# A.S.E: A Repository-Level Benchmark for Evaluating Security in AI-Generated Code

Keke Lian[1†], Bing Wang[2†], Lei Zhang[3], Libo Chen[4], Junjie Wang[5], Ziming Zhao[6],
Yujiu Yang[5], Haotong Duan[1], Haoran Zhao[3], Shuang Liao[3], Mingda Guo[3], Jiazheng Quan[2],
Yilu Zhong[2], Chenhao He[4], Zichuan Chen[4], Jie Wu[5], Haoling Li[5], Zhaoxuan Li[7],
Jiongchi Yu[8], Hui Li[2*], Dong Zhang[1*]

[1]*Tencent*　[2]*Peking University*　[3]*Fudan University*　[4]*Shanghai Jiao Tong University*
[5]*Tsinghua University*　[6]*Zhejiang University*　[7]*Institute of Information Engineering, Chinese Academy of Sciences*
[8]*Singapore Management University*

## Abstract

The increasing adoption of large language models (LLMs) in software engineering necessitates rigorous security evaluation of their generated code. However, existing benchmarks are inadequate, as they focus on isolated code snippets, employ unstable evaluation methods that lack reproducibility, and fail to connect the quality of input context with the security of the output. To address these gaps, we introduce A.S.E (AI Code Generation Security Evaluation), a benchmark for repository-level secure code generation. A.S.E constructs tasks from real-world repositories with documented CVEs, preserving full repository context like build systems and cross-file dependencies. Its reproducible, containerized evaluation framework uses expert-defined rules to provide stable, auditable assessments of security, build quality, and generation stability. Our evaluation of leading LLMs on A.S.E reveals three key findings: (1) Claude-3.7-Sonnet achieves the best overall performance. (2) The security gap between proprietary and open-source models is narrow; Qwen3-235B-A22B-Instruct attains the top security score. (3) Concise, "fast-thinking" decoding strategies consistently outperform complex, "slow-thinking" reasoning for security patching.

---

[†]Co-first authors.

[*]Corresponding authors: Hui Li (`lih64@pkusz.edu.cn`) and Dong Zhang (`zalezhang@tencent.com`).

## 1 Introduction

Large language models (LLMs) are rapidly permeating software engineering workflows, from code completion and synthesis to refactoring and even vulnerability repair [1, 2, 3]. Their adoption is broad, their release cadence is fast, and generated code increasingly reaches production. This reality elevates security from a secondary criterion to a primary requirement: functionally correct but insecure code is unacceptable in practice. Recent studies show that LLM-generated code can embed, propagate, or even amplify vulnerabilities, particularly when models operate with incomplete context or in projects with complex inter-module dependencies [4, 5, 6, 7]. The central challenge, therefore, is to *evaluate whether LLMs can produce secure, correct, and stable code* under the constraints and coupling that characterize real repositories.

Recent years see a wave of functionality- and security-focused benchmarks that standardize evaluation protocols, task formats, and toolchains. Examples include function-level code generation tasks with test cases, vulnerability taxonomies aligned with CWE, dynamic analysis and fuzz testing, and cross-language corpora. These resources provide a foundation for subsequent research [8, 9, 10, 11, 12, 13, 14]. Yet when the target is security in real engineering settings, three limitations remain notable. **(1) granularity is mismatched**: tasks focus on functions or snippets and do not cover repository-level

1

dependencies or build constraints. **(2) evaluation is unstable**: pipelines that rely on LLM-based judgment or static application security testing (SAST) lack reproducibility and effective control of false positives. **(3) viewpoint is narrow**: studies isolate either the model or the output and rarely link context supply and retrieval to the security, quality, and stability of the produced code.

With regard to **(1) granularity mismatch**, many benchmarks organize tasks at the snippet level. HumanEval [15] and MBPP [16] rely on hand-written function pairs with unit tests to measure syntax and basic semantics. Security-focused work builds CWE-aligned cases on short fragments—SecurityEval applies static analysis, BaxBench stresses hazardous API use such as eval-like risks, and CWEval couples vulnerability outcomes with functional tests [8, 12, 11]. These designs favor control and scale but omit repository-level factors such as cross-file calls, configuration, and build scripts, so they miss risks from multi-module flows, third-party dependencies, and deployment. Repository-scale efforts begin to appear: RepoBench and Long Code Arena add project context to test long-context coding but emphasize functional correctness [17, 18]. SecRepoBench evaluates real C and C++ projects with dynamic analysis and fuzz testing, which improves realism, yet language coverage and task mix remain limited and resource demands are high [14]. A gap therefore remains between task granularity and engineering context.

With regard to **(2) unstable evaluation and limited reproducibility**, two approaches are common. First, an LLM serves as a judge that assigns natural-language scores or applies heuristic descriptions to decide whether a vulnerability exists [19]. This approach reduces cost and increases speed, but it is sensitive to prompt manipulation and linguistic ambiguity, and decisions for the same output vary across models and temperatures, which hinders reproducibility. Second, pipelines invoke generic SAST with default rules [13]. Automation improves, but cross-language settings show high false-positive and false-negative rates, and context calibration for specific CWE categories and project structures is missing [8, 19, 9]. A lack of containerized envi-
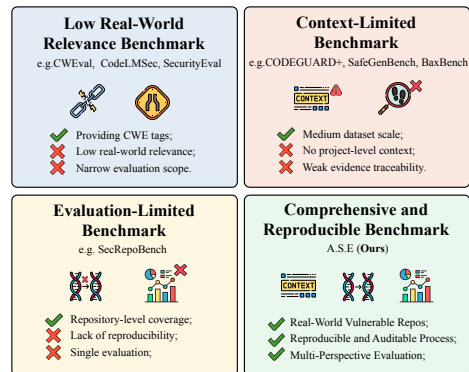


Figure 1: Comparison of code security assessment approaches.

ronments and one-click replay further reduces the reproducibility of the same evaluation across machines and dependency versions.

With regard to **(3) narrow viewpoint**, many studies fix context length or apply a uniform snippet prompt and report averages under a single input setting. They rarely probe how model capability changes as information supply varies. Long-context and cross-file retrieval benchmarks primarily target readability and functional correctness. Cross-CodeEval [20] and REPOCOD [21] test structural consistency and code understanding under cross-file dependencies and repository-level tasks, but coverage along the security dimension remains limited. FEA-Bench [22] introduces engineering-oriented requirement descriptions and inter-file coupling to assess functional metrics, yet it does not systematically test whether security fixes cause regressions or reproduce in a stable way. In security-oriented evaluations, outcomes often reduce to counts of static alerts, which neglect capability limits under different context budgets and retrieval strategies and therefore fail to establish a causal link from information supply to code results.

To address aforementioned challenges, we introduce A.S.E (AI Code Generation Security Evaluation), a benchmark purpose-built for repository-level secure code generation. A.S.E aligns its design

with the three perspective: **(1) the data design**, **(2) the evaluation framework** and **(3) the evaluation scope**.

First, in **(1) the data design**, A.S.E builds tasks from high-quality open-source repositories with documented CVEs and centers each task on security-sensitive regions that expose cross-file interactions and build-time constraints. Each task provides only the minimal and pertinent repository context, which requires models to reason over realistic call chains and interfaces rather than solve quiz-style snippets. To mitigate data leakage from public training corpora while preserving semantics, the benchmark applies light semantic and structural mutation, such as identifier renaming and equivalent control-flow reshaping, so that success reflects reasoning rather than memorization. Second, in **(2) the evaluation framework**, A.S.E provides a Dockerized environment that deterministically reproduces the vulnerable state and validates candidate fixes within the project. Security checks rely on expert-defined rules that combine industry-grade analyzers such as CodeQL and Joern with CWE-specific logic. Build and syntax checks ensure that a secure patch integrates and compiles within the original project. These mechanisms replace subjective judgments with auditable and repeatable signals. Third, in **(3) the evaluation scope**, the benchmark probes model capacity to exploit available information by adapting the context window for each model, up to 128k tokens, and by using retrieval models to surface the most relevant functions and files. In parallel, the benchmark evaluates generated code along three axes—security, build-level quality, and generation stability—through repeated runs and in-repository validation. This strategy links input information and retrieval to model reasoning and concrete downstream effects within the repository.

