	24.09%	47.48%	0.31M
DeepSeek-R1	14.83%	95.90%	79.34%	10.53%	25.78%	44.89%	0.57M
GPT-4o	9.31%	97.00%	71.14%	8.68%	17.81%	43.58%	0.11M
Average	10.65%	98.03%	75.79%	11.50%	20.86%	43.85%	0.30M

USD 5,124) running on NoCode-bench Full, we focus on NoCode-bench Verified for those models.

Scaffold Selection. Due to the complexity of NL-driven feature addition tasks, directly prompting LLMs to generate the required changes is impractical. Therefore, we adopt two state-of-the-art general scaffolds for solving software engineering tasks: a pipeline-based method, Agentless (Xia et al., 2025), and an agent-based method, OpenHands (Wang et al., 2025). They respectively represent two distinct paradigms to automate SE tasks using LLMs, and both have attained top scores on SWE-bench, detailed in Appendix B.2.

Evaluation Metrics. Following prior studies (Jimenez et al., 2024; Zan et al., 2025; Guo et al., 2025b), our end-to-end evaluation primarily considers the success rate (i.e, the proportion of instances for which given patch passes the accompanied test suite) of feature addition tasks (**Success%**), the success rate of patch application (**Applied%**), and the average token cost (**#Token**). Additionally, two other metrics are used: **File Matched Rate (File%)** measures alignment between model-edited and reference files, and **Regression Tests Pass Rate (RT%)** assesses the preservation of existing functionality. We introduce a new metric named **Feature Validation Rate (FV)** to evaluate whether the added feature passes the F→P tests, which is reported at both micro and macro levels to reflect overall test-level performance and instance-level robustness, respectively. Detailed definitions are

provided in Appendix B.1.

5 Results

5.1 Performance on NoCode-bench Verified

The evaluation results are shown in Table 2. We find that even the most advanced LLMs still exhibit limited performance in NL feature addition. Specifically, among the evaluated models, Qwen3-Coder-480B achieves the highest success rate of 37.72% on NoCode-bench Verified with the OpenHands scaffold. The closed-source models GPT-4o and Gemini-2.5-Pro exhibit lower scores, with success rates of 5.26% and 0%. These results highlight significant performance differences across models of varying sizes and architectures under the same scaffold. Notably, the success rate on NoCode-bench Verified is much lower than that on SWE-bench Verified. For example, Claude-4-Sonnet combined with OpenHands achieves 70.4% on SWE-bench Verified but only 25.44% on NoCode-bench Verified, implying that NL-driven feature addition tasks present a substantially greater challenge compared to issue-solving tasks. A detailed comparison of scaffold performance and further analysis can be found in Appendix C.1. We further discuss the unique challenges associated with NL-driven feature addition in Section 5.3.

5.2 Performance on NoCode-bench Full

Considering the better average performance of Agentless and the high cost of OpenHands, we use Agentless to evaluate the models on NoCode-bench Full, and report the results in Table 3. We find that all four models show decreased performance on NoCode-bench Full compared to NoCode-bench Verified. Specifically, the average success rate of the four used models drops from 18.20% to 10.65%, and the file matched rate drops from 49.49% to

Table 4: Performance across different LLMs on NoCode-bench for *single-file* and *multi-file* tasks.

Model	Scaffold	Single-File Modification		Multi-File Modification	
		Success%	Applied%	Success%	Applied%
NoCode-bench Verified					
Total Instances		61		53	
Qwen3-235B	Agentless	21.31% (13)	100.00% (61)	3.77% (2)	100.00% (53)
	OpenHands	13.11% (8)	65.57% (40)	1.89% (1)	64.15% (34)
DeepSeek-v3	Agentless	31.15% (19)	100.00% (61)	9.43% (5)	96.23% (51)
	OpenHands	14.75% (9)	67.21% (41)	7.55% (4)	64.15% (34)
DeepSeek-R1	Agentless	40.98% (25)	98.36% (60)	7.55% (4)	88.68% (47)
	OpenHands	11.48% (7)	47.54% (29)	1.89% (1)	47.17% (25)
Qwen3-Coder-480B	Agentless	40.98% (25)	100.00% (61)	13.21% (7)	100.00% (53)
	OpenHands	52.46% (32)	93.44% (57)	20.75% (11)	90.57% (48)
GPT-4o	Agentless	19.67% (12)	100.00% (61)	5.66% (3)	100.00% (53)
	OpenHands	6.56% (4)	68.85% (42)	3.77% (2)	56.60% (30)
Gemini-2.5-Pro	Agentless	19.67% (12)	100.00% (61)	3.77% (2)	100.00% (53)
	OpenHands	0.00% (0)	52.46% (32)	0.00% (0)	56.60% (30)
Claude-4-Sonnet	Agentless	44.26% (27)	100.00% (61)	9.43% (5)	100.00% (53)
	OpenHands	39.34% (24)	68.85% (42)	9.43% (5)	67.92% (36)
GPT-5-mini	Agentless	39.34% (24)	96.72% (59)	11.32% (6)	100.00% (53)
	OpenHands	52.46% (32)	91.80% (56)	18.87% (10)	83.02% (44)
NoCode-bench Full					
Total Instances		290		344	
Qwen3-235B	Agentless	11.38% (33)	99.66% (289)	2.62% (9)	99.13% (341)
DeepSeek-v3	Agentless	20.00% (58)	100.00% (290)	4.94% (17)	99.71% (343)
DeepSeek-R1	Agentless	24.83% (72)	96.55% (280)	6.40% (22)	95.35% (328)
GPT-4o	Agentless	15.17% (44)	98.97% (287)	4.36% (15)	95.35% (328)

43.85%. This drop may be primarily due to the increased scale and complexity, as shown in Section 3, as well as the presence of some noise (e.g., ambiguous document changes or overly specific tests) in NoCode-bench Full. These factors make it more challenging. A detailed analysis of this phenomenon is provided in Appendix C.2. It is worth noting that the ranking and relative performance of the models on NoCode-bench Verified are preserved on NoCode-bench Full except Qwen3-235B and GPT-4o. A larger decline of Qwen3-235B than GPT-4o suggests that model robustness can vary as the evaluation shifts from a curated, high-quality subset to a larger set in real-world settings. We further explain this phenomenon in Appendix C.2.

