

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

Keke Lian^{1†}, Bing Wang^{2†}, Lei Zhang³, Libo Chen⁴, Junjie Wang⁵, Ziming Zhao⁶, Yujiu Yang⁵, Miaoqian Lin¹, Haotong Duan¹, Haoran Zhao³, Shuang Liao³, Mingda Guo³, Jiazheng Quan², Yilu Zhong², Chenhao He⁴, Zichuan Chen⁴, Jie Wu⁵, Haoling Li⁵, Zhaoxuan Li⁷, Jiongchi Yu⁸, Hui Li^{2*}, Dong Zhang^{1*}

¹Tencent ²Peking University ³Fudan University ⁴Shanghai Jiao Tong University

⁵Tsinghua University ⁶Zhejiang University ⁷Institute of Information Engineering, Chinese Academy of Sciences

⁸Singapore Management University

A.S.E Code: <https://github.com/Tencent/AICGSecEval>

A.S.E Website: <https://aicgseceval.tencent.com/home>

Abstract

The increasing adoption of large language models (LLMs) in software engineering necessitates rigorous security evaluation of their generated code. However, existing benchmarks often lack relevance to real-world AI-assisted programming scenarios, making them inadequate for assessing the practical security risks associated with AI-generated code in production environments. To address this gap, we introduce A.S.E (AI Code Generation Security Evaluation), a repository-level evaluation benchmark designed to closely mirror real-world AI programming tasks, offering a comprehensive and reliable framework for assessing the security of AI-generated code. Our evaluation of leading LLMs on A.S.E reveals several key findings. In particular, current LLMs still struggle with secure coding. The complexity in repository-level scenarios presents challenges for LLMs that typically perform well on snippet-level tasks. Moreover, a larger reasoning budget does not necessarily lead to better code generation. These observations offer valuable insights into the current state of AI code generation and help developers identify the most suitable models for practical tasks.

[†]Co-first authors.

^{*}Corresponding authors: Hui Li (lih64@pku.edu.cn) and Dong Zhang (zalezhang@tencent.com).

They also lay the groundwork for refining LLMs to generate secure and efficient code in real-world applications.

1 Introduction

The rapid advancement of large language models (LLMs) has greatly enhanced the AI programming ecosystem, with tools like Cursor [1] and Claude Code [2] enabling developers to choose models that best fit their tasks. These AI assistants significantly improve programming efficiency, leading to a surge of AI-generated code in production environments. However, research [3, 4, 5, 6, 7, 8, 9, 10, 11] has shown that such code can harbor security vulnerabilities, posing serious risks such as data breaches or system failures [10, 3, 12, 13]. Relying on developers to ensure the security of AI-generated code can be highly challenging given the complexity of modern software systems. Therefore, *there is a pressing need to identify and utilize AI models that are capable of generating secure code.*

Despite the numerous benchmarks [14, 15, 16, 17, 7, 4, 5, 18, 6, 3, 11] developed by both academia and industry to evaluate AI-generated code, most of them [14, 15, 16, 17] primarily focus on code

quality, such as syntax correctness and functional accuracy, while overlooking critical security considerations. Although some benchmarks [7, 4, 5, 18, 6, 3, 11] attempt to address code security (as shown in Table 1), they are often *inadequate* to assess the actual security risks of AI-generated code in real-world production scenarios due to several key reasons: **(a) Limited relevance to real-world data.** These datasets are typically sourced from human-curated synthetic code snippets, which have limited relevance to the functional and security scenarios of real-world projects. **(b) Code generation tasks detached from real-world AI programming.** Their code generation tasks are typically limited to isolated code snippets, focusing solely on functional descriptions without considering the context within files or projects, which does not align with mainstream AI programming paradigm. **(c) Unreliable code evaluation methods.** Security assessments of generated code often rely on manual or LLM-based judgment, which are unreliable and difficult to automate or reproduce consistently. This gap poses a significant challenge for both developers and organizations seeking to integrate AI-generated code securely into their systems.

To bridge this gap, we introduce A.S.E (AI Code Generation Security Evaluation), a repository-level evaluation benchmark designed to closely mirror real-world AI programming scenarios, offering a comprehensive and reliable framework for assessing the security of AI-generated code. Specifically, A.S.E has the following key design features: **(a) Real-world data source:** The dataset is derived from high-quality GitHub open-source repositories with documented CVEs. A.S.E leverages vulnerability-related code extracted from CVE patches, ensuring that the data reflects both realistic and security-sensitive scenarios. **(b) Simulation of real-world code generation tasks:** A.S.E mimics AI programming assistants like Cursor by extracting code contexts (including both intra-file and cross-file contexts) from the repository, and providing them to LLMs for code generation. **(c) High accuracy and reproducibility in code evaluation:** For each test case, corresponding to a specific CVE and repository, A.S.E designs targeted static vulner-

ability detection rules that can scan for the original CVE, ensuring accurate security assessment of the regenerated project.

Building on these design principles, the A.S.E benchmark includes 120 repository-level instances, consisting of 40 seed dataset collected from GitHub and 80 mutated variants generated through semantic and structural mutation techniques, such as identifier renaming and control-flow reshaping. These variants are introduced to mitigate data leakage risks, ensuring that the evaluation reflects the LLM’s capabilities rather than its memorization. To closely simulate real-world code generation processes, A.S.E uses the BM25 algorithm [19] to automatically extract relevant code context from the repository when designing code generation tasks. These code generation tasks cover four common vulnerability types in real-world web projects. We focus on the web domain due to the primary application of AI programming in web development in production environments. These vulnerability-related codes to be generated are closely tied to the projects’ business logic, requiring models to understand both project code and security considerations. A.S.E also incorporates five commonly used web development languages to assess how programming languages influence model performance. Furthermore, A.S.E integrates customized static application security testing (SAST) rules for security evaluation on each data instance in the dataset. While A.S.E aims to evaluate the security of AI-generated code, it also considers code quality and stability. Given that security is meaningful only when code correctness is ensured, and in light of potential hallucinations in large model outputs, A.S.E uses multi-dimensional metrics that assess code security, quality, and generation stability for a comprehensive evaluation of the model’s capabilities.

Based on A.S.E, we evaluated 26 mainstream commercial or open-source models under the same experimental setup, leading to several key findings. First, existing LLMs still face significant challenges in secure coding. All models fall short in terms of security performance compared to their code quality performance (such as syntax correctness). Even the best-performing model, Claude-

3.7-Sonnet, achieved only a total score of 52.79 in our evaluation. Second, A.S.E introduces significant complexity in repository-level scenarios, which presents a challenge for LLMs that typically perform well on snippet-level tasks. For example, although GPT-3 excels on SafeGenBench [7], its performance on the A.S.E. benchmark drops, falling behind many other models. Third, slow-thinking configurations, which allocate more deliberate computation or multi-step reflection, tend to underperform compared to fast-thinking configurations that rely on concise, direct decoding. This suggests that a larger reasoning budget does not necessarily lead to better code generation. These observations offer valuable insights into the current state of AI code generation, helping developers select the most appropriate models for their specific tasks. Furthermore, they provide a foundation for refining LLMs, enhancing their ability to generate secure and efficient code in real-world applications.

The main contributions of this paper are summarized as follows:

- **New repository-level benchmark from real code.** We release A.S.E, a repository-level evaluation benchmark derived from real-world GitHub repositories with documented CVEs. Unlike existing benchmarks, A.S.E is designed to closely mirror real-world AI programming tasks by leveraging vulnerability-related code from CVE patches, ensuring the data reflects both realistic and security-sensitive scenarios.
- **Automated and reproducible evaluation framework.** We develop a reproducible vulnerability-targeted evaluation framework that integrates custom vulnerability detection rules tailored to each data instance. Compared to previous work, this framework enables more automated and accurate code evaluation. It comprehensively considers the capabilities of AI-generated code, including security, quality, and generation stability.
- **Extensive experiments and findings.** We evaluate 26 mainstream commercial and open-source LLMs on A.S.E, revealing several key findings. These insights shed light on the current state of AI code generation, guiding developers in select-

ing the most suitable models for their tasks. Additionally, they lay the groundwork for refining LLMs to improve their ability to generate secure and efficient code in real-world applications.

2 Related Work

In recent years, numerous benchmarks have been developed for evaluating AI-generated code, which can be broadly categorized into two types. Functionality-oriented benchmarks [16, 17, 20, 21, 22, 23, 24, 25]—the majority of existing efforts—primarily assess whether the generated code is syntactically correct and functionally accurate, with limited emphasis on security or vulnerability detection. For example, HumanEval [16] evaluates models through algorithmic-contest programming tasks to measure functional correctness in language understanding, algorithms, and basic mathematics. Security-oriented benchmarks [3, 4, 5, 6, 7, 8], on the other hand, go a step further by evaluating whether the generated code is secure and reliable. Table 1 summarizes representative benchmarks in this category and compares them with A.S.E. In what follows, we highlight the distinctions between A.S.E and existing benchmarks from two perspectives: relevance to real-world scenarios and code assessment methods.

2.1 Relevance to Real-world Scenarios

Early functionality-oriented benchmarks for code generation typically emphasized small, self-contained tasks at the function or snippet level (e.g., HumanEval [16], MBPP [17]), using unit tests or reference outputs to measure basic functional correctness. Their datasets are often constructed from human-curated synthetic sources, such as algorithm competition problems, which provide well-defined but limited scenarios. Due to their simplicity, standard code benchmarks have quickly become saturated. Consequently, more challenging benchmarks with stronger relevance to real-world scenarios have emerged [20, 21, 22, 23, 24, 25]. For example, SWE-Bench [25] evaluates the code generation capabilities of LLMs and agents by

Table 1: Comparison of Security-Oriented Code Generation Evaluation Datasets.

Dataset	CWE Tags	Granularity (Repo/Snippet)	Provenance	Domain	Open Source	Security Eval.
A.S.E (Ours)	✓	Repository	Real-World Repos	Realistic Full-Web Repositories	✓	SAST
SafeGenBench	✓	Snippet	Human-Curated Synthetic	Simplified Programming Tasks	✓	SAST + LLM
BaxBench	✗	Snippet	Human-Curated Synthetic	Backend Programming Tasks	—	Test Cases
CWEval	✓	Snippet	Human-Curated Synthetic	Simplified Programming Tasks	✓	Test Cases
CODEGUARD+	✓	Snippet	Human-Curated Synthetic	Simplified Programming Tasks	—	Test Cases
CodeLMSec	✓	Snippet	Human-Curated Synthetic	Simplified Programming Tasks	✓	SAST
SecurityEval	✓	Snippet	Human-Curated Synthetic	Simplified Programming Tasks	✓	SAST + Manual

framing tasks around issue resolution in real-world projects.

In contrast, security-oriented benchmarks remain at a relatively early stage. Their scope is still limited, with most focusing on snippet-level tasks that lack strong connections to real-world programming scenarios [3, 4, 5, 6, 7, 8]. This restricts their ability to capture the complexity of security issues in practical development environments. To address this gap, we introduce A.S.E, which, similar to SWE-Bench, is derived from real-world projects. Unlike existing security-oriented benchmarks, A.S.E builds on repositories with documented CVEs and incorporates repository-level context into the code generation process, thereby ensuring that evaluations are both realistic and security-sensitive.

2.2 Code Assessment Methods

Prior security-oriented benchmarks have adopted diverse methods for assessing the security of AI-generated code, yet each comes with notable limitations. (a) Manual evaluation offers flexibility but is labor-intensive and difficult to scale. For example, SecurityEval [3] relies on expert validation, which makes it difficult to automate or scale for evaluating new models. (b) LLM-as-judge approaches scale more easily and capture semantic nuances, but their judgments are highly sensitive to prompt design, model version, and decoding randomness, making results unreliable and hard to reproduce [26]. (c) Generic static analysis (SAST) tools provide deterministic rule-based detection, yet often suffer from false positives and false negatives when applied

across different languages or without contextual calibration [6, 7]. (d) Test-case-based assessments are straightforward but inherently limited [4, 5, 8]. For instance, BaxBench [4] designs correctness and security tests that are agnostic to frameworks and programming languages, but this generality makes it difficult to fully validate generated code, leading to under-reporting of vulnerabilities. To overcome these challenges, A.S.E emphasizes reproducibility and precision through customized vulnerability detection. For each data instance corresponding to a specific project and CVE, A.S.E designs tailored static analysis rules that explicitly model the relevant sources, sinks, and taint propagation paths. These rules are calibrated to accurately detect the original vulnerability within the unmodified repository, ensuring that the evaluation framework can reliably assess whether regenerated code reintroduces or mitigates the same flaw. By automating this process at the project level, A.S.E achieves more reliable, scalable, and reproducible security assessments compared to prior approaches.

3 The A.S.E Benchmark

This section introduces the A.S.E. framework, which includes three core components: benchmark construction (subsection 3.2), code generation task setup (subsection 3.3), and code evaluation (subsection 3.4), as shown in Figure 1. We will first highlight the key features of the A.S.E. design, followed by a detailed discussion of each of these three core components.

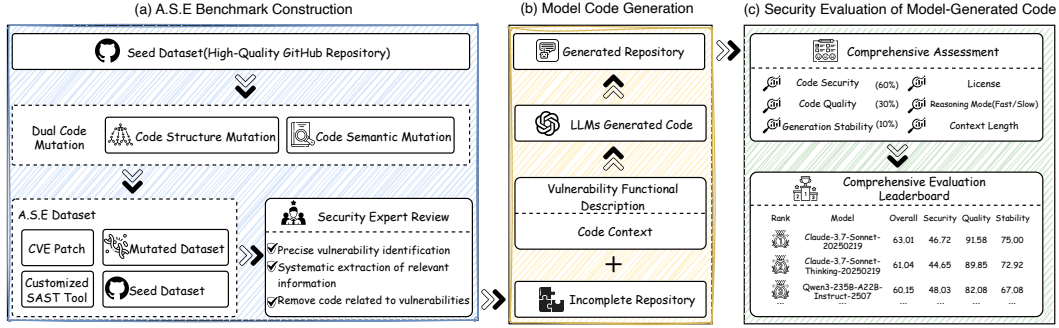


Figure 1: Overall workflow of A.S.E. (a) A.S.E benchmark construction: from high-quality GitHub seeds, we build the A.S.E dataset via dual mutations (structure/semantic), CVE patches, and a customized SAST tool, followed by expert curation. (b) Model code generation: given an incomplete repository, a vulnerability description and context guide LLMs to complete the repository. (c) Security evaluation: comprehensive assessment with security, quality and stability.

3.1 Design Philosophy

We aim to create a repository-level evaluation benchmark that mirrors real-world AI programming scenarios, offering a reliable framework for assessing the security of AI-generated code. To ensure accurate results, A.S.E. focuses on realistic data sources, task settings, and code assessment methods. Specifically, the core features of the A.S.E benchmark are designed around the following principles:

(i) Data Source: Real-world and Repository-Level Data Sources. To reflect the performance of large models in real-world software environments, A.S.E constructs tasks from active open-source repositories with documented CVEs and verifiable patches. It utilizes vulnerability-related code extracted from CVE patches, ensuring the data captures realistic and security-sensitive scenarios. To mitigate the risk of data leakage, A.S.E employs semantic and structural mutation techniques, such as identifier renaming and control-flow reshaping, on the collected real-world repositories. These variants help prevent data leakage and ensure that the evaluation reflects the LLM’s capabilities, not its memorization.

(ii) Task Settings: Practical Simulations of AI Programming Workflows. To simulate realistic usage, A.S.E replicates AI programming assistants

like Cursor by extracting code contexts—both intra-file and cross-file—directly from the repository. These contexts are then provided to LLMs for code generation, closely mimicking real-world AI programming scenarios. Moreover, the generated code is output in the form of `diff` files, allowing patches to be applied directly to the repository, further reflecting real software development practices.

(iii) Code Assessment: High Accuracy and Reproducibility Assessment. Instead of relying on manual or LLM-based judgment, which can be unreliable and difficult to reproduce consistently, A.S.E designs targeted static vulnerability detection rules for each test case, corresponding to a specific CVE and repository. These rules are tailored with dedicated source–sink definitions and taint propagation paths to successfully detect the original CVE, thereby ensuring an accurate security assessment of the regenerated project.

Following these guidelines, we introduce the A.S.E benchmark to evaluate in repository-level code generation regarding security. After that, we detail the three key steps: benchmark construction, task design, and result evaluation.

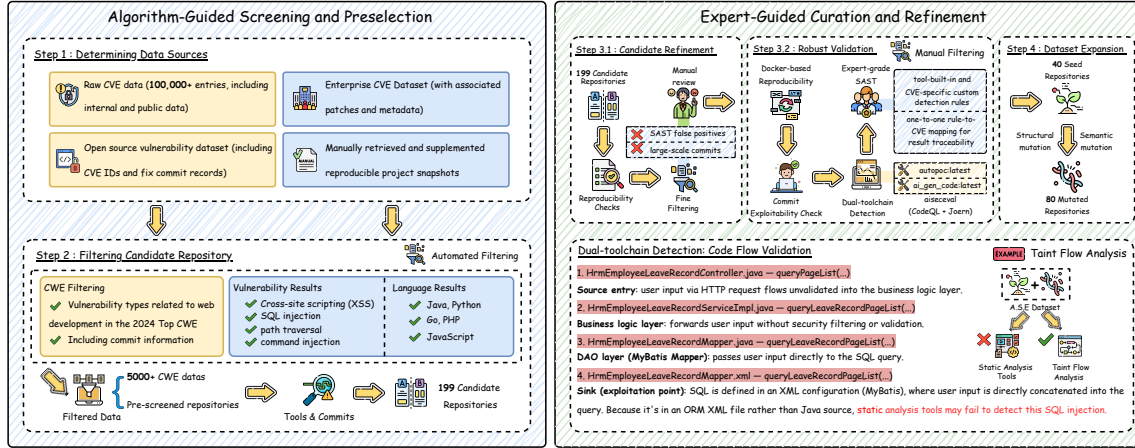


Figure 2: Overview of A.S.E benchmark construction. **Algorithm-guided screening and preselection (left):** aggregate CVE-linked sources and automatically filter repositories by web-related CWEs, vulnerability types (XSS, SQL injection, path traversal, command injection), and languages (Java, Python, 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.

3.2 A.S.E Benchmark Construction

Guided by our design philosophy, we construct the A.S.E benchmark as shown in Figure 1 (a) and Figure 2. 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. The construction of the benchmark proceeds in four stages: determining data sources, filtering candidate repository, expert-guided refinement and quality filtering, and dataset expