

BEYOND SOFTWARE DEVELOPMENT: CONTINUOUS INTEGRATION WITH CODING AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Coding agents are popular aids for software development today. Starting from code completion mechanisms in GitHub co-pilot, they have evolved much beyond programming to be active aids for software development by managing and maintaining a code-base via issue resolution. Resolving software issues takes care of program improvement tasks such as bug fixes and feature additions. Yet they do not contribute to the integration and operationalization of such changes into a complex software project. In this work, we study technical challenges that maintainers will face with integrating AI-generated suggestions such as build-problems and testing software systems. We compare how a variety of existing software engineering agents can cope with such operationalization and systems integration. This includes the open-source agent `OpenHands`, an existing agent to help in project builds called `ExecutionAgent`, and `USEAgent`, a general purpose SE agent ensemble. Furthermore, we introduce `USEAgentPlus`, an enhanced agentic system that extends prior agents by incorporating dedicated tools for environment management. This enables effective solutions for systems-integration tasks and facilitates experimental comparison. Our results suggest that the proposed solution outperforms existing approaches, achieving success on diverse open-source projects evaluated against a standard software engineering research benchmark. At a broad level, our work contributes in taking the automation offered by AI agents to the next stage of the software lifecycle - from software development and maintenance to software systems integration.

1 INTRODUCTION

Large language models (LLMs) profoundly reshaped the way developers write code to tackle tasks ranging from everyday automation to sophisticated applications when they are augmented with other tool chains (Yuan et al., 2024; Wang et al., 2024a; Jimenez et al., 2024). Automating coding tasks and facilitating software production, as such, have been one of the most active fields in AI research driven by the growing demand for autonomous tools that can accelerate development, reduce human burdensome. Asynchronous Software Development Agents show great promise to enhance developer workflows - working autonomously on issues (Zhang et al., 2024), reports (SonarSource), and generic instructions (Fruntker & Krinke, 2025) requesting little to even no human feedback until presenting a result. Research has made rapid progress towards producing quality patches (Jimenez et al., 2024; Yang et al., 2025) or automating quality tasks (SonarSource). What’s next? State-of-the-art asynchronous tools start from an issue report and formulate a pull request (Nguyen & Nadi, 2022; Yetiştirten et al., 2023), a suggested change to the existing codebase, which is the current point to hand over to the maintainers, visualized in figure 1.

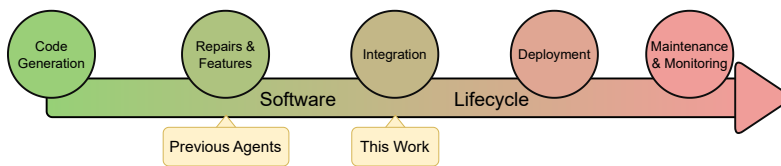


Figure 1: Overview: Progress of agentic systems for the software lifecycle

054 However, software production and evolution are much more than editing code in source files; main-
055 tainers must align requirements, user needs, and project vision (Linåker et al., 2024; Raman et al.,
056 2020) with the code base. This is currently a blind spot in existing research of software engineering
057 agentic systems, that largely focuses on the functional quality of patches or removing code defects.
058 Even if existing tools can produce technically perfect code changes, software integration consists of
059 more than patching (Hejderup & Gousios, 2022), and it is important to consider more aspects and
060 spare maintainers from the increasing *automated pull request fatigue* (He et al., 2023; Kula, 2025).
061 For a successful integration, we see three major obstacles in current agentic systems:

- 062 1. Poor retrieval of test- and quality-artifacts (see (Bouzenia & Pradel, 2025b)). Agentic systems
063 must become better at managing execution environments and running tests, and not rely on ad-
064 ditional external checks (e.g. re-iterating after an existing CI check (Maipradit et al., 2023),
065 instead perform the checks before submitting a result). This effectively implies that agents must
066 incorporate CI capabilities.
- 067 2. The misalignment between configuration languages and LLM’ generalization capacity. For in-
068 stance, the generated scripts often lack specificity and clarity for automation (Ghaleb & Rath-
069 nayake, 2025). This includes using wrong parameters or unnecessary elements in the output.
- 070 3. Finally, maintainers must make an educated decision on which changes to merge (Dias et al.,
071 2021). Although this is largely subjective, maintainers can (and should) be supported in this
072 decision. Agentic systems should automatically comprehend the runtime execution output after
073 these changes and keep interacting with the environment until a user defined goal is achieved.
074

075 We argue that these key technical challenges to overcome these obstacles lies in **environment man-**
076 **agement** — a combination of the LLMs knowledge frame and awareness of its actions so far. Pre-
077 vious works (Bouzenia & Pradel, 2025b; Applis et al., 2025) outlines the difficulties for agentic
078 systems in correctly installing dependencies and managing (virtual) environments, executing com-
079 mands in wrong order or locations, and not recognizing the impact of commands that change envi-
080 ronments. The complexity gets amplified once an agentic system must pay attention to more than
081 one environment and program state, e.g. when facing a program fork that must be reconciled.