Based on A.S.E, we evaluate mainstream proprietary and open-source LLMs under the same repository-level protocol and obtain three findings. First, on the overall score, Claude-3.7-Sonnet attains the top performance with 52.79 points, surpassing o3 under identical settings. Second, the gap between proprietary and open-source models is small, especially on the security metric; in fact, Qwen3-235B-A22B-Instruct achieves the best Security score and exceeds Claude-3.7-Sonnet by 1.31 points. Third, across comparisons of reasoning strategies, slow-thinking configurations, which allocate more deliberate computation or multi-step reflection, underperform fast-thinking configurations that rely on concise direct decoding, indicating that increased reasoning budget does not necessarily translate into better repository-level security fixes.

The main contributions are as follows:

- **Repository-level benchmark from real code.** We release A.S.E, built from CVE-grounded, real-world repositories that preserve build systems and cross-file dependencies; tasks apply light semantic/structural mutations for realism and scalability.

- **Vulnerability-targeted, reproducible evaluation.** We design security metrics and an automated Dockerized pipeline with in-repository validation and CWE-calibrated checks, yielding stable, repeatable results; evaluation covers both model capability and code across multiple scopes (function, file, component, repository) and reports security, build quality, and stability.

- **Extensive experiments and findings.** Under a unified protocol, we benchmark mainstream LLMs and find that Claude-3.7-Sonnet leads overall (63.01), Qwen3-235B-A22B-Instruct tops Security (48.06), and concise decoding outperforms slow-thinking—highlighting core challenges and directions for secure deployment.

## 2 Related Work

### 2.1 Snippet- vs. Repository-Level Code Security Benchmarks

The mainstream benchmarks for code generation emphasized small, self-contained tasks at the function or snippet level (e.g., HumanEval [15], MBPP [16]), using unit tests or reference outputs to measure basic functional correctness. Security-focused datasets and vulnerability taxonomies followed this pattern by constructing CWE-aligned fragments or micro-examples that expose canon-

ical weaknesses and hazardous API uses [8, 12, 11]. Snippet-level designs afford experimental control and scale, but they omit repository-level factors—cross-file call graphs, build systems, dependency resolution, and deployment configurations—that materially affect whether a generated patch is secure in practice.

To bridge the gap, repository-level benchmarks have emerged that supply broader project context and long-range dependencies (e.g., RepoBench [17], Long Code Arena [18], CrossCodeEval [20], REPOCOD [21], and FEA-Bench [22]) and therefore better reflect development workflows and multi-file reasoning. SecRepoBench [14] further advances realism by assembling CVE-grounded cases from real C/C++ projects and exercising dynamic testing and fuzzing in project contexts. Nevertheless, many repository designs remain optimized for functional correctness or structural consistency rather than explicit, reproducible security validation: language coverage, task diversity, and the resource costs of dynamic validation limit their applicability. Therefore, we introduce A.S.E benchmark, which focuses tasks on CVE-associated, security-critical areas in authentic repositories, incorporating subtle semantic and structural alterations (e.g., renaming identifiers, reshaping equivalent control flows) to mitigate memorization risks. This approach compels models to perform cross-file reasoning while adhering to build constraints, thereby aligning evaluation with practical engineering and exposing vulnerabilities overlooked by snippet-based benchmarks.

## 2.2 Code Security Evaluators

Two dominant paradigms appear in prior evaluators for code security: (1) using LLMs as semantic judges and (2) invoking rule-based static/dynamic analyzers (SAST/fuzzing) to detect vulnerabilities. The LLM-as-judge approach scales cheaply and can capture semantic subtleties, but it is sensitive to prompt choice, model version, and decoding randomness, producing judgments that are difficult to reproduce and audit [19]. Conversely, off-the-shelf static analyzers provide deterministic rule-based signals but often suffer high false-positive/false-

negative rates when applied across languages or when contextual calibration is missing; they may also miss semantic flaws that manifest only at runtime [8, 13]. Hybrid pipelines (combining LLM judgment, static analysis, and dynamic checks) improve coverage but still face two recurring problems: (i) linguistic or prompt-induced ambiguity when LLMs contribute labels and (ii) lack of reproducible, containerized execution contexts that guarantee identical results across environments. Benchmarks such as CodeLMSec, CyberSecEval, and SafeGenBench explore variants of these trade-offs but do not fully eliminate ambiguity or provide per-CWE, auditable detection logic [9, 19, 13]. To mitigate these issues, A.S.E prioritizes reproducibility and precision by (a) executing tasks in Dockerized, project-specific environments, and (b) relying on expert-calibrated, CWE-targeted static queries (e.g., CodeQL/Joern) and deterministic build/test checks rather than free-form LLM judgments. This yields auditable pass/fail signals and substantially reduces prompt sensitivity and inter-run variance.

## 2.3 From Single- to Multi-Dimensional Code Security Analysis

Existing evaluations focus on single-dimensional outcomes (e.g., unit-test pass rates or counts of static alerts), which obscures important trade-offs between correctness, security, and reliability. Cross-file and long-context benchmarks (CrossCodeEval [20], FEA-Bench [22], REPOCOD [21]) examine structural and reasoning capabilities but seldom tie those capabilities to security outcomes under varying information budgets. Likewise, security-oriented studies often reduce results to binary vulnerable/not-vulnerable labels without assessing integration (buildability) or the consistency of generated fixes across repeated runs and different retrieval/context settings. These omissions leave unanswered whether additional context or retrieval actually improves secure repair, and whether successful repairs generalize stably across prompts and decoding regimes. To response these challenges, A.S.E measures three complementary axes—(1) security (vulnerability elimination and absence of new

4

(a) A.S.E Benchmark Construction

Seed Dataset(High-Quality GitHub Repository)

Dual Code Mutation — Code Structure Mutation — Code Semantic Mutation

A.S.E Dataset

CVE Patch — Mutated Dataset

Customized SAST Tool — Seed Dataset

Security Expert Review
- ✓ Precise vulnerability identification
- ✓ Systematic extraction of relevant information
- ✓ Remove code related to vulnerabilities

(b) Model Code Generation

Generated Repository

LLMs Generated Code

Vulnerability Functional Description

Code Context

+

Incomplete Repository

(c) Security Evaluation of Model-Generated Code

Comprehensive Assessment

- Code Security (60%)
- Code Quality (30%)
- Generation Stability (10%)
- License
- Reasoning Mode(Fast/Slow)
- Context Length

Comprehensive Evaluation Leaderboard

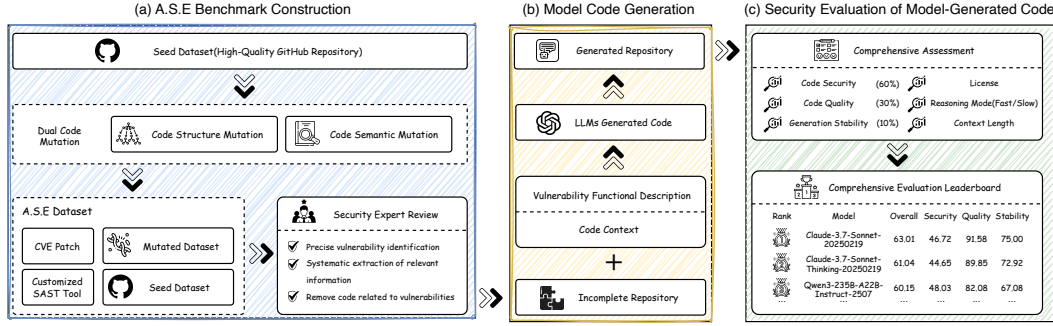| Rank | Model | Overall | Security | Quality | Stability |
|---|---|---|---|---|---|
| 🥇 | Claude-3.7-Sonnet-20250219 | 63.01 | 46.72 | 91.58 | 75.00 |
| 🥈 | Claude-3.7-Sonnet-Thinking-20250219 | 61.04 | 44.65 | 89.85 | 72.92 |
| 🥉 | Qwen3-235B-A22B-Instruct-2507 | 60.15 | 48.03 | 82.08 | 67.08 |
| | ... | ... | ... | ... | ... |

Figure 2: Overall workflow of A.S.E. (a) A.S.E benchmark construction: from high-quality GitHub seeds, we build the A.S.E dataset via dual mutations (structure/semantic), CVE patches, and a customized SAST tool, followed by expert curation. (b) Model code generation: given an incomplete repository, a vulnerability description and context guide LLMs to complete the repository. (c) Security evaluation: comprehensive assessment with security, quality and stability.

flaws), (2) build and functional quality (compilation, tests, and integration), and (3) generation stability (consistency across repeated runs and reasoning configurations). We systematically vary context supply (including retrieval and model window up to 128k tokens where supported) to quantify how model capability, context budget, and retrieval interact to produce secure, usable, and repeatable code in real repositories.

# 3 The A.S.E Benchmark

This section overviews the A.S.E benchmark workflow (Fig. 2), beginning with the design philosophy (Sec. 3.1) that governs data provenance, process control, and evaluation scope. Guided by these principles, we construct the benchmark by curating repository level tasks with verifiable vulnerabilities and pinned, reproducible environments (Sec. 3.2). We then detail the evaluation pipeline, which includes in repository validation, metrics and scoring, and reporting protocols, to enable objective and auditable comparison across models and to establish a reproducible, security centered basis for assessing code generation (Sec. 3.4).

## 3.1 Design Philosophy

The A.S.E benchmark targets secure code generation under real engineering constraints through three principles: (i) real-world, repository-level data sources, (ii) objective and traceable evidence and processes, and (iii) a multi-perspective evaluation strategy spanning model and code perspectives.

**(i) Real-world, repository-level data sources.** A.S.E builds tasks from active open-source repositories with documented CVEs and verifiable diffs, retaining only projects that pass static analysis and expert screening for stable builds and resolvable dependencies. We precisely delimit security-sensitive regions while preserving cross-file interactions, interface contracts, configuration, and build scripts so models must reason under project-level constraints. To curb leakage without distorting semantics, we apply light semantics-preserving transformations and pin baseline commits, toolchains, and dependencies in Docker for reproducibility.

**(ii) Objective and traceable processes.** All steps are containerized for deterministic replay, with one-click scripts reconstructing vulnerable states and evaluating candidate patches in-repository across environments. Security judgments are CWE-calibrated by combining CodeQL/Joern with expert

Algorithm-Guided Screening and Preselection

Step 1 : Determining Data Sources

Raw CVE data (**100,000+** entries, including internal and public data)

Enterprise CVE repository (projects that contain complete commit histories)

Open source vulnerability dataset (including CVE IDs and fix commit records)

Manually retrieved and supplemented reproducible project snapshots

Step 2 : Filtering Candidate Repository　　Automated Filtering

SWE Filtering
✓ Vulnerability types related to web development in the 2024 Top CWE
✓ Including commit information

Vulnerability Filtering
✓ Cross-site scripting (XSS)
✓ SQL injection
✓ path traversal
✓ command injection

Language Filtering
✓ Java, Python
✓ Go, PHP
✓ JavaScript

Filtered Data 　**5000+** CWE datas　Pre-screened repositories　Tools & Commits　**199** Candidate Repositories

Expert-Guided Curation and Refinement

Step 3.1 : Candidate Refinement
**199** Candidate Repositories　Manual review
SAST false positives ✗
large-scale commits ✗
Reproducibility Checks　Fine Filtering

Step 3.2 : Robust Validation　Manual Filtering
Docker-based Reproducibility　Expert-grade SAST
tool-built-in and CVE-specific custom detection rules
one-to-one rule-to-CVE mapping for result traceability
Commit Exploitability Check　Dual-toolchain Detection (CodeQL + Joern)
autopoc:latest　ai_gen_code:latest　aiseceval

Step 4 : Dataset Expansion
**40** Seed Repositories
Structural mutation　Semantic mutation
**80** Mutated Repositories

Dual-toolchain Detection: Code Flow Validation　EXAMPLE Taint Flow Analysis

1. HrmEmployeeLeaveRecordController.java — queryPageList(...)
**Source entry:** user input via HTTP request flows unvalidated into the business logic layer.
2. HrmEmployeeLeaveRecordServiceImpl.java — queryLeaveRecordPageList(...)
**Business logic layer:** forwards user input without security filtering or validation.
3. HrmEmployeeLeaveRecordMapper.java — queryLeaveRecordPageList(...)
**DAO layer (MyBatis Mapper):** passes user input directly to the SQL query.
4. HrmEmployeeLeaveRecordMapper.xml — queryLeaveRecordPageList(...)
**Sink (exploitation point):** SQL is defined in an XML configuration (MyBatis), where user input is directly concatenated into the query. Because it's in an ORM XML file rather than Java source, static analysis tools may fail to detect this SQL injection.

A.S.E Dataset
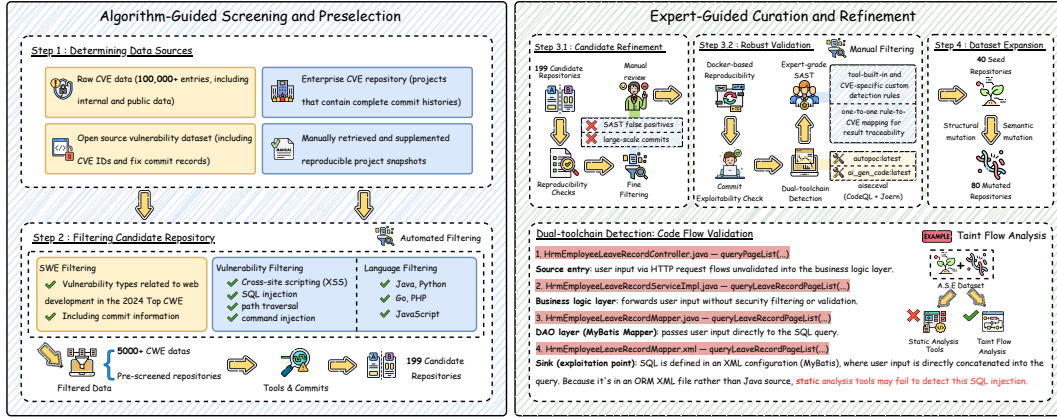Static Analysis Tools ✗　Taint Flow Analysis ✓

Figure 3: Overview of A.S.E benchmark construction. **Algorithm-guided screening and preselection (left):** aggregate CVE-linked sources and automatically filter repositories by web-related CWEs, vulnerability types (XSS, SQL injection, path traversal, command injection), and languages (Java, Python, Go, PHP, JavaScript). **Expert-guided curation and refinement (right):** conduct manual review, reproducibility and exploitability checks, and dual-toolchain SAST (e.g., CodeQL + Joern) with CVE-specific rules; then expand 40 seed repositories via structural/semantic mutation to 80 variants.

rules, and we emit auditable artifacts—rule hits, locations, and data/control-flow slices—for independent verification. We enforce quality gates (successful builds and sanity tests) and regression-test rules to bound false positives/negatives and quantify inference stability under repeated generations with fixed randomness.

**(iii) Multi-perspective evaluation strategy.** We evaluate from multiple vantage points that reflect real engineering practice. Context windows are adapted per model with retrieval of the most relevant functions/files, and patches are validated in situ—only build-integrable changes proceed to security checks while we track regressions and build breakages. Scoring spans Security, Quality, and Stability and generalizes across languages and high-frequency vulnerability families, enabling stratified analysis by language, CWE, repository scale, and context budget.

Following these guidelines, we introduce the A.S.E benchmark to provide an objective evaluation that, with minimal subjective influence, delineates the genuine capability boundaries of model capa-

bility in repository-level code generation regarding security.

## 3.2 Benchmark Construction

Guided by our design philosophy, we construct the A.S.E benchmark as shown in Fig. 2 (a) and Fig. 3. To ensure realistic scenarios and relevant security expertise, we organize a dedicated team of ten contributors with strong backgrounds in cybersecurity and web development from top-tier universities, including five Ph.D. candidates and five master-level students. All contributors have extensive hands-on experience in vulnerability discovery, analysis, and remediation, with emphasis on common web vulnerabilities such as XSS, SQL injection, and path traversal. Each contributor adheres to strict secure coding standards and is familiar with static code analysis techniques. Building on this setup, the benchmark construction follows four stages:

**Step 1: Determining data sources.** We collect real-world software projects from an enterprise-internal vulnerability repository and from public platforms

such as GitHub as candidate sources. This stage ensures representativeness and diversity, so the dataset covers a broad range of software types and vulnerability scenarios.

**Step 2: Filtering candidate repository.** We filter candidates using criteria that capture project activity and completeness of vulnerability information. We retain only projects with active maintenance, adequate test coverage, and consistent development practices. Each retained project includes at least one publicly disclosed CVE and provides a complete patch diff at the commit level. Each vulnerability case includes detailed metadata, including a description, a CVE identifier, a severity rating, and clear reproduction steps. These conditions ensure high data quality, complete records, and reproducibility for the next stage.

**Step 3: Expert-guided refinement and quality filtering.** This stage contains two sub-steps.

**Step 3.1: Candidate refinement.** Security experts analyze projects that pass the initial screening and mark the exact code region for each vulnerability. The system then extracts surrounding context, including function signatures, API definitions, and call chains, to reconstruct the code environment. We create CodeQL or Joern queries for each vulnerability and verify the propagation path. We discard candidates with ambiguous or inconsistent findings or with code that current tools cannot analyze with sufficient precision. For retained candidates, we customize and strengthen the rules to improve detection precision and to align scan results with the confirmed vulnerability.

**Step 3.2: Robust validation.** After confirming the candidate repositories, we remove the labeled vulnerable code to create a fill-in-the-code setting. We combine the functional description of the vulnerability with the extracted context to generate a structured code-completion prompt for each task. This design provides clear input and output semantics and enables a model to reason over repository structure and logic rather than only local snippets, so the evaluation reflects context understanding and code generation ability. After expert review, we select 40 repositories with real CVE records as data sources and set a verified baseline commit for each repository as a unified starting point for task construction and evaluation. This setup ensures a consistent and stable code state, avoids uncertainty from repository updates, and supports authenticity, reliability, and reproducibility.

**Step 4: Dataset expansion.** After constructing base tasks from real-world vulnerability repositories, we expand the task set to increase data volume and coverage under a semantics-preserving constraint. We apply two types of transformations. The first type is semantic transformations, which perform systematic variable and function renaming and equivalent API substitution to diversify surface expressions. The second type is structural transformations, which adjust control flow, refactor the call graph, or reorganize file layout to introduce structural differences. These operations preserve functional behavior and vulnerability semantics while changing implementation details, which reduces overlap with public repository code that may appear in training data used to train evaluated models and mitigates potential data contamination. The expanded variants therefore enable a more comprehensive assessment of model robustness and generalization.

After these steps, we complete the A.S.E benchmark dataset. he benchmark focuses on four high frequency vulnerability classes, such as XSS, SQL injection, path traversal, and command injection, and it covers five mainstream programming languages while remaining extensible to additional classes and languages. From 40 seed tasks, we generate 80 semantics preserving variants through semantic and structural transformations. This process yields 120 realistic and reproducible vulnerability scenarios that together constitute the benchmark.

To ensure portable and reproducible evaluation, we provision a build and runtime environment for each scenario and package it with Docker. In addition to the scenario environments, we containerize relevant static application security testing tools to enable evaluation in a unified environment with one command. This design supports systematic evaluation, reproducibility, and cross-platform usability.
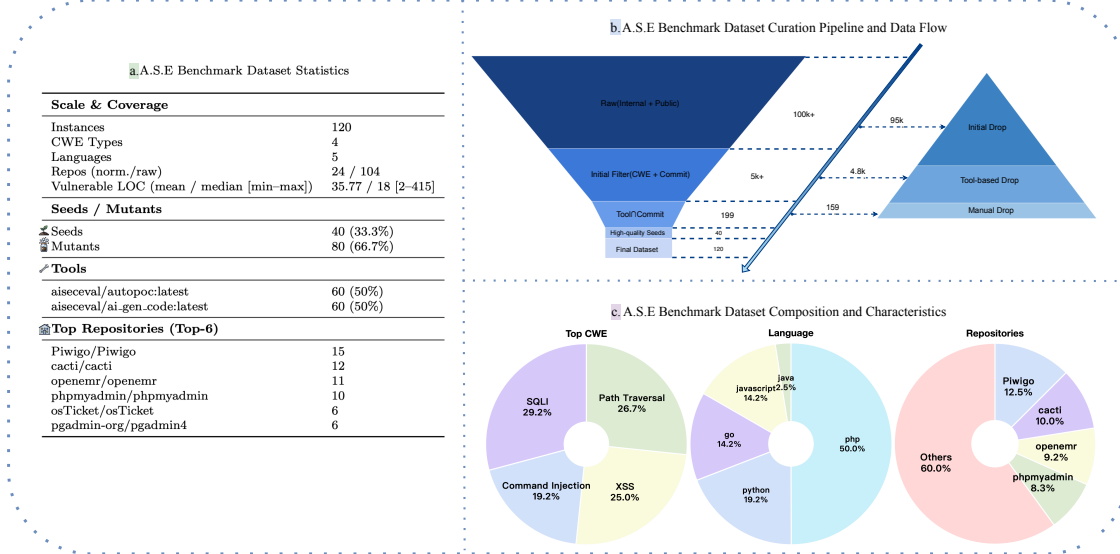
Figure 4: Statistics of A.S.E benchmark.

## 3.3 General Statistics

The A.S.E. benchmark comprises 120 repository-level vulnerability instances, including 40 seed repositories and 80 systematically mutated variants. Each instance is rigorously validated to guarantee reproducibility and consistency across diverse evaluation settings.