5.3 Failure Analysis

To investigate why LLMs fail to solve tasks on NoCode-bench, we manually examine 160 failed instances (randomly sample 10 instances per model per scaffold) and identify the following primary reasons for LLM failures on NoCode-bench.

Lack of cross-file editing capability. In NoCode-bench Verified and NoCode-bench Full, only 53.5% and 45.7% of the instances can be solved by editing a single file, respectively, indicating that NoCode-bench requires the LLM to be good at performing cross-file edits. To further investigate how cross-file editing capability affects performance on feature addition tasks, we compute the number of modified files in the golden patch

for each instance and divide the dataset into two groups: *single-file* (#modified_files = 1) and *multi-file* (#modified_files > 1). We analyze the performance of LLMs on *single-file* and *multi-file* tasks, and present the results in Table 4. On NoCode-bench Verified and NoCode-bench Full, the best-performing model solves only 20.75% and 6.40% of the *multi-file* instances, respectively. These results indicate that the lack of cross-file editing capability is one of the main reasons for LLM failures on NoCode-bench.

Lack of comprehensive understanding of existing code modules. Through manual inspection, we find that LLMs often attempt to implement new features by directly modifying the code of existing features, and 39.38% of the examined instances exhibit this error. This can lead to the new feature overriding existing features, resulting in failure in regression tests. A typical example is the patch generated by DeepSeek-v3+Agentless on the task “*pylint-7869*” (detailed in Appendix C.3). This implies that LLMs lack a comprehensive understanding of existing code modules and face challenges in figuring out what and how existing code modules should be edited.

Parsing failure in LLM’s tool-calling. This type of error occurs during the tool-calling phase of the OpenHands CodeAct Agent. For Gemini-2.5-Pro, all the examined instances exhibit this error. OpenHands relies on a set of predefined tools that allow the model to access the repository, requiring the LLM to return detailed tool-call information in a strictly defined JSON format. This makes the system highly dependent on the output formatting capability of LLMs. However, nearly all LLMs fail to invoke tools in the correct output format for some samples, with Gemini-2.5-Pro being particularly prone to this issue, which directly results in a success rate of 0 on NoCode-bench Verified. A typical case is shown in Appendix C.4, where Gemini repeatedly fails to format tool calls correctly, eventually causing the context window to overflow and preventing successful patch generation. This case demonstrates that the model’s capability to generate formatted output during the tool-calling significantly impacts the agent’s ability to operate on the codebase and can even directly lead to failure.

6 Related Work

6.1 Benchmarks for Evaluating LLMs in SE

Evaluating the capabilities of LLMs in software engineering (SE) has attracted increasing attention, leading to the development of several benchmarks. Among these, SWE-bench (Jimenez et al., 2024) has emerged as one of the most influential benchmarks for assessing LLMs on issue resolution tasks. A line of follow-up work has extended SWE-bench to multilingual (Zan et al., 2024, 2025), multimodal (Yang et al., 2025b; Guo et al., 2025b), and industrial (Rondon et al., 2025) scenarios, promoting broader language and scenario coverage in issue-solving tasks. FEA-Bench (Li et al., 2025b) aims to evaluate LLMs on feature addition tasks, with detailed code-level specifications, docstrings, and pull request descriptions as task inputs. However, these benchmarks rely on developer-centric inputs, such as issue reports, and assume familiarity with the codebase and the implementation-level terminology (Chaparro et al., 2019; GitHub, Inc., 2025a). As a result, they provide limited insight into LLMs’ capabilities in user-centric scenarios. Also, most of them focus on bug fixing or issue solving, while largely overlooking feature addition, which accounts for about 60% of software maintenance activities (Rahman, 2019).

Table 5 provides a detailed comparison between NoCode-bench and related benchmarks, highlighting several distinctive features of NoCode-bench: (1) NoCode-bench adopts a more user-centric no-code setting where user-facing documentation changes are provided as inputs, rather than code-level specifications; (2) NoCode-bench does not constrain instances to those introducing new components, recognizing that many feature additions extend existing code; (3) NoCode-bench identifies feature-adding PRs based on developer-maintained release notes to reduce noise (Xu et al., 2024; Zhang et al., 2025). In addition, we manually construct a verified subset to enable lightweight yet reliable evaluation.

6.2 Automating SE tasks with LLMs

Recently, LLM-based methods (Jiang et al., 2025) have been proposed for automating SE tasks, mainly including pipeline-based and agent-based methods. Pipeline-based methods typically decompose tasks into stages (e.g., localization, repair, and validation) based on prior knowledge and require LLMs to follow a fixed workflow. Examples

Table 5: Comparison with existing benchmarks.

Aspect	SWE-bench	GITS-Eval	FEA-Bench	NoCode-bench
Development scenario	Developer-centric	Developer-centric	Developer-centric	User-centric
Task type	Issue resolution	Issue resolution	Feature addition	Feature addition
Task input	Issue description	Issue description	Code-level specs, PR description	Documentation changes
Instance selection	N/A	N/A	PRs introducing new components	PRs containing doc changes
Feature identification	N/A	N/A	Prompting LLMs	Mining release notes
Verified subset	Yes	Yes	No	Yes
Open-source	Yes	No	Yes	Yes

include Agentless (Xia et al., 2025), and PatchPilot (Li et al., 2025a). Agent-based methods equip agents with various tools (e.g., bash and language servers) to access and modify the codebase, and rely on the decision-making and tool invocation capabilities of LLMs. Examples include OpenHands (Wang et al., 2025), RepairAgent (Bouzenia et al., 2025), and ExecutionAgent (Bouzenia and Pradel, 2025). While these methods have achieved significant advances, they are typically evaluated in a developer-centric way, where the inputs include implementation-level details and assume programming expertise. In contrast, feature addition tasks in NoCode-bench are specified with user-facing documentation changes. This mismatch introduces new and unique challenges for applying prior methods to NoCode-bench.