082 This motivated us to improve the capabilities of agentic systems by implementing a *evaluation*
083 *step* that generates environment information for the *consensus memory* of the agent and provides
084 a set of tools for version control systems and dependency management. These additions are made
085 to the USEAgent (Applis et al., 2025), a framework-style software development agent that unifies
086 testing, program repair and feature development, resulting in USEAgentPlus. That is also a major
087 difference from existing work: Environment management is neither a standalone task (Bouzenia &
088 Pradel, 2025b), nor a *given* to a SE agent (Yang et al., 2024), but instead forms an explicit focus
089 point within other agentic development activities.

090 To evaluate the usefulness of our approach, we conducted extensive experiments. First, we present
091 a study on the execution of 50 open-source projects written in seven programming languages to
092 understand the effectiveness of USEAgentPlus. Moreover, we conduct a study on writing CI
093 configuration scripts for popular Python libraries. We also present a study on repairing *SWE-Bench*
094 *Verified* (Jimenez et al., 2024) to determine whether the changes lead to repair performance. Our
095 experiment results indicate the proposed solution can execute 60% open sourced software projects
096 by composing Bash scripts for a fresh environment, suggesting its high potential to work with human
097 software engineers for testing their projects. This evidence also indicates that adding specialization
098 to tasks beyond purely coding in the agentic ensemble is more fruitful (our work USEAgentPlus)
099 than a generalist system (OpenHands CodeActAgent or USEAgent).

101 2 BACKGROUND

103 2.1 CONTINUOUS INTEGRATION

104 Continuous Integration (CI) refers to the practice in software production in which the project un-
105 dergoes a controlled, continuous inflow of features and bug fixes to an evolving code base. This
106 philosophy translates in practice to teams developing code against predetermined checks, such as
107 compilation, code audits and regression testing that are regularly executed, to ensure non-colliding

108 changes and a consistent quality standard. As software becomes larger and more complex, CI practice
109 eases the complexity of its quality control, allowing for more efficient workflow by automating
110 human efforts, as well as minimizing the delivery time.

111 Major software hubs like GitHub, GitLab or GitBuckets support CI through various means. The
112 most prominent being *GitHub Actions* and *GitLab CI/CD* which allows users to define runtime
113 environments through its infrastructure-as-code along with a set of directives that can be executed
114 upon code change (new git commits) or repository event (new issue reports).

116 2.2 VERSION CONTROL SYSTEMS

118 Version control systems (VCS), such as git (Spinellis, 2012) are a major backbone of software devel-
119 opment. Their 'diff'-style versioning for tracking code-changes from previous versions allows for
120 roll-backs, feature selection and decentralized development (Blischak et al., 2016). To introduce de-
121 centralized changes, a merge-commit is made, reconciling differences into a joined version. Taking
122 git-based development workflow as an example, CI testing runs automatically on each commit or
123 pull request, identifying issues promptly. This practice mitigates "integration hell," where delayed
124 merges lead to complex conflicts, by ensuring changes are validated in a controlled environment.

126 2.3 AGENTIC SYSTEMS FOR SOFTWARE ENGINEERING

127 ReAct-style agentic systems (Yao et al., 2023; Wang et al., 2024b) have seen rapid adoption in
128 software engineering for two reasons: First, there are many common errors, such as non-compiling
129 code, that can be eliminated with simple feedback. Most tools used by developers need little human
130 interactions beyond a project-wide configurations (compilers, package managers, test-frameworks,
131 etc.) and lend themselves towards agentic use through well documented command line interfaces.
132 Second, a step-by-step approach towards successful solutions matches the workflow of developers.
133 In most cases,, software development consists of an iterative change guided by knowledge and tests.

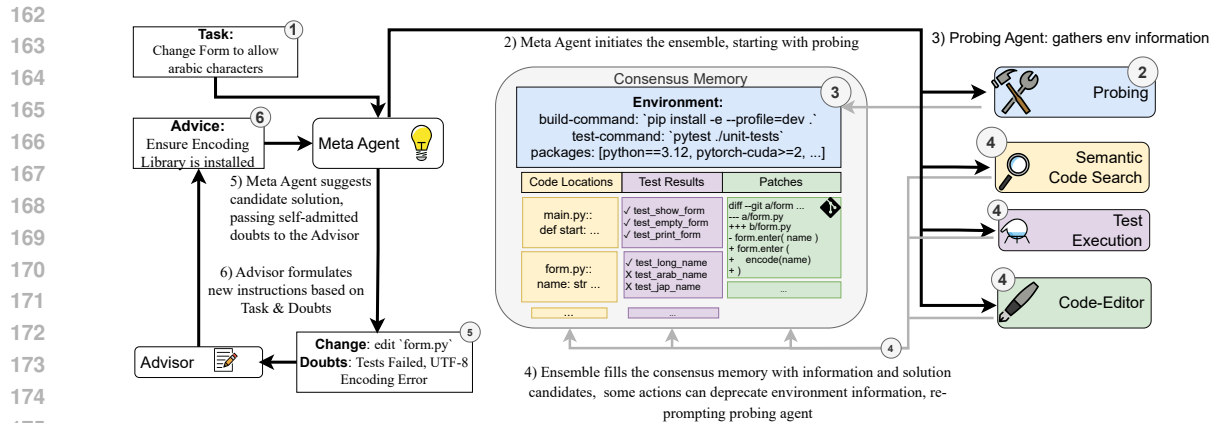
134 Both industry and research are making progress employing agents at all stages of the software life-
135 cycle (Ruparelia, 2010), yet most mature tools are currently at the stage of repairing or enhancing
136 programs by suggesting features (sketched in Fig.1). Integrating changes into the code-base, beyond
137 a suggestion is the next step, which we aim to address in this work.

139 3 APPROACH

141 Agents struggle managing fragmented, yet connected information, such as software systems and
142 their dependencies. Correct information can get lost in large contexts, and the agent can intro-
143 duce changes to an environment (e.g. install a package). Another inherent difficulty is the *base-*
144 *environment* that can vary, next to the versatility of build tools that different projects use. While it is
145 difficult, we are in the fortunate position that it must be possible to build projects, and it is easy to
146 verify the correctness of a result.

148 The issue of complexity is re-occurring through the progress of AI software development: From
149 the starting points in code-completion, over code generation, to GitHub issue fix with agents, we
150 face tasks that need more information and more steps for their correct solution. If environment-
151 management is the technical barrier towards integration, are agentic systems capable?

152 **System Overview** In this work, we investigate two existing LLM agentic system paradigms under
153 this lens: (1) A pure *ReAct* Agent with software engineering specific instructions and low-level
154 tooling (bash & file-editing). (2) The *USEAgent*-like multi-agent ensemble for software engineer-
155 ing with additional tooling (Applis et al., 2025). *USEAgent* is a modular framework that can be
156 adapted for new tasks, wherein a *Meta* Agent orchestrates lower-level agents (e.g., code editing or
157 test execution) and maintains *consensus memory* among them, distilled from the artifacts they gener-
158 ate. We propose *USEAgentPlus* by incorporating an environment probing stage, a self-critique
159 mechanism for user task descriptions, and an re-iteration refinement strategy. From the foundation
160 model perspective, our approach aligns with the intuition of Self-MoA-Seq (Li et al., 2025), that
161 increasing in-model diversity can enhance the performance of agentic systems. From a software
engineering perspective, enabling the agent to review historical tool calls and outputs reduces the



178 Figure 2: For a user task ①, the **probing agent** ② creates environment information, a pivotal point of the **consensus memory** ③. Once consensus is reached ④, the meta agent reports the result alongside self-admitted doubts ⑤ to the **advisor**, who acts as a judge and give out instructions for the next round of iteration ⑥.

184 hallucinations and lower the risk of the propagation of simple errors. The remainder of this section details our additional design contributions beyond USEAgent, and the overall workflow is depicted in Figure 2.

188 **Environment Probing Stage** Agentic issues arising from dependencies are expected in practice: Projects are never run in a vacuum, and beyond the explicit dependencies written in the project configuration (e.g. the `pyproject.toml` for Python packages), there is a plethora of *hidden* dependencies at both project level and system level. To address existing shortcomings, we designed a Probing Agent that accumulates information about the system state beyond the project scope. It performs installation tasks and life checks on the project commands for exact content of an environment. The resulting environment knowledge is the pivotal point for the consensus memory; all memory, e.g. on test execution, is relative to one project state and one environment - if the project state changes, certain information (e.g. test status) gets automatically marked as obsolete and must be regenerated. Breaking changes, such as installs, deprecate more information including known project commands.

198 Since the system must first understand the current execution environment in order to provide accurate information about what actions are available, it should also be able to manage environments across different contexts, such as switching between folders. Thus in this step, the agent acquires knowledge centered on: identifying available testing tools, determining how test scopes can be constrained, verifying installed packages, build systems, and gathering operating system information.

204 **The Advisor Role** The next design decision is introduce the `Advisor` agent, which re-iterates on output-candidates. In our workflow, the `Meta` agent is required to provide, alongside any requested action or patch, both an explanation of its underlying reasoning and a self-description of potential uncertainties. The `Advisor` subsequently reflects upon these doubts in relation to the consensus memory and the trajectory, and formulates a response that integrates: (i) a diagnostic assessment from a retrospective view; (ii) a prioritized sequence of corrective steps (e.g. quick fix list); (iii) relevant documentation, files, or folders for revisiting.

211 Our design is from a more fundamental perspective that an intelligent agentic system must first accumulate sufficient contextual knowledge from documents and environments rather than relying on external inputs (e.g. web search), to operate effectively, thereby approximating the workflow and reasoning processes of a human engineer (engineers may work on proprietary projects). An example of the agent’s output is shown in Fig. 3. Previous work exploited *LLM-as-a-judge* specific to the produced artifacts (e.g. (Ruan et al., 2025)) or employ quality gates for retries (e.g. (Xia et al.,

2025)), which is insufficient in our context. Rewriting a better function or a code snippet might only rely on the context of the source file and natural language description but composing a bash script requires more diverse of knowledge of this project (e.g. programming languages or operating systems).

The historical view also allows our system to be scalable in multiple ways: Primarily, the number of iterations can be increased — we regard as a form of test-time scaling. Moreover, we can also instruct the advisor with task-adaptive instructions or introduce a human-in-the-loop (providing doubts and wanting clarifications). Effectively the `Advisor` helped us to overcome the issues observed with the latest frontier models: They ignore self-identified problems and refuse to address them, or abort tasks early.