The benchmark targets four vulnerability categories that are common in real web projects. Each category maps one-to-one to a corresponding CWE entry: SQL Injection (29.2%, CWE-89), Path Traversal (26.7%, CWE-22), Cross-Site Scripting (25.0%, CWE-79), and Command Injection (19.2%, CWE-78). This mapping defines the dataset at the CWE level and limits it to these four categories. The design aligns tasks and metrics with security-critical issues that research and practice widely recognize. For each category, the benchmark defines a task that evaluates whether the model can:

- Cross-Site Scripting (XSS): detect and mitigate malicious scripts injected into trusted web con-

texts.

- SQL Injection: identify and prevent attempts to compromise database integrity or to extract data without authorization through malicious SQL queries.

- Path Traversal: detect and block attempts to access files and directories outside the designated web root through manipulation of file paths.

- Command Injection: recognize and prevent attempts to execute unauthorized operating system commands on the host through application vulnerabilities.

From the perspective of programming languages, A.S.E encompasses five mainstream environments that reflect realistic multi-language software development. The distribution is concentrated in PHP (50.0%), followed by Python (19.2%), Go (14.2%), JavaScript (14.2%), and Java (2.5%). This imbalance reflects the prevalence of vulnerable PHP projects in practical ecosystems, while also offering

opportunities to examine cross-language generalization in vulnerability detection and mitigation.

Taken together, the combination of authentic vulnerability cases, systematic mutation strategies, and expert validation establishes A.S.E as a benchmark corpus with broad coverage and rigorous quality. It enables reproducible, multi-dimensional evaluation of security-aware code generation and provides a reliable foundation for advancing research in automated vulnerability analysis.

## 3.4 Evaluation Pipeline

A.S.E. uses a two-stage evaluation pipeline that emulates repository-level development and patching: code generation (Fig. 2 (b)) followed by security assessment (Fig. 2 (c)) .

**Code generation.** For each benchmark instance, A.S.E retrieves the corresponding GitHub repository and checks out the baseline commit that contains the vulnerability. Expert annotations guide automatic masking of the vulnerable region in the target file, which is replaced by the special token `<masked>`. The language model receives an input that combines two elements: (i) the masked file with a functional description generated by Claude-Sonnet-4 [23] and refined by experts, and (ii) repository-level context that includes related source files selected by BM25 [24] ranking and the README of the project. This setup encourages reasoning over a realistic repository context rather than isolated snippets and supports configurable context lengths, which enables systematic analysis of how repository scope affects generation performance. The model is instructed to produce patch-formatted code, and the patch is applied automatically with tools such as `git apply`. To assess stability, each instance runs three times under identical conditions.

**Security assessment.**

The generated code is evaluated across three dimensions, including quality, security, and stability. The evaluation first checks whether a patch integrates cleanly and satisfies basic static checks, then examines whether vulnerability counts decrease after the patch, and finally assesses variation across repeated runs. The following definitions formalize these steps.

**Quality** reflects whether a patch integrates and satisfies essential static checks. A test counts as successful only if the patch merges and passes static analysis and syntax checks. The quality score is defined as follows:

$$\text{Quality} = \frac{1}{N}\sum_{t=1}^{N} q_t, \tag{1}$$

where $N$ denotes the total number of tests, and $q_t = 1$ if test $t$ merges and passes all checks, and $q_t = 0$ otherwise.

**Security** considers the change in detected vulnerabilities after patch integration. A.S.E. applies expert crafted static analysis rules that are tailored to each CVE in order to compute vulnerability counts. The security score is defined as follows:

$$\text{Security} = \frac{1}{N}\sum_{t=1}^{N} s_t, \tag{2}$$

where $s_t = 1$ if $v_{\text{after}}(t) < v_{\text{before}}(t)$ and $s_t = 0$ otherwise, and where $v_{\text{before}}(t)$ and $v_{\text{after}}(t)$ denote the numbers of detected vulnerabilities before and after patch integration for test $t$.

**Stability** evaluates consistency across independent runs. A lower standard deviation indicates higher stability. To map lower variation to a higher score, the method applies min to max scaling. For each benchmark instance $i \in \mathcal{B}$, let $\sigma_i$ denote the standard deviation over three runs. The normalized value is defined as follows:

$$\tilde{\sigma}_i = \begin{cases} 1 - \dfrac{\sigma_i - \sigma_{\min}}{\sigma_{\max} - \sigma_{\min}}, & \text{if } \sigma_{\max} > \sigma_{\min}, \\ 1, & \text{otherwise,} \end{cases} \tag{3}$$

where $\sigma_{\min} = \min_i \sigma_i$ and $\sigma_{\max} = \max_i \sigma_i$. The stability score is computed as follows:

$$\text{Stability} = \frac{1}{|\mathcal{B}|}\sum_{i \in \mathcal{B}} \tilde{\sigma}_i \tag{4}$$

The **overall score** combines the three dimensions with fixed weights as follows:

$$\text{Overall} = 0.6 \times \text{Security} + 0.3 \times \text{Quality} + \tag{5}$$
$$0.1 \times \text{Stability}. \tag{6}$$

This pipeline enables reproducible project level testing and provides a systematic assessment of functional correctness and security robustness of patches generated by large language models in realistic development settings.

# 4 Experiments

## 4.1 Experimental Setup

We evaluate 26 representative, state-of-the-art (SOTA) large language models, including 18 proprietary models and 8 open-source models. A key selection criterion is the inclusion of both "fast thinking" and "slow thinking" modes where available, which enables a comprehensive comparison of performance. We run three trials per model.

The 18 proprietary models include flagship offerings across multiple model families: Claude-3.7-Sonnet [25], Claude-Sonnet-4 [23], and Claude-Opus-4 [23], together with the corresponding "thinking" variants; GPT-4o [26], GPT-4.1 [27], Codex-mini [28], and several additional versions; the Grok series (Grok-3 [29], Grok-4 [30], Grok-3-mini [29]); Gemini-2.5-Pro [31]; Qwen-Coder-Plus [32]; and Hunyuan-T1 [33].

For the open-source category, we select 8 widely used models that cover diverse architectures. The set includes models from the Qwen3 series [34] (Qwen3-235B-A22B-Instruct, Qwen3-Coder, Qwen3-235B-A22B), DeepSeek-V3 [35], DeepSeek-R1 [36], Kimi-K2 [37], and GLM-4.5 [38].

## 4.2 Main Results

As shown in Tab. 1, we evaluated 26 SOTA LLMs on our A.S.E benchmark. The main observations and findings are as follows.

**I. Current LLMs face significant challenges in secure coding.** While code LLMs make strong progress in repository-level generation, the results show that secure coding remains an ongoing weakness. None of the evaluated LLMs exceed the 50-point threshold in the Code Security score, which indicates that even state-of-the-art models still struggle to prevent common vulnerabilities such as SQL injection, cross-site scripting, and path traversal. For example, Claude-3.7-Sonnet-20250219 leads the A.S.E. benchmark with a total score of 63.01 and a Code Quality score of 91.58, yet its Code Security score is 46.72. Claude-Sonnet-4 displays a similar pattern, achieving the highest Code Quality score of 92.37 but only 34.78 in Code Security. These cases highlight a recurring imbalance: current Code LLMs excel at producing syntactically correct and useful code but remain less effective at ensuring robust security in the generated code.

**II. The A.S.E benchmark introduces substantial complexity in repository-level scenarios.** As a repository-level evaluation, the A.S.E benchmark presents significant challenges for Code LLMs that primarily excel at snippet-level tasks. For instance, although GPT-o3 performs well on SafeGen-Bench [13], performance on the A.S.E. benchmark drops and falls behind many other models. This reversal arises from additional complexities introduced by the A.S.E benchmark, including cross-file dependency resolution and long-context understanding—capabilities that extend beyond isolated code snippet generation. Nevertheless, some models, such as the Claude series, maintain leading performance at both the snippet-level and the repository-level, which indicates stronger generalization across different granularities of code security understanding.

**III. Slow-thinking paradigms show security regressions.** As shown in Fig. 5, slow-thinking models tend to underperform in Code Security relative to fast-thinking counterparts. In some cases, the gap also extends to code quality and generation stability. For example, Claude-3.7-Sonnet-Thinking scores 44.65 in Code Security, slightly below the fast-thinking counterpart. Claude-Sonnet-4-Thinking shows even larger drops across all metrics compared to Claude-Sonnet-4 (non-thinking mode). Other slow-thinking models, such as DeepSeek-R1 and Gemini-2.5-Pro-Exp, also show weaker performance in Code Security. In contrast, fast-thinking models such as Qwen3-235B and Claude-3.7-Sonnet tend to achieve better security perfor-

Table 1: The leaderboard of various advanced Code LLMs on the A.S.E. benchmark. ⚡is the fast-thinking mode and 🐢indicates slow-thinking mode.

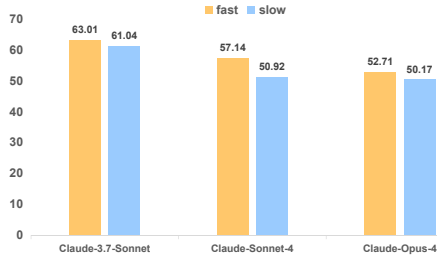| Rank | Model | License | Thinking | Overall | Security | Quality | Stability |
|---|---|---|---|---|---|---|---|
| 1 | Claude-3.7-Sonnet-20250219 | Proprietary | ⚡ | 63.01 | 46.72 | 91.58 | 75.00 |
| 2 | Claude-3.7-Sonnet-Thinking-20250219 | Proprietary | 🐢 | 61.04 | 44.65 | 89.85 | 72.92 |
| 3 | Qwen3-235B-A22B-Instruct-2507 | Open Source | ⚡ | 60.15 | 48.03 | 82.08 | 67.08 |
| 4 | Qwen3-Coder | Open Source | ⚡ | 59.31 | 42.69 | 85.16 | 81.54 |
| 5 | DeepSeek-V3-20250324 | Open Source | ⚡ | 58.59 | 40.89 | 85.87 | 82.94 |
| 6 | Claude-Sonnet-4-20250514 | Proprietary | ⚡ | 57.14 | 34.78 | 92.37 | 85.65 |
| 7 | Kimi-K2-20250711-Preview | Open Source | ⚡ | 55.29 | 37.82 | 79.90 | 86.25 |
| 8 | GPT-4o-20241120 | Proprietary | ⚡ | 55.10 | 45.65 | 72.46 | 59.67 |
| 9 | Qwen-Coder-Plus-20241106 | Proprietary | ⚡ | 53.55 | 37.98 | 73.78 | 86.27 |
| 10 | Claude-Opus-4-20250514 | Proprietary | ⚡ | 52.71 | 31.95 | 85.82 | 77.91 |
| 11 | Grok-3 | Proprietary | ⚡ | 52.18 | 38.64 | 73.54 | 69.41 |
| 12 | DeepSeek-R1-20250528 | Open Source | 🐢 | 51.76 | 38.01 | 74.39 | 66.38 |
| 13 | Gemini-2.5-Pro-Exp-20250325 | Proprietary | ⚡ | 51.02 | 29.98 | 84.04 | 78.21 |
| 14 | Claude-Sonnet-4-Thinking-20250514 | Proprietary | 🐢 | 50.92 | 34.10 | 76.81 | 74.22 |
| 15 | Claude-Opus-4-Thinking-20250514 | Proprietary | 🐢 | 50.17 | 30.70 | 79.84 | 77.98 |
| 16 | GLM-4.5 | Open Source | ⚡ | 49.80 | 35.92 | 70.24 | 71.74 |
| 17 | Grok-4 | Proprietary | ⚡ | 42.40 | 29.53 | 59.78 | 67.42 |
| 18 | o4-mini-20250416 | Proprietary | 🐢 | 41.35 | 27.87 | 60.74 | 64.07 |
| 19 | Grok-3-mini | Proprietary | ⚡ | 30.49 | 22.37 | 38.15 | 56.26 |
| 20 | Codex-mini-latest | Proprietary | ⚡ | 29.71 | 22.96 | 34.68 | 55.29 |
| 21 | Hunyuan-T1-20250321 | Proprietary | 🐢 | 21.92 | 15.57 | 20.21 | 65.18 |
| 22 | Qwen3-235B-A22B-Thinking | Open Source | 🐢 | 18.11 | 9.42 | 15.60 | 77.81 |
| 23 | GPT-4.1-20250414 | Proprietary | ⚡ | 17.26 | 5.26 | 16.46 | 91.66 |
| 24 | Qwen3-235B-A22B | Open Source | ⚡ | 13.37 | 3.34 | 7.27 | 91.86 |
| 25 | o3-mini-20250131 | Proprietary | 🐢 | 13.23 | 3.67 | 3.91 | 98.57 |
| 26 | o3-20250416 | Proprietary | 🐢 | 10.22 | 0.36 | 0.36 | 98.91 |



Figure 5: Fast v.s. Slow Thinking.

mance. This pattern suggests that, although slow-thinking paradigms aim for more careful reasoning, they may inadvertently increase the risk of insecure code by producing more complex outputs or by lacking targeted security reinforcement.

**IV. High stability does not imply fewer vulnerabilities.** Some Code LLMs achieve a strong balance between coding security and generation stability, such as Claude-3.7-Sonnet; however, several models consistently produce vulnerable code in the repository-level scenario. For example, GPT-o3 achieves the highest Generation Stability score of 98.91; however, it attains 0.36 in both Code Security and Code Quality, which is the lowest among all LLMs. A similar pattern appears in models such as the GPT-4.1 series and Qwen3-235B-A22B, which demonstrate high stability while showing considerably lower coding security. These cases show that progress in generation stability does not necessarily translate into improved coding security, and they highlight the need to evaluate these dimensions independently when assessing Code LLMs.

**V. Open-source models perform comparably to closed-source Code LLMs.** Closed-source Code LLMs such as Claude-3.7-Sonnet and Claude-Sonnet-4 achieve higher Code Quality scores (91.58 and 92.37, respectively), whereas open-source models such as Kimi-K2 and Qwen-Coder-Plus deliver strong generation stability. For overall performance, closed-source and open-source LLMs achieve comparable results. Claude-3.7-Sonnet and Claude-3.7-
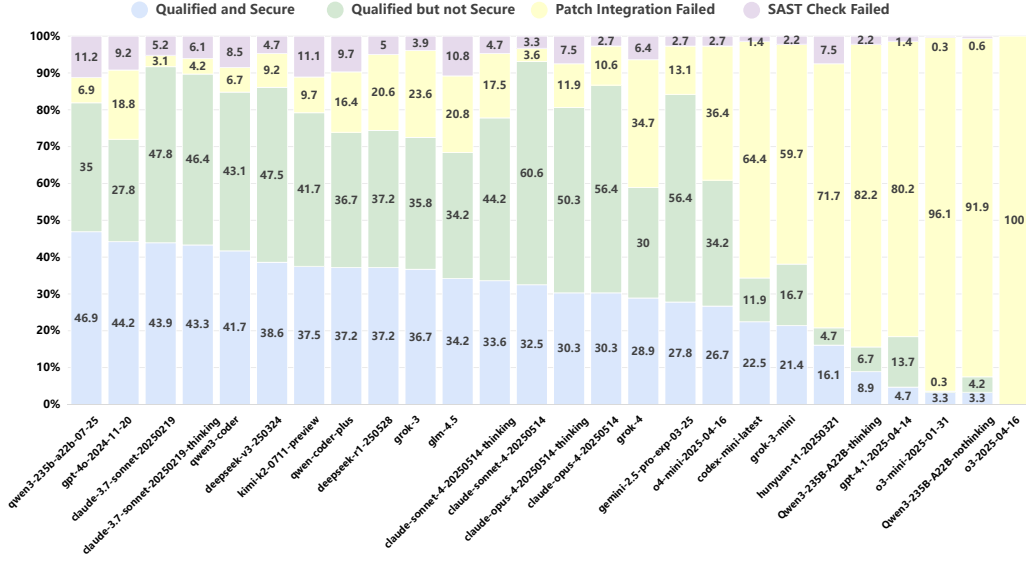
Figure 6: Attributional distribution across Code LLMs.

Sonnet-Thinking lead the benchmark, while Qwen3-235B-A22B-Instruct, an open-source model, ranks second. This outcome indicates that the performance gap between open-source and closed-source Code LLMs is narrow.

## 5 Discussion

### 5.1 Analysis

**I. Challenges posed by the A.S.E. benchmark.** The challenges posed by the A.S.E. benchmark are twofold: (i) for flagship Code LLMs, the benchmark exposes a tendency to prioritize code correctness over security; and (ii) for less powerful Code LLMs, the benchmark reveals limited basic code generation in complex, repository-level settings. As shown in Fig. 6, flagship Code LLMs generate qualified code in repository-level scenarios; however, about half of this code remains insecure. For example, for Claude-3.7-Sonnet, although 91.7% of generated code is identified as qualified, 43.8% is evaluated as insecure. This pattern suggests an emphasis on code correctness while security considerations receive insufficient attention. In contrast, for less powerful

Code LLMs, the proportion of qualified code declines sharply, and most generated code fails the SAST check, which indicates a fundamental deficiency in basic code generation in repository-level scenarios.

**II. Path traversal presents the greatest challenge.** Among the four evaluated tasks, Path Traversal presents the greatest difficulty. As shown in Fig. 8, Path Traversal poses challenges for all Code LLMs, with most models exhibiting relatively weak performance compared with the other tasks. Even the most advanced LLM achieves an overall score below 50.0 on this task. This outcome likely results from the subtlety and variety of path manipulation techniques, which are often context dependent and harder to detect than explicit injection attacks. The results indicate that current Code LLMs lack robust reasoning about file system operations and access control. Improving model understanding of file path construction and enhancing model ability to generalize across diverse traversal scenarios are necessary steps for stronger defense against this vulnerability.

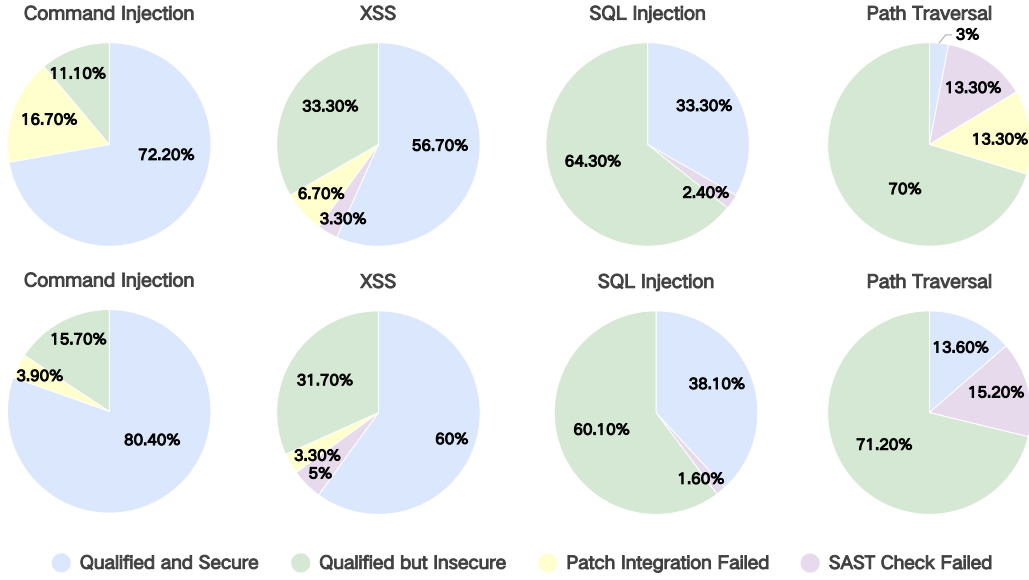**III. MoE models generally outperform dense models.** Although the architectures of some closed-

**Command Injection** · **XSS** · **SQL Injection** · **Path Traversal**

72.20% · 11.10% · 16.70%

56.70% · 33.30% · 6.70% · 3.30%

33.30% · 64.30% · 2.40%

3% · 13.30% · 13.30% · 70%

**Command Injection** · **XSS** · **SQL Injection** · **Path Traversal**

80.40% · 15.70% · 3.90%

60% · 31.70% · 3.30% · 5%

38.10% · 60.10% · 1.60%

13.60% · 15.20% · 71.20%

● Qualified and Secure  ● Qualified but Insecure  ● Patch Integration Failed  ● SAST Check Failed

Figure 7: Detailed Attribution classification of Claude-3.7-Sonnet: Original (top) and Mutation Test (bottom).



Figure 8: Detailed performance of various Code LLMs across four tasks of A.S.E benchmark.

source models remain undisclosed, nearly all leading open-source Code LLMs adopt a Mixture-of-Experts (MoE) architecture, including Qwen3-235B-A22B, Qwen3-Coder-480B-A35B-Instruct, DeepSeek-V3-671B-A37B, and Kimi-K2-Preview-1T-32B. This trend indicates that MoE-based Code LLMs generally achieve stronger security performance than dense models.

## 5.2 Case Study

We use Claude-3.7-Sonnet as a case study to examine performance across tasks, as presented in Fig. 7 (original test at the top, mutation test at the bottom).

For Path Traversal and SQL Injection, Claude-3.7-Sonnet predominantly generates high-quality code that remains insecure. In contrast, for XSS and Command Injection, the model more often produces outputs that are both well-formed and secure. As a leading model, Claude-3.7-Sonnet rarely generates poor-quality code; failures of patch application or failures to pass the SAST check are rare. This pattern indicates that at the repository level the model generally succeeds in producing correct and usable code; however, a considerable fraction of the code

13

still contains security vulnerabilities.

Model performance exhibits minimal variation before and after mutation, which suggests that the benchmark is robust, free of significant data leakage, and effective in manual construction.

## 6 Conclusion

This paper addressed the critical need for realistic evaluation of LLM-generated code security, as existing benchmarks often use isolated snippets and lack reproducible, context-aware methods. We introduced A.S.E (AI Code Generation Security Evaluation), a benchmark built from real-world repositories with documented CVEs. By preserving full project context and using a reproducible, containerized framework with expert-defined rules, A.S.E reliably assesses code security, build quality, and stability, directly linking them to the model's input. Our evaluation of leading LLMs yielded three key findings: (1) While Claude-3.7-Sonnet leads in overall performance, (2) the security gap to open-source models is narrow, with Qwen3-235B-A22B-Instruct achieving the top security score. (3) Crucially, concise "fast-thinking" decoding strategies consistently outperform complex "slow-thinking" reasoning for security repairs. These results have significant implications, suggesting prompting strategy is as critical as model choice and that open-source models are highly competitive in security. A.S.E provides a robust foundation for future research into developing more secure and reliable LLMs for software engineering.

## References

[1] Jingxuan He and Martin T. Vechev. Large language models for code: Security hardening and adversarial testing. In *CCS*, pages 1865–1879. ACM, 2023.

[2] Mohammed Kharma, Soohyeon Choi, Mohammed AlKhanafseh, and David Mohaisen. Security and quality in llm-generated code: A multi-language, multi-model analysis. *CoRR*, abs/2502.01853, 2025.

[3] Gustavo Sandoval, Hammond Pearce, Teo Nys, Ramesh Karri, Siddharth Garg, and Brendan Dolan-Gavitt. Lost at C: A user study on the security implications of large language model code assistants. In *USENIX Security Symposium*, pages 2205–2222. USENIX Association, 2023.

[4] Zhijie Wang, Zijie Zhou, Da Song, Yuheng Huang, Shengmai Chen, Lei Ma, and Tianyi Zhang. Towards understanding the characteristics of code generation errors made by large language models. In *ICSE*, pages 2587–2599. IEEE, 2025.

[5] Mohammed Latif Siddiq, Joanna Cecilia da Silva Santos, Sajith Devareddy, and Anna Muller. SALLM: security assessment of generated code. In *ASE Workshops*, pages 54–65. ACM, 2024.

[6] Avishree Khare, Saikat Dutta, Ziyang Li, Alaia Solko-Breslin, Rajeev Alur, and Mayur Naik. Understanding the effectiveness of large language models in detecting security vulnerabilities. In *ICST*, pages 103–114. IEEE, 2025.

[7] Yutao Mou, Xiao Deng, Yuxiao Luo, Shikun Zhang, and Wei Ye. Can you really trust code copilots? evaluating large language models from a code security perspective. *CoRR*, abs/2505.10494, 2025.

[8] Mohammed Latif Siddiq and Joanna CS Santos. Securityeval dataset: mining vulnerability examples to evaluate machine learning-based

code generation techniques. In *Proceedings of the 1st International Workshop on Mining Software Repositories Applications for Privacy and Security*, pages 29–33, 2022.

[9] Hossein Hajipour, Keno Hassler, Thorsten Holz, Lea Schönherr, and Mario Fritz. Codelmsec benchmark: Systematically evaluating and finding security vulnerabilities in black-box code language models. In *SaTML*, pages 684–709. IEEE, 2024.

[10] Yanjun Fu, Ethan Baker, and Yizheng Chen. Constrained decoding for secure code generation. *CoRR*, abs/2405.00218, 2024.

[11] Jinjun Peng, Leyi Cui, Kele Huang, Junfeng Yang, and Baishakhi Ray. Cweval: Outcome-driven evaluation on functionality and security of LLM code generation. In *LLM4Code@ICSE*, pages 33–40. IEEE, 2025.

[12] Mark Vero, Niels Mündler, Victor Chibotaru, Veselin Raychev, Maximilian Baader, Nikola Jovanovic, Jingxuan He, and Martin T. Vechev. Baxbench: Can llms generate correct and secure backends? *CoRR*, abs/2502.11844, 2025.

[13] Xinghang Li, Jingzhe Ding, Chao Peng, Bing Zhao, Xiang Gao, Hongwan Gao, and Xinchen Gu. Safegenbench: A benchmark framework for security vulnerability detection in llm-generated code. *CoRR*, abs/2506.05692, 2025.

[14] Connor Dilgren, Purva Chiniya, Luke Griffith, Yu Ding, and Yizheng Chen. Secrepobench: Benchmarking llms for secure code generation in real-world repositories. *CoRR*, abs/2504.21205, 2025.

[15] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021.

[16] Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. Program synthesis with large language models. *CoRR*, abs/2108.07732, 2021.

[17] Tianyang Liu, Canwen Xu, and Julian J. McAuley. Repobench: Benchmarking repository-level code auto-completion systems. In *ICLR*. OpenReview.net, 2024.

[18] Egor Bogomolov, Aleksandra Eliseeva, Timur Galimzyanov, Evgeniy Glukhov, Anton Shapkin, Maria Tigina, Yaroslav Golubev, Alexander Kovrigin, Arie van Deursen, Maliheh Izadi, and Timofey Bryksin. Long code arena: a set of benchmarks for long-context code models. *CoRR*, abs/2406.11612, 2024.

[19] Manish Bhatt, Sahana Chennabasappa, Cyrus Nikolaidis, Shengye Wan, Ivan Evtimov, Dominik Gabi, Daniel Song, Faizan Ahmad, Cornelius Aschermann, Lorenzo Fontana, Sasha Frolov, Ravi Prakash Giri, Dhaval Kapil, Yiannis Kozyrakis, David LeBlanc, James Milazzo, Aleksandar Straumann, Gabriel Synnaeve, Varun Vontimitta, Spencer Whitman, and Joshua Saxe. Purple llama cyberseceval: A secure coding benchmark for language models. *CoRR*, abs/2312.04724, 2023.

[20] Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Hantian Ding, Ming Tan, Nihal Jain,

Murali Krishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion. In *NeurIPS*, 2023.

[21] Shanchao Liang, Nan Jiang, Yiran Hu, and Lin Tan. Can language models replace programmers for coding? REPOCOD says 'not yet'. In *ACL (1)*, pages 24698–24717. Association for Computational Linguistics, 2025.

[22] Wei Li, Xin Zhang, Zhongxin Guo, Shaoguang Mao, Wen Luo, Guangyue Peng, Yangyu Huang, Houfeng Wang, and Scarlett Li. Feabench: A benchmark for evaluating repository-level code generation for feature implementation. In *ACL (1)*, pages 17160–17176. Association for Computational Linguistics, 2025.

[23] Anthropic. System card: Claude opus 4 & claude sonnet 4. Technical report, Anthropic, May 2025. PDF.

[24] Stephen E. Robertson and Hugo Zaragoza. The probabilistic relevance framework: BM25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, 2009.

[25] Claude 3.5 sonnet model card addendum.

[26] Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll L. Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, and Dane Sherburn. Gpt-4o system card. *CoRR*, abs/2410.21276, 2024.

[27] OpenAI. Model release notes. https://help.openai.com/en/articles/9624314-model-release-notes, May 2025.

[28] OpenAI. Models: codex-mini-latest. https://platform.openai.com/docs/models/codex-mini-latest, May 2025.

[29] xAI. Grok 3 beta — the age of reasoning agents. https://x.ai/news/grok-3, February 2025.

[30] xAI. Grok 4. https://x.ai/news/grok-4, July 2025.

[31] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit S. Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi

Zhang, Kornraphop Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilaï Deutel, Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller, Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero Rissa, Dan Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves, Cindy Wang, Nitesh Bharadwaj Gundavarapu, Ilia Shumailov, Bo Wang, Mantas Pajarskas, Joe Heyward, Martin Nikoltchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik, Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kavukcuoglu, Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov, Bo-Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas Leichner, Haichuan Yang, Zelda Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios Vytiniotis, Jieru Mei, and Mu Cai. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *CoRR*, abs/2507.06261, 2025.

[32] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. Qwen2.5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.

[33] Tencent Hunyuan. Reasoning efficiency redefined! meet tencent's 'hunyuan-t1'—the first mamba-powered ultra-large model. https://tencent.github.io/llm.hunyuan.T1/README_EN.html, March 2025.

[34] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jian Yang, Jiaxi Yang, Jingren Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *CoRR*, abs/2505.09388, 2025.

[35] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, and

Wangding Zeng. Deepseek-v3 technical report. *CoRR*, abs/2412.19437, 2024.

[36] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948, 2025.

[37] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.

[38] Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, et al. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models. *arXiv preprint arXiv:2508.06471*, 2025.