7 CONCLUSION

We introduce NoCode-bench, a challenging benchmark for evaluating LLMs’ software engineering capabilities in NL-driven feature addition. NoCode-bench is constructed with a systematic five-phase pipeline. It consists of 634 real-world feature addition tasks identified from developer-maintained release notes. Given documentation changes within a repository as input, NoCode-bench requires systems to generate a patch that adds a new feature. We evaluate 8 advanced LLMs using two representative scaffolds, Agentless and OpenHands. Our results show that the best-performing model succeeds in only 37.72% of instances in NoCode-bench Verified and 14.83% in NoCode-bench Full. This demonstrates that NL-driven feature addition is a uniquely difficult and unsolved problem. Moreover, we analyze failure cases on NoCode-bench to guide future improvements in cross-file editing, code module understanding, and tool-calling capabilities.

LIMITATIONS

NoCode-bench, while designed to reflect realistic documentation-driven feature additions, has several limitations.

First, the benchmark is constructed from historical GitHub repositories, raising the possibility of data leakage, as such data may appear in the training corpora of modern large language models. We mitigate this risk by masking PR numbers and other metadata, and the consistently low success rates across strong LLMs suggest limited practical impact, though the risk cannot be completely ruled out.

Second, the current benchmark focuses on well-maintained open-source Python projects with structured release notes. This choice ensures clean and interpretable task definitions but limits language and ecosystem diversity. Although the construction methodology is language-agnostic and can be extended to other settings, future work is needed to broaden project coverage.

In addition, our evaluation relies on two scaffolds (i.e., Agentless and OpenHands), which represent strong baselines but do not cover the full spectrum of possible code-execution or agent frameworks. Performance under other scaffolds remains unexplored. Nonetheless, the modest results observed across multiple LLMs indicate that NoCode-bench poses meaningful challenges that extend beyond scaffold-specific behaviors.

Ethics Statement

NoCode-bench is constructed entirely from publicly available open-source software projects, all of which are used in accordance with their original licenses. The benchmark is intended for research purposes and does not involve any sensitive personal data. Both the dataset and the accompanying code will be released for academic use. Users should be aware that code produced by models when performing benchmark tasks could potentially be harmful to computer systems, and we recommend evaluating model outputs in controlled environments such as Docker containers.

References

Emad Aghajani, Csaba Nagy, Olga Lucero Vega-Márquez, Mario Linares-Vásquez, Laura Moreno, Gabriele Bavota, and Michele Lanza. 2019. Software documentation issues unveiled. In *2019 IEEE/ACM*

41st International Conference on Software Engineering (ICSE), pages 1199–1210.

Anthropic. 2024. Claude 4. <https://www.anthropic.com/news/claude-4>. Accessed: 2025-01.

Tingting Bi, Xin Xia, David Lo, John Grundy, and Thomas Zimmermann. 2020. An empirical study of release note production and usage in practice. *IEEE Transactions on Software Engineering*, 48(6):1834–1852.

Tegawendé F. Bissyandé, David Lo, Lingxiao Jiang, Laurent Réveillère, Jacques Klein, and Yves Le Traon. 2013. Got issues? who cares about it? a large scale investigation of issue trackers from github. In *2013 IEEE 24th International Symposium on Software Reliability Engineering (ISSRE)*, pages 188–197.

Islem Bouzenia, Premkumar Devanbu, and Michael Pradel. 2025. RepairAgent: An autonomous, LLM-based agent for program repair. In *International Conference on Software Engineering (ICSE)*.

Islem Bouzenia and Michael Pradel. 2025. You name it, i run it: An llm agent to execute tests of arbitrary projects. In *Proceedings of the 34rd ACM SIGSOFT International Symposium on Software Testing and Analysis*.

Oscar Chaparro, Carlos Bernal-Cárdenas, Jing Lu, Kevin Moran, Andrian Marcus, Massimiliano Di Penta, Denys Poshyvanyk, and Vincent Ng. 2019. Assessing the quality of the steps to reproduce in bug reports. In *Proceedings of the 2019 27th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering*, pages 86–96.

Neil Chowdhury, James Aung, Chan Jun Shern, Oliver Jaffe, Dane Sherburn, Giulio Starace, Evan Mays, Rachel Dias, Marwan Aljubeih, Mia Glaese, Carlos E. Jimenez, John Yang, Leyton Ho, Tejal Patwardhan, Kevin Liu, and Aleksander Madry. 2024. [Introducing SWE-bench verified](#).

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

GitHub, Inc. 2025a. GitHub Issues · Project planning for developers. <https://github.com/features/issues>.

GitHub, Inc. 2025b. GitHub Spec Kit: Toolkit to help you get started with Spec-Driven Development. <https://github.com/github/spec-kit>.

R.L. Glass. 2001. Frequently forgotten fundamental facts about software engineering. *IEEE Software*, 18(3):112–111.

768	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,	OpenAPI. 2025. Usage openapi generator. https://openapi-generator.tech/docs/usage .	824
769	Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong		825
770	Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a.		
771	Deepseek-r1: Incentivizing reasoning capability in	Pallets. 2025. Flask: Micro framework for building web	826
772	llms via reinforcement learning. <i>arXiv preprint</i>	applications. https://flask.palletsprojects.com/en/stable/ .	827
773	<i>arXiv:2501.12948</i> .		828
774	Lianghong Guo, Wei Tao, Runhan Jiang, Yanlin Wang,	Tom Preston-Werner. 2010. Readme driven develop-	829
775	Jiachi Chen, Xilin Liu, Yuchi Ma, Mingzhi Mao,	ment. https://tom.preston-werner.com/2010/08/23/readme-driven-development .	830
776	Hongyu Zhang, and Zibin Zheng. 2025b. Omnigirl:		831
777	A multilingual and multimodal benchmark for github	Rohith Pudari and Neil A Ernst. 2023. From copilot to	832
778	issue resolution. In <i>Proceedings of the 34rd ACM</i>	pilot: Towards ai supported software development.	833
779	<i>SIGSOFT International Symposium on Software Test-</i>	<i>arXiv preprint arXiv:2303.04142</i> .	834
780	<i>ing and Analysis</i> .		
781	Martin Hirzel. 2023. Low-code programming models.	Mohammad Masudur Rahman. 2019. Supporting code	835
782	<i>Communications of the ACM</i> , 66(10):76–85.	search with context-aware, analytics-driven, effec-	836
783	Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong	tive query reformulation. In <i>2019 IEEE/ACM 41st</i>	837
784	Wang, Li Li, Xiapu Luo, David Lo, John Grundy,	<i>International Conference on Software Engineering:</i>	838
785	and Haoyu Wang. 2024. Large language models	<i>Companion Proceedings (ICSE-Companion)</i> , pages	839
786	for software engineering: A systematic literature re-	226–229.	840
787	view. <i>ACM Transactions on Software Engineering</i>	Nikitha Rao, Jason Tsay, Kiran Kate, Vincent Hellen-	841
788	<i>and Methodology</i> , 33(8):1–79.	doorn, and Martin Hirzel. 2024. Ai for low-code for	842
789	Zhonghao Jiang, David Lo, and Zhongxin Liu. 2025.	ai. In <i>Proceedings of the 29th International Confer-</i>	843
790	Agentic software issue resolution with large language	<i>ence on Intelligent User Interfaces</i> , pages 837–852.	844
791	models: A survey. <i>arXiv preprint arXiv:2512.22256</i> .		
792	Carlos E Jimenez, John Yang, Alexander Wettig,	Damien Rolon-Mérette, Matt Ross, Thaddé Rolon-	845
793	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R	Mérette, and Kinsey Church. 2016. Introduction to	846
794	Narasimhan. 2024. SWE-bench: Can language mod-	anaconda and python: Installation and setup. <i>Quant.</i>	847
795	els resolve real-world github issues? In <i>The Twelfth</i>	<i>Methods Psychol</i> , 16(5):S3–S11.	848
796	<i>International Conference on Learning Representa-</i>	Pat Rondon, Renyao Wei, José Cambronero, Jürgen	849
797	<i>tions</i> .	Cito, Aaron Sun, Siddhant Sanyam, Michele Tufano,	850
798	Yulia Kumar, Israel Akinwunmi, and Dov Kruger. 2025.	and Satish Chandra. 2025. Evaluating agent-based	851
799	Evaluating the advantage of an ai-native ide cursor on	program repair at google. In <i>2025 IEEE/ACM 47th</i>	852
800	programmer performance. In <i>2025 IEEE Integrated</i>	<i>International Conference on Software Engineering:</i>	853
801	<i>STEM Education Conference (ISEC)</i> , pages 1–8.	<i>Software Engineering in Practice (ICSE-SEIP)</i> .	854
802	Hongwei Li, Yuheng Tang, Shiqi Wang, and Wenbo	Apurvanand Sahay, Arsene Indamutsa, Davide Di Rus-	855
803	Guo. 2025a. Patchpilot: A cost-efficient software	cio, and Alfonso Pierantonio. 2020. Supporting the	856
804	engineering agent with early attempts on formal veri-	understanding and comparison of low-code develop-	857
805	fication. In <i>Forty-second International Conference</i>	ment platforms. In <i>2020 46th Euromicro Conference</i>	858
806	<i>on Machine Learning</i> .	<i>on Software Engineering and Advanced Applications</i>	859
807	Wei Li, Xin Zhang, Zhongxin Guo, Shaoguang Mao,	(SEAA), pages 171–178.	860
808	Wen Luo, Guangyue Peng, Yangyu Huang, Houfeng	Mozhan Soltani, Felienne Hermans, and Thomas Bäck.	861
809	Wang, and Scarlett Li. 2025b. FEA-bench: A bench-	2020. The significance of bug report elements. <i>Em-</i>	862
810	mark for evaluating repository-level code generation	<i>pirical Software Engineering</i> , 25(6):5255–5294.	863
811	for feature implementation. In <i>Proceedings of the</i>	SymPy Development Team. 2025. Sympy: Python	864
812	<i>63rd Annual Meeting of the Association for Computa-</i>	library for symbolic mathematics. https://docs.sympy.org/latest/index.html .	865
813	<i>tional Linguistics (Volume 1: Long Papers)</i> .		866
814	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu,	867
815	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi	868
816	Deng, Chenyu Zhang, Chong Ruan, and 1 others.	Song, Bowen Li, Jaskirat Singh, Hoang H. Tran,	869
817	2024. Deepseek-v3 technical report. <i>arXiv preprint</i>	Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian,	870
818	<i>arXiv:2412.19437</i> .	Yanjun Shao, Niklas Muennighoff, Yizhe Zhang,	871
819	OpenAI. 2024. <i>Gpt-4o system card</i> . <i>Preprint</i> ,	Binyuan Hui, and 5 others. 2025. Openhands: An	872
820	<i>arXiv:2410.21276</i> .	open platform for AI software developers as gener-	873
821	OpenAI. 2025. Introducing gpt-5. https://openai.com/index/introducing-gpt-5/ . Ac-	alist agents. In <i>The Thirteenth International Confer-</i>	874
822	cessed: 2025-10-17.	<i>ence on Learning Representations</i> .	875
823			

876	Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2025. Demystifying llm-based software engineering agents. <i>Proceedings of the ACM on Software Engineering</i> , 2(FSE):801–824.	A.1.1 Project and Release Note Processing	931
877		While the main text describes the project selection criteria, here we detail how release notes are processed in practice. We treat each feature-labeled item in a release note as an <i>update entry</i> . Because release note formats vary substantially across projects (e.g., Markdown, HTML, or Jupyter Notebook), we implement lightweight, project-specific parsers to normalize release note structures. These parsers extract: (1) the textual description of each feature entry, and (2) references to the corresponding GitHub PRs. This information will be used for subsequent crawling and filtering.	932
878		As illustrated in Phase 1 of Figure 4, we can locate the release notes on the project’s homepage. The change log for each version contains clearly labeled feature entries, each accompanied by a link to the corresponding PR. We extracted the text part (i.e., Added the objects.Text stat) and PR number (i.e., #3051) from this example.	933
879			934
880	Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022. A systematic evaluation of large language models of code. In <i>Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming</i> , pages 1–10.		935
881			936
882			937
883			938
884			939
885	Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is inevitable: An innate limitation of large language models. <i>arXiv preprint arXiv:2401.11817</i> .		940
886			941
887			942
888			943
889	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .		944
890			945
891			946
892			947
893			948
894	John Yang, Carlos E. Jimenez, Alex L. Zhang, Kilian Lieret, Joyce Yang, Xindi Wu, Ori Press, Niklas Muennighoff, Gabriel Synnaeve, Karthik R. Narasimhan, Diyi Yang, Sida I. Wang, and Ofir Press. 2025b. SWE-bench multimodal: Do ai systems generalize to visual software domains? In <i>The Thirteenth International Conference on Learning Representations</i> .	A.1.2 PR Crawling and Metadata Extraction	951
895		For each PR linked from release notes, we automatically crawl the full PR metadata using the GitHub API. We extract the base commit (i.e., the commit on which the feature branch was originally developed) by identifying the merge base between the PR branch and the main branch at merge time. This base commit is used to reconstruct the pre-feature codebase. To construct benchmark inputs and evaluation oracles, we compute three diffs for each retained PR: (1) the documentation diff, which serves as the task input; (2) the feature patch, which contains source code changes; and (3) the test patch, which is used to verify the correctness of the feature implementation.	952
896			953
897			954
898			955
899			956
900			957
901			958
902	John Yang, Kilian Lieret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. 2025c. SWE-smith: Scaling data for software engineering agents. In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	A.1.3 Environment Reconstruction and Dependency Repair	966
903		Dependency Repair	967
904		This section details the environment construction process summarized in the main text. For each project, we construct a single base Docker image and manage instance-level version isolation using Anaconda environments. Instances are first grouped by minor versions using automated scripts, which parse version-related files such as <i>setup.py</i> , <i>__init__.py</i> , and release notes. We then manually inspect the format and location of environment configuration files, such as <i>README.md</i> , <i>requirements.txt</i> , <i>pyproject.toml</i> , or <i>pom.xml</i> to identify how dependencies are declared. Automated scripts are developed to extract and install these dependen-	968
905			969
906			970
907			971
908			972
909	Liguo Yu. 2009. Mining change logs and release notes to understand software maintenance and evolution. <i>CLEI Electronic Journal</i> , 12(2):1–1.		973
910			974
911			975
912	Daoguang Zan, Zhirong Huang, Wei Liu, Hanwu Chen, Linhao Zhang, Shulin Xin, Lu Chen, Qi Liu, Xiaojian Zhong, Aoyan Li, and 1 others. 2025. Multi-swe-bench: A multilingual benchmark for issue resolving. <i>arXiv preprint arXiv:2504.02605</i> .		976
913			977
914			978
915			979
916			980
917	Daoguang Zan, Zhirong Huang, Ailun Yu, Shaoxin Lin, Yifan Shi, Wei Liu, Dong Chen, Zongshuai Qi, Hao Yu, Lei Yu, and 1 others. 2024. Swe-bench-java: A github issue resolving benchmark for java. <i>arXiv preprint arXiv:2408.14354</i> .		
918			
919			
920			
921			
922	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, and 1 others. 2025. Siren’s song in the ai ocean: A survey on hallucination in large language models. <i>Computational Linguistics</i> .		
923			
924			
925			
926			
927	A BENCHMARK DETAILS		
928	A.1 Construction Details		
929	This appendix provides implementation-level details of the NoCode-bench construction pipeline.		
930			

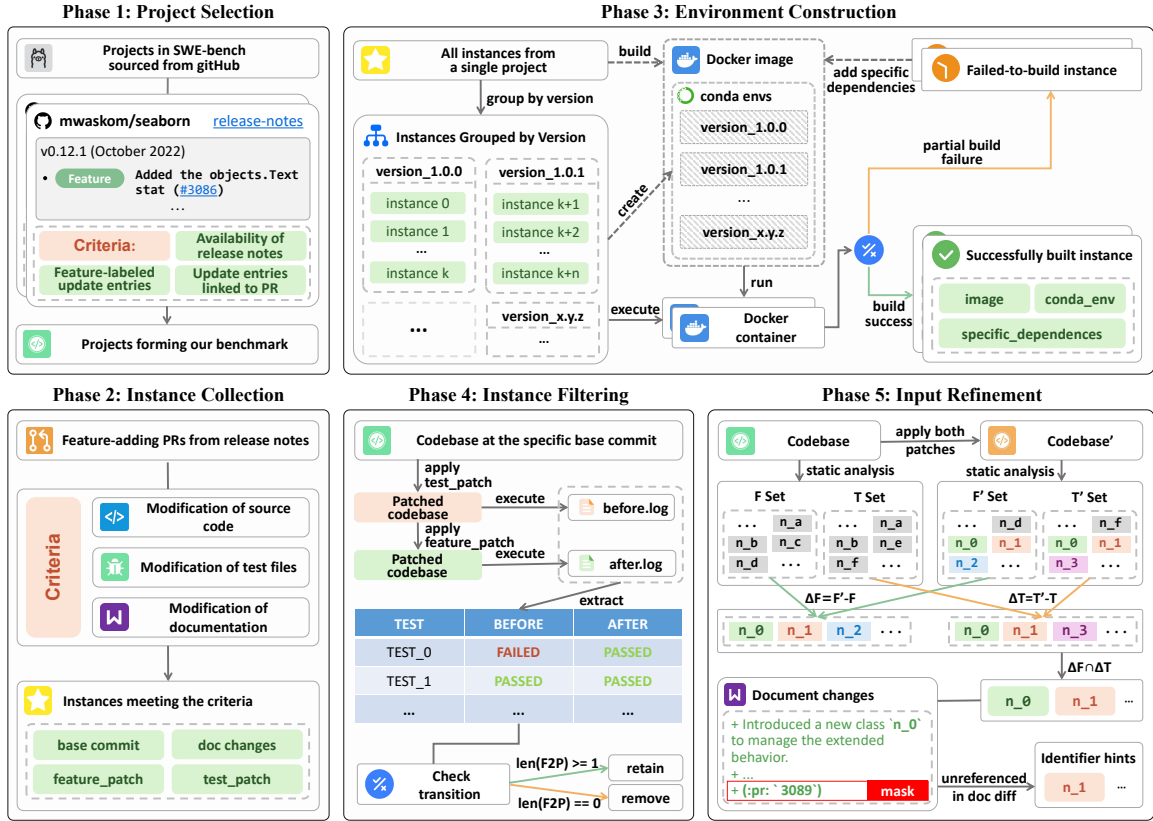


Figure 4: The workflow of building our benchmark.

cies accordingly. Some instances fail to build or execute due to loosely pinned or outdated dependencies. For such cases, we log build and runtime errors and apply minimal manual fixes, such as pinning a dependency version or adding a missing library. Each fix is recorded as a reproducible shell script. Because instances within the same project often share similar dependency issues, a single fix typically applies to many instances, resulting in limited additional human effort.

A.1.4 Test Execution and Status Extraction

To verify each instance as a valid feature addition, we execute all test files that are modified or added in the *test_patch*. For each instance, we collect two sets of logs, i.e., *before.log* and *after.log*. *before.log* is obtained by applying the *test_patch* to the base commit and then executing the corresponding test suite, while *after.log* is obtained by applying both the *test_patch* and the *feature_patch* to the base commit, followed by test execution.

We support multiple test frameworks (e.g., *pytest*, *unittest*) by normalizing their outputs into a unified representation of per-test execution status. Test case identifiers and their corresponding outcomes are extracted using framework-specific

parsers. These normalized logs are then used to determine test status transitions, which serve as the basis for instance filtering as described in the main text.

A.1.5 Identifier Hint Extraction

This section provides implementation details for identifier hint extraction. Given an instance, we first switch the codebase to its base commit. We then identify the list of files modified by the *feature_patch* and the *test_patch*, denoted as *feature files* and *test files*, respectively. For each file list, we extract all file names and the identifiers (e.g., classes, functions, attributes) in these files through static analysis, producing a set of entities. We refer to the entity sets extracted from feature files and test files as F and T , respectively. Next, we apply both *feature_patch* and *test_patch* to the codebase and repeat the same process to obtain the post-patch entity sets, denoted F' and T' . We compute the differences between the pre- and post-patch sets to obtain the sets of newly introduced entities in feature files and test files, respectively, i.e., $\Delta_F = F' - F$ and $\Delta_T = T' - T$. We then compute the intersection of Δ_F and Δ_T to isolate entities that are newly introduced and referenced in both test and

1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030

Table 6: Guidelines for the accuracy of the evaluation.

Score	Accuracy of the Evaluation
0	The tests perfectly cover all aspects of the new feature mentioned in the doc changes.
1	The tests cover the majority of aspects of the feature involved in doc changes, but may overlook some special cases.
2	The tests are too narrow/broad, only test a few special cases, or test content far beyond the description in the document.
3	The tests look for something different from what the doc changes is about.

implementation code. This step ensures we only consider entities relevant to the verification process. Finally, we remove all entities that are already mentioned in the documentation changes, yielding the final set of undocumented but test-referenced entity names.

A.2 NoCode-bench Verified

A.2.1 Sampling and Annotation Protocol

To ensure project diversity, we sample up to 35 task instances per project (or fewer if insufficient candidates exist), yielding a total of 238 candidate tasks. Considering the subjectivity of human judgment, each sample is annotated independently by two different annotators. For quality assurance, we conservatively select the more stringent label when disagreements arise, following SWE-bench Verified. As a result, each annotator evaluated approximately 95 task instances. We construct NoCode-bench Verified by randomly sampling candidate instances from NoCode-bench Full. To ensure diversity across projects, we cap the number of sampled instances at 35 per project (or fewer if insufficient candidates exist), resulting in a total of 238 candidate tasks.

Each candidate task is independently annotated by two annotators. All annotators have at least three years of Python development experience and hold a bachelor’s degree or higher in a relevant field. In cases of disagreement, we conservatively adopt the stricter label, following the protocol of SWE-bench Verified (Chowdhury et al., 2024). Under this setup, each annotator evaluates approximately 95 tasks.

A.2.2 Annotation Rubric

Clarity of the Task. Annotators assess whether the documentation change alone provides a suf-

Table 7: Human verification results for task clarity and test accuracy. For the *Other Problem* column, a score of 0 indicates no issues were identified, while a score of 3 indicates the presence of significant issues.

Score	Task Clarity	Test Accuracy	Other Problem
0 (best)	59	68	220
1	99	94	–
2	59	64	–
3 (worst)	18	9	18

ficiently clear and unambiguous specification of the intended feature. Annotators act as experienced software engineers with full access to the codebase and judge whether a reasonable implementation can be inferred from the documentation change.

Accuracy of the Evaluation. Annotators assess whether the test patch provides appropriate and targeted functional coverage of the behavior described in the documentation change. Unlike SWE-bench Verified, which focuses on whether tests can reliably evaluate diverse correct implementations (e.g., avoiding false negatives due to mismatches in naming or structure), we mitigate such concerns during benchmark construction in Phase 5. As a result, our evaluation emphasizes whether the provided test cases faithfully and comprehensively cover the functionality described in the documentation changes. Scoring is based on the criteria in Table 6.

Additionally, if the annotator believes that there are other major problems with the sample that have not been considered, i.e., any other reasons this sample should not be used in our setup for evaluating coding ability, the annotator needs to flag the sample and provide a basic explanation.

A.2.3 Additional Filtering Details

After all the samples have been annotated, we filter out any sample from the original test set where the severity of the clarity of the task or the accuracy of the evaluation is labeled as 2 or above. We also filter out all samples that have a major problem flagged by annotators. For example, in *pytest-3576*, the associated PR includes a large number of unrelated documentation changes, which are not indicative of genuine feature additions and thus do not meet the criteria of our benchmark. As summarized in Table 7, a total of 106 instances were removed due to unclear task descriptions or inaccurate evaluations, and 18 instances were excluded due to other major problems.

B DETAILED EXPERIMENT SETUP

B.1 Evaluation Metrics

- **File Matched Rate (File%)**: The percentage of modified files in the submitted patch that match one of the file paths edited by the ground truth patch.
- **Regression Tests Pass Rate (RT%)**: The proportion of tasks where all regression tests pass after applying the generated patch.
- **Feature Validation (FV)**: This metric aims to evaluate the completeness of the LLM’s implementation of a specific feature, which we measure based on the pass rate of F→P tests. Specifically, we calculate FV from both macro and micro perspectives.
 - **Micro-level FV (FV-Micro)**: The overall ratio of passed F→P tests to the total number of F→P tests across all instances: $FV\text{-Micro} = \frac{\sum_{i=1}^N \#Passed_i}{\sum_{i=1}^N \#Total_i}$, where $\#Passed_i$ and $\#Total_i$ denote the number of passed and total F→P tests for instance i , respectively.
 - **Macro-level FV (FV-Macro)**: The average F→P pass rate across instances: $FV\text{-Macro} = \frac{1}{N} \sum_{i=1}^N \frac{\#Passed_i}{\#Total_i}$, where N is the total number of instances.

B.2 Used Scaffolds

- **Agentless** adopts a hierarchical approach to sequentially identify the relevant files, classes or methods, and lines of code that require modification, and then generates patches based on the localized content with LLM. Following Guo et. al. (Guo et al., 2025b), we use Agentless-1.0 for evaluation.
- **OpenHands** is a widely adopted platform for building software development agents. We adapt the document changes in NoCode-bench as inputs to the OpenHands CodeAct Agent and utilize the standard evaluation docker images provided by NoCode-bench as a sandbox environment to ensure the correct execution of tasks.

C FURTHER ANALYSIS OF EXPERIMENT RESULTS

C.1 Detailed Analysis of NoCode-bench Verified Results

Performance Across Scaffolds. The Agentless scaffold achieves an average patch application rate

of 98.80%, whereas OpenHands produces applicable patches for only 67.98% of the tasks on average. For regression tests, 76.97% of the patches generated by Agentless passed all regression tests on average, compared to only 62.06% for OpenHands. In terms of full F2P test suite pass rate, Agentless achieves 8.74% (FV-Micro), with an average of 31.25% (FV-Macro) of F2P tests passed per instance, while OpenHands achieves 4.92% and 23.35%, respectively. The results show that Agentless outperforms OpenHands on NoCode-bench Verified.

Why Agentless Outperforms OpenHands.

Through manual inspection, we find that the reason lies in how OpenHands maintains a dialogue history to record interactions between the agent and the repository. Since solving tasks in NoCode-bench Verified requires editing an average of 2.39 files, the observed context often exceeds the model’s context window, resulting in a large number of incomplete or empty patches. In contrast, Agentless uses hierarchical localization to limit the length of retrieved context and applies post-processing to verify patch completeness, leading to better overall performance.

Gemini-2.5-Pro’s Complete Failure on OpenHands.

By inspecting the failure cases, we find that Gemini-2.5-Pro is almost unable to generate correct invocation formats for any of the tools used by OpenHands. This directly prevents it from accessing and editing the repository.

C.2 Detailed Analysis of NoCode-bench Full Results

Why Performance Drops on NoCode-bench Full.

Since NoCode-bench Full is a large-scale dataset involving more extensive file modifications and more complex code changes, it presents a diverse range of challenges. In addition, NoCode-bench Full may contain some noise, including unclear documentation changes or low-coverage regression tests, which make it more representative of real-world, NL feature addition tasks, but harder for LLMs to solve.

Model Robustness Comparison.

For DeepSeek-R1, DeepSeek-v3 and Qwen3-235B, success rates drop from 25.44% to 14.83%, 21.05% to 11.83%, and 13.16% to 6.62%, respectively, while GPT-4o only drops from 13.16% to 9.31%. Qwen3-235B and GPT-4o achieve similar performance

```

"""DeepSeek-v3 Patch"""
diff --git a/pylint/reporters/text.py b/pylint/reporters/text.py
--- a/pylint/reporters/text.py
+++ b/pylint/reporters/text.py
@@ -294,14 +295,15 @@ class ColorizedTextReporter(TextReporter):
    """Manage message of different types, and colorize output
    using ANSI escape codes.
    """
-
-   if msg.module not in self._modules:
-       msg_style = self._get_decoration("S")
-       if msg.module:
-           modsep = colorize_ansi(f"***** Module {msg.module}", msg_style)
-       else:
-           modsep = colorize_ansi(f"***** {msg.module}", msg_style)
-       self.writeln(modsep)
-       self._modules.add(msg.module)
+   if getattr(self.linter.config, "output_format", None) != "no-header":
+       if msg.module not in self._modules:
+           msg_style = self._get_decoration("S")
+           if msg.module:
+               modsep = colorize_ansi(f"***** Module {msg.module}", msg_style)
+           else:
+               modsep = colorize_ansi(f"***** {msg.module}", msg_style)
+           self.writeln(modsep)
+           self._modules.add(msg.module)
+   msg_style = self._get_decoration(msg.C)

"""Golden Patch"""
diff --git a/django/forms/models.py b/django/forms/models.py
--- a/django/forms/models.py
+++ b/django/forms/models.py
@@ -212,6 +213,18 @@ def _display(self, layout: Section) -> None:
    TextWriter().format(layout, self.out)
+class NoHeaderReporter(TextReporter):
+    """Reports messages and layouts in plain text without a module header."""
+
+    name = "no-header"
+
+    def handle_message(self, msg: Message) -> None:
+        """Write message(s) without module header."""
+        if msg.module not in self._modules:
+            self._modules.add(msg.module)
+        self.write_message(msg)
+class ParseableTextReporter(TextReporter):
+    """A reporter very similar to TextReporter, but display messages in a form
+    recognized by most text editors :

```

Figure 5: A failure example of `pylint-dev__pylint-7869` caused by the LLM’s lack of architectural understanding.

on NoCode-bench Verified, while GPT-4o outperforms Qwen3-235B on NoCode-bench Full. These results demonstrate that compared to DeepSeek-v3 and Qwen3-235B, GPT-4o performs more robustly in the scenarios that are closer to real-world conditions.

C.3 Case Study: Lacking Understanding of Existing Code Modules in `pylint-7869`

Figure 5 highlights a critical difference between the patch generated by DeepSeek and the golden patch for the task “`pylint-7869`”. This task requires introducing a new “no-header” option to Pylint. When this option is enabled, Pylint’s output will omit (i.e., not display) the header information that normally shows the name of the module being analyzed. DeepSeek-v3 directly modifies the member functions of the `ColorizedTextReporter` class to implement the “no-header” option, whereas the

golden patch introduces a new `NoHeaderReporter` class to achieve the same goal. Compared to the golden patch, DeepSeek-v3’s modification breaks the original functionality. Although it passes all F→P tests related to the new feature, it results in numerous regression test failures. This example demonstrates how LLMs may lack understanding of existing code modules needed to extend functionality without breaking existing features. Users may alleviate such phenomena by providing the LLM with information about the differences between existing modules or by guiding the LLM to analyze existing modules through a predefined plan.

C.4 Case Study: Tool-Calling Format Error in `astropy__astropy-9425`

Figure 6 shows a tool-calling dialogue from OpenHands while attempting to solve the task

1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238

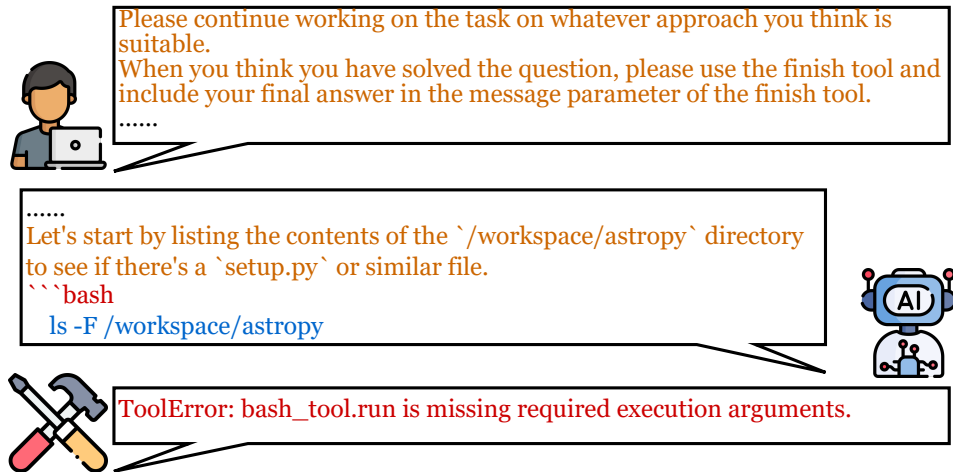


Figure 6: A failure example of *astropy__astropy-9425* caused by the LLM's incorrect format of tool calling.

1239 *astropy__astropy-9425*. In the first round, the agent
 1240 chooses to use the “bash” tool to inspect the reposi-
 1241 tory structure. However, since the tool requires the
 1242 LLM to return commands in JSON format, Gemini
 1243 fails the call by responding with the intended
 1244 command in plain text instead. As a result, the tool
 1245 invocation fails, and the LLM repeatedly retries
 1246 the same tool call. Eventually, the conversation
 1247 history fills the context window, preventing Open-
 1248 Hands from successfully generating a patch for this
 1249 instance. This example illustrates how formatting
 1250 errors in tool calls can cascade into complete task
 1251 failure.