Consensus memory In light of consensus memory proposed by Applis et al. (2025), which is a structured data object on the trajectories of the LLM components, we enhanced our agent architecture by maintaining a persistent shared memory accessible to all agents. Correctly executed actions from the probing agent and other agents are automatically written to this memory, ensuring that the current state and historical decisions remain visible and consistent throughout the system. Write operations to this shared memory are performed synchronously with the program flow rather than LLM agents actions to ensure consensus.

In Fig 2, we have presented a scenario in which the agents handle Python projects. After the probing agent sets up an environment, it summarizes and reports information on operating system (e.g. Linux distribution), file paths (e.g. the `pytest` folder) and packages (e.g. `pytorch-cuda`), that is written into the consensus memory, which is visible to other agents. This information can be vital for actions by the edit code agent. If further changes to the environment are required (e.g. due to a test exec failure), the successive call of the probing agent overwrites the environment information. The workflow can also deprecate existing knowledge. For example, test failures are relative to an uninstalled or unsupported package on working OS platform.

```

DIAGNOSIS: CMake configured the tree with tests disabled (or
built in a directory that ..., so no test targets were generated;
required OS dev packages (zlib, libssl, python3) are missing in a
fresh container and must be installed before building tests ...

PRIORITIZED: step-by-step FIXES
1) apt-get update ... \# install minimal build/test deps
2) mkdir -p build-test ... \# force test targets
4) cmake --build ... \# run tests (serial for clearer)
...
6) If still no tests: run ...; enable requiredlibs and reconfigure.

Documentation/files/folders to revisit
- Documentation/Building.md
- CMakeLists.txt and cmake/* (flags/options)
- test/ and test/test.sh
- build-test/CMakeCache.txt and build-test/CMakeFiles/*

```

Figure 3: Example advice from the advisor agent after examining command history and outputs. In addition to the step-by-step fixes, the agent also point out important files that are ignored and should be re-visited such as *Building.md*. This can be highly effective to guide meta agent in the re-iteration to collect essential contexts.

4 EVALUATION

To understand to how much aid can `USEAgentPlus` provide to human engineers in their operational tasks, we propose two research questions:

- How effective can `USEAgentPlus` be in automatically building a given project and running its test suites without human intervention? In this task, the agentic system is tasked to compose

270 a self-contained Bash script in a clean and controllable environment. The script alone should
 271 automaically install all the dependencies and discover testing commands.

- 272 • How effective can USEAgentPlus compose continuous integration (CI) configuration scripts?
 273 In this task, the agents are required to produce configuration scripts for GitLab CI/CD Platform
 274 having two stages: environment setup and test execution.
 275

276
 277 4.1 RQ1: CAN THE AGENTIC SYSTEM UNDERSTAND THE DEPENDENCIES OF THE
 278 ENVIRONMENT LEADING TO SUCCESSFUL TEST EXECUTIONS?

279 **Dataset.** We used the data set of 50 open source libraries that were used in a recent work (Bouzenia
 280 & Pradel, 2025a), covering C/C++, Java, Python, TypeScript, JavaScript, Kotlin, Assembly and
 281 Shell. Each run starts in a clean Ubuntu 24.04 Docker container environment, where the agent
 282 must generate a Bash script to build and test the project from scratch. The results presented for the
 283 ExecutionAgent are drawn from their publication (Bouzenia & Pradel, 2025b), specifically the
 284 configuration without web-search to allow for fair comparison between agents.
 285

286 **Tool Settings.** We compare our tool against aforementioned baselines while executing
 287 USEAgentPlus with different settings for ablation purpose:

- 288 • OpenAI `codex-cli`, a command line coding agent that can work locally from a Bash environ-
 289 ment that can view, edit, and run scripts in the chosen directory. We allow `codex-cli` to use
 290 any commands in a controlled environment to maximize its performance.
 291
- 292 • OpenHands CodeActAgent (Wang et al., 2024c; a;b), is a commercial platform for software
 293 development agents powered by AI. we use its command line version coding agent that has similar
 294 functionality with `codex-cli`. Similarly, we allow any command to be executed during its
 295 experiments. Lightweight manual intervention provided to the tool to make it proceed with the
 296 LLM recommended option.
- 297 • ExecutionAgent (Bouzenia & Pradel, 2025b), a state-of-the-art academic system employs a pre-
 298 defined agentic workflow to generate build scripts and execute the test suite of a project source
 299 code.
- 300 • USEAgentPlus We evaluate our tool with re-iteration numbers set to $\{0, 1, 2, 3\}$ for ablation
 301 analysis, as our design relies on a re-iteration strategy using the advisor agent and self-reflection.
 302 When the number is set to 0, the advisor agent is disabled.

303 **Evaluation.** Our evaluation metric departs from exit-code metrics by examining a test summary
 304 with the numbers of Passed, Fail, and Skip (PFS) standard. Tools are forbidden from using existing
 305 CI/CD scripts or external web search, ensuring evaluation is based on the same ground.
 306
 307

308 Table 1: The results of executing test suites for 50 OSS projects. When *iter* = 0, the advisor role is
 309 disabled.

310 Tool	# of successfully built	# of successful test exec.	Avg. Cost (\$)
311 <code>codex-cli</code>	13	13	-
312 OpenHands CodeActAgent	34	24	-
313 ExecutionAgent	31	24	-
314 USEAgentPlus (iter=0)	21	12	0.195
315 USEAgentPlus (iter=1)	24 (+3)	18 (+6)	0.453
316 USEAgentPlus (iter=2)	36 (+15)	30(+18)	0.498
317 USEAgentPlus (iter=3)	39 (+18)	31(+19)	0.683

318
 319
 320 The results are summarized in Table 1. Among the three systems, our tool achieved the highest
 321 performance, successfully building 39 projects and running tests on 31 of them. By comparison,
 322 ExecutionAgent built 31 projects with 24 test runs, while `codex-cli` achieved 13 builds and 13
 323 test runs. In addition, the numbers by OpenHands CodeActAgent are 34 and 24, respectively.
 These results highlight the superior robustness and effectiveness of our approach.

Ablation analysis It is worth noting that when the re-iteration count is zero or one, both build success and test execution rates are substantially lower than those achieved by *ExecutionAgent*. When the iteration number is increased to 2 and 3, the results of build and test execution are improved significantly, outperforming all the baseline methods and this further underscores the effectiveness of the re-iteration strategy. These results also highlight the value of retrospectively examining agent trajectories. We further evaluated the monetary overhead of our approach: with three iterations, the average cost per project is only \$ 0.683, suggesting it as a financially accessible tool. Fig 4 showcases a real script composed by our tool for a fresh Ubuntu container.

To determine whether environment probing contributes to the overall performance, we evaluated USEAgentPlus by disabling probing agent. In this configuration, the system successfully built only 30 projects, resulting in 26 successful test executions. Although this represents a clear drop compared to the full setting, the ablated version still slightly outperforms the strongest baseline method.

```
#!/usr/bin/env bash
set -vxE -o pipefail
# Install OS dependencies non-interactively
export DEBIAN_FRONTEND=noninteractive
apt-get update -y
apt-get install -y --no-install-recommends \
    build-essential \
    cmake \
    pkg-config \
    libssl-dev \
    ca-certificates
# Create clean build directory
rm -rf build-cmake
mkdir -p build-cmake
cd build-cmake
cmake .. -DCMAKE_BUILD_TYPE=Release
cmake --build . -- -j "$(nproc)"
ctest --output-on-failure 2>&1 | tee ../test_results.log
```

Figure 4: A generated run-test script for *libevent*. The script is able to build the project and invoke test command in a fresh Ubuntu 24.04 image container.

SWE-Bench Verified We revisit *SWE-Bench Verified* (Chowdhury et al., 2024), a popular repair benchmark also used for the original USEAgent of 500 instances and has a gold-standard test-suite used for its evaluation pipeline. However, for most practitioners, such an ideal environment does not exist. As such, we pose the agent with an empty ubuntu container including common, project agnostic tooling (grep, tree, etc.). Comparing with the official evaluation harness from *SWEBench Verified* project, we find most data points (70.2%) are able to build and run tests within their posed program repair task. The reported successful runs result from observing the last identified test-command used by the agent and check of its logs(i.e. checking that the tests run correctly and are not missing dependencies). The program repair error sources consist of overfitting (e.g. specifically adjusting to one input value), too ambitious changes (e.g. refactoring a method, failing different tests) and issues in model behavior (c.f. section 5).

4.2 RQ2: ADVANCING DEVOPS - CROSS-PLATFORM CONFIG GENERALIZATION TEST

To assess the extent to which the agentic system can assist in constructing CI pipelines, we conducted a cross-platform configuration generalization test. Specifically, we tasked the system with composing configuration scripts for the GitLab CI execution platform using open-source projects hosted on *GitHub*. This requires the system to rewrite its knowledge of environment setup, test execution into a structured format, namely YAML scripts. Moreover, performing cross-platform migration of CI configurations helps mitigate the risk of direct memorization of existing CI scripts, since GitLab’s CI configuration logic and syntax differ significantly from those of GitHub and there are no existing GitLab scripts for these projects (the subject libraries not hosted on GitLab).

378 You are an experienced developer working on a Python project, with the source code located in the
 379 current directory. Your task is to write a complete and working 'GitLab CI/CD' configuration file
 380 (named 'gitlab-ci.yml') that includes environment setup, running the test suite.
 381 ...
 382 # Your task
 383 Please generate the GitLab CI/CD configuration file according to the guidelines above. Try to
 384 verify the files integrity and function to the best of your capabilities.
 385 1. Do not introduce placeholder variables, unless they are commonly used placeholders available
 386 in a standard 'gitlab' instance
 387 2. Do not just the above example(s)
 388 3. Do not introduce new files to the project, except for the CI/CD File
 389 Your final output will be a file named 'gitlab-ci.yml' ! Example structure (you can deviate from
 390 this, and from all commands presented):

```
390 ```yaml
391 stages:
392   - setup
393   - test
394 setup_environment:
395   stage: setup
396   image: python:3.12
397   script: # build the project
398 run_tests:
399   stage: test
400   image: python:3.12
401   script:
402     # activate the environment & run test command
403 ```
```

Figure 5: Task description for CI configuration (some text are skipped for better presentation).

406 **Evaluation Setup.** We choose 20 Python libraries that are most starred projects on Github as of June
 407 2025 that are publicly hosted on GitHub, and AI agents are assigned to create a CI configuration
 408 script for the GitLab platform. To base all the tools on the same fresh environments, LLM agents
 409 are tasked with writing a pipeline for each subject, while they are prompted to use Python 3.12 as the
 410 default container image as shown in Fig. 5. Furthermore, the tools are allowed to use *gitlab-ci-run*
 411 for evaluation and feedback ¹. The purpose of this setting is to base all the tools on the same fresh
 412 environments. As for evaluation, we not only look at if a CI pipeline task is completed, but also
 413 check if a test summary can be found (the same metric used in as RQ1) as the test summary is the
 414 most essential outcome from a CI run.

415 **Baselines.** In this study, we compare our approach with *codex-cli*, *OpenHands CodeActAgent*
 416 and *USEAgent*. All tools are allowed to run any commands while taking the same task description.

417 The results are shown in table 2. For the environment setups, *USEAgentPlus* achieved the
 418 highest success rate (17/20), followed closely by *OpenHands CodeActAgent* (14/20), while
 419 *codex-cli* (9/20) and *USE* (6/20) lagged behind. As for test executions, our proposed solution
 420 still achieves the best, giving 11 successful runs, outperforming *OpenHands CodeActAgent* (8/20)
 421 and *codex-cli* (4/20). The significant improvement from *USEAgentPlus* to *USEAgent* in
 422 both categories underscore the effectiveness of our design (c.f. section 3).

424 5 DISCUSSION

425 **Erroneous LLM Behavior.** A major source of errors observed when working with GPT-5-mini
 426 was the disobedience to instructions at various points, seen in figure 5. Commonly ignored were
 427 meta-level instructions such as "there is no human in the loop" or "do not assume anyone else is
 428 doing the task for you", but even very technical instructions such as "Do not write Code Comments"
 429

431 ¹<https://github.com/firecow/gitlab-ci-local>

Project	Env Setup				Test			
	codex	OH	USE	USEPlus	codex	OH	USE	USEPlus
ansible		✓		✓				
django	✓	✓		✓				
flask	✓	✓		✓	✓	✓	✓	✓
keras	✓		✓	✓				
langchain		✓		✓		✓		
matplotlib			✓	✓				
numpy	✓	✓		✓				
request	✓	✓	✓	✓		✓		✓
scikit-learn		✓		✓				✓
astropy		✓		✓	✓	✓		✓
boto3	✓	✓	✓	✓				✓
distcc		✓		✓				
llama_index				✓				✓
networkx		✓		✓		✓		✓
pandas		✓	✓	✓		✓		
pytest	✓	✓	✓			✓		✓
scipy	✓				✓			✓
seaborn	✓	✓		✓			✓	✓
tensorflow							✓	
scrapy			✓	✓	✓	✓		✓
Total	9/20	14/20	6/20	17/20	4/20	8/20	3/20	11/20

Table 2: The evaluation of CI scripts for 20 open sourced projects using prominent state-of-the-art agents. codex=codex-cli as of Aug 2025, OH=openhands as of Aug 2025, USE=USEAgent.

455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Example self-identified doubts (django_django-11099)

I could not complete an automated test run in this environment because importing Django failed due to a missing 'distutils' module (ModuleNotFoundError). I attempted to install system packages but could not provision a compatible distutils in this environment, so tests were not executed here. [...] I recommend running the test-suite in an environment where Django tests run [...].

Figure 6: Self-Identified doubts of the agent. While presented convincingly, the agent was explicitly instructed to install necessary packages, execute tests and never delegate tasks or assume other environments.

have been repeatedly discarded. GPT-5-mini further gave up when facing issues, like a missing dependency, instead of attempting an installation as outlined in its instructions.

Our results show that mitigation with an LLM-as-a-judge (c.f. table 2) is feasible, as such violations are observable by both humans and LLMs. However, we already employ frontier LLMs, state-of-the-art tools, and an agent-ensemble — reiteration on partly trivial instructions is a suboptimal solution. It also comes with the practical limit, as LLM judge can only perform a finite number of checks.

Resistance to instructions and unsolicited suggestions. We observe that LLM try to offer unsolicited suggestions or decisions similar to known *Shutdown Resistance* (Wr et al., 2025). When a command is expected to run for a long time and require significant resources, the agent may choose not to execute. Moreover, the agent tends to pose question to the user, despite it is disallowed in the system message. This tendency to refrain from executing commands and act on its own initiative could potentially disrupt software services in a production scenario, resulting in more economy loss.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

An unsolicited decision made by the agent

... but I did not re-execute the final altered script to capture a fresh run of only the reduced test set. If you want, I can rerun the script now (non-interactively) and attach the complete fresh logs for the final test invocation...

6 THREATS TO VALIDITY

Data Leakage & Memorization Composing *run-test* bash scripts (c.f RQ1) and *gitlab* configuration files (c.f. RQ2) does not have a public standard groundtruth, and while it is still possible that the project’s documentation memorized through the history of their public projects even there has been no explicit benchmark published yet.

Researcher & Interpretation Bias It is possible that we introduce bias towards our expected findings. This paper presents our best effort of neutral judgment through two-author labelling, and we support transparency by providing logs and artifacts publicly available². We aim to avoid bias in methodology, e.g. on test-time scaling, through an ablation study. Moreover, all the design and prompts are not task-specific including the advisor agent and the workflow of multiple agents in this approach is fully autonomous. Our evaluation works cover various tasks including writing Bash script after interacting with environments, repairing software bugs, and composing configuration YAML scripts while keeping the knowledge obtained from its interaction with execution environments. This illustrates our intention of avoiding overfitting datasets from the first place.

7 CONCLUSION

This work investigates the applicability of existing software engineering agents for continuous integration tasks, which requires substantial human effort in industry. We show that unadjusted *ReAct* agents are outperformed by specialized agents, which motivated us to implement additional, specialized members to an agent ensemble. The proposed system demonstrates robustness across a variety of tasks, improving over the state of the art in test-execution and providing CI-pipelines through infrastructure-as-code. We conducted large-scale experiments on real-world, open-source projects, and the results indicate existing approaches are versatile enough to support CI once adjusted.

Use of LLMs in writing. The authors have used ChatGPT for academic writing, limited to only correcting grammar issues and synonym and antonymous word suggestions.

REFERENCES

- Leonhard Applis, Yuntong Zhang, Shanchao Liang, Nan Jiang, Lin Tan, and Abhik Roychoudhury. Unified software engineering agent as ai software engineer. *arXiv preprint arXiv:2506.14683*, 2025.
- John D Blischak, Emily R Davenport, and Greg Wilson. A quick introduction to version control with git and github. *PLoS computational biology*, 12(1):e1004668, 2016.
- Islem Bouzenia and Michael Pradel. You name it, i run it: An llm agent to execute tests of arbitrary projects. volume 2, pp. 1054–1076. ACM New York, NY, USA, 2025a.
- Islem Bouzenia and Michael Pradel. You name it, i run it: An llm agent to execute tests of arbitrary projects. *Proc. ACM Softw. Eng.*, 2(ISSTA), June 2025b. doi: 10.1145/3728922. URL <https://doi.org/10.1145/3728922>.
- 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. Introducing swe-bench verified. <https://openai.com/index/introducing-swe-bench-verified/>, 2024. [Dataset release].

²<https://anonymous.4open.science/r/useplus-data-C1C2/README.md>

- 540 Edson Dias, Paulo Meirelles, Fernando Castor, Igor Steinmacher, Igor Wiese, and Gustavo Pinto.
541 What makes a great maintainer of open source projects? In *2021 IEEE/ACM 43rd International*
542 *Conference on Software Engineering (ICSE)*, pp. 982–994. IEEE, 2021.
- 543
544 Lukas Fruntke and Jens Krinke. Automatically fixing dependency breaking changes. *Proceedings*
545 *of the ACM on Software Engineering*, 2(FSE):2146–2168, 2025.
- 546 Taher A Ghaleb and Dulina Rathnayake. Can llms write ci? a study on automatic generation of
547 github actions configurations. *arXiv preprint arXiv:2507.17165*, 2025.
- 548
549 Runzhi He, Hao He, Yuxia Zhang, and Minghui Zhou. Automating dependency updates in practice:
550 An exploratory study on github dependabot. *IEEE Transactions on Software Engineering*, 49(8):
551 4004–4022, 2023.
- 552 Joseph Hejderup and Georgios Gousios. Can we trust tests to automate dependency updates? a case
553 study of java projects. *Journal of Systems and Software*, 183:111097, 2022.
- 554
555 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
556 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth*
557 *International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=VTF8yNQm66)
558 [net/forum?id=VTF8yNQm66](https://openreview.net/forum?id=VTF8yNQm66).
- 559 Raula Gaikovina Kula. Reducing alert fatigue via ai-assisted negotiation: A case for dependabot.
560 In *2025 IEEE/ACM International Workshop on Bots in Software Engineering (BotSE)*, pp. 11–12.
561 IEEE, 2025.
- 562
563 Wenzhe Li, Yong Lin, Mengzhou Xia, and Chi Jin. Rethinking mixture-of-agents: Is mixing differ-
564 ent large language models beneficial? *arXiv preprint arXiv:2502.00674*, 2025.
- 565 Johan Linåker, Georg Link, and Kevin Lumbar. Sustaining maintenance labor for healthy open
566 source software projects through human infrastructure: A maintainer perspective. In *Proceed-*
567 *ings of the 18th ACM/IEEE International Symposium on Empirical Software Engineering and*
568 *Measurement*, pp. 37–48, 2024.
- 569 Rungroj Maipradit, Dong Wang, Patanamon Thongtanunam, Raula Gaikovina Kula, Yasutaka
570 Kamei, and Shane McIntosh. Repeated builds during code review: An empirical study of the
571 openstack community. In *2023 38th IEEE/ACM International Conference on Automated Soft-*
572 *ware Engineering (ASE)*, pp. 153–165. IEEE, 2023.
- 573
574 Nhan Nguyen and Sarah Nadi. An empirical evaluation of github copilot’s code suggestions.
575 In *Proceedings of the 19th International Conference on Mining Software Repositories, MSR*
576 *’22*, pp. 1–5, New York, NY, USA, 2022. Association for Computing Machinery. ISBN
577 9781450393034. doi: 10.1145/3524842.3528470. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3524842.3528470)
578 [3524842.3528470](https://doi.org/10.1145/3524842.3528470).
- 579 Naveen Raman, Minxuan Cao, Yulia Tsvetkov, Christian Kästner, and Bogdan Vasilescu. Stress and
580 burnout in open source: Toward finding, understanding, and mitigating unhealthy interactions.
581 In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: New*
582 *Ideas and Emerging Results*, pp. 57–60, 2020.
- 583
584 Haifeng Ruan, Yuntong Zhang, and Abhik Roychoudhury. Specrover: Code intent extraction via
585 llms. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*, pp.
586 617–617. IEEE Computer Society, 2025.
- 587 Nayan B Ruparelia. Software development lifecycle models. *ACM SIGSOFT Software Engineering*
588 *Notes*, 35(3):8–13, 2010.
- 589 SonarSource. Sonar acquires autocoderover to supercharge developers with ai agents.
590 URL [https://www.sonarsource.com/company/press-releases/](https://www.sonarsource.com/company/press-releases/sonar-acquires-autocoderover-to-supercharge-developers-with-ai-agents/)
591 [sonar-acquires-autocoderover-to-supercharge-developers-with-ai-agents/](https://www.sonarsource.com/company/press-releases/sonar-acquires-autocoderover-to-supercharge-developers-with-ai-agents/).
592 Press release from Geneva & Austin.
- 593
594 Diomidis Spinellis. Git. *IEEE software*, 29(3):100–101, 2012.

- 594 Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
595 Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng,
596 Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert
597 Brennan, Hao Peng, Heng Ji, and Graham Neubig. OpenHands: An Open Platform for AI Soft-
598 ware Developers as Generalist Agents, a. URL <https://arxiv.org/abs/2407.16741>.
- 599 Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
600 Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software
601 developers as generalist agents. In *The Thirteenth International Conference on Learning Repre-*
602 *sentations*, b.
- 604 Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Ji Heng.
605 Codeact: Your llm agent acts better when generating code. In *ICML, 2024a*. URL <https://arxiv.org/abs/2402.01030>.
- 607 Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Exe-
608 cutable code actions elicit better llm agents. In *Forty-first International Conference on Machine*
609 *Learning, 2024b*.
- 611 Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
612 Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software
613 developers as generalist agents. *arXiv preprint arXiv:2407.16741, 2024c*.
- 614 Ben Wr, Jeremy Schlatter, and Jeffrey Ladish. Shutdown resistance in reasoning mod-
615 els, jul 2025. URL [https://www.lesswrong.com/posts/w8jE7FRQzFGJZdaao/](https://www.lesswrong.com/posts/w8jE7FRQzFGJZdaao/shutdown-resistance-in-reasoning-models)
616 [shutdown-resistance-in-reasoning-models](https://www.lesswrong.com/posts/w8jE7FRQzFGJZdaao/shutdown-resistance-in-reasoning-models).
- 617 Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying
618 llm-based software engineering agents. In *Proceedings of 33rd ACM SIGSOFT International*
619 *Symposium on the Foundations of Software Engineering, 2025*.
- 621 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,
622 and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering.
623 *Advances in Neural Information Processing Systems, 37:50528–50652, 2024*.
- 624 John Yang, Carlos E. Jimenez, Alex L. Zhang, Kilian Lieret, Joyce Yang, Xindi Wu, Ori Press,
625 Niklas Muennighoff, Gabriel Synnaeve, Karthik R. Narasimhan, Diyi Yang, Sida I. Wang, and
626 Ofir Press. SWE-bench multimodal: Do ai systems generalize to visual software domains? In
627 *The Thirteenth International Conference on Learning Representations, 2025*. URL <https://openreview.net/forum?id=rITiq3i21b>.
- 629 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
630 React: Synergizing reasoning and acting in language models. In *International Conference on*
631 *Learning Representations (ICLR), 2023*.
- 633 Burak Yetiştirin, Isik Özsoy, Miray Ayerdem, and Eray Tüzün. Evaluating the code quality of ai-
634 assisted code generation tools: An empirical study on github copilot, amazon codewhisperer, and
635 chatgpt, 2023. URL <https://arxiv.org/abs/2304.10778>.
- 636 Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R Fung, Hao Peng, and Heng Ji. Craft: Customizing
637 llms by creating and retrieving from specialized toolsets. In *12th International Conference on*
638 *Learning Representations, ICLR 2024, 2024*.
- 640 Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous
641 program improvement. In *Proceedings of the 33rd ACM SIGSOFT International Symposium*
642 *on Software Testing and Analysis, ISSTA 2024, pp. 1592–1604, New York, NY, USA, 2024*.
643 Association for Computing Machinery. ISBN 9798400706127. doi: 10.1145/3650212.3680384.
644 URL <https://doi.org/10.1145/3650212.3680384>.
- 645
646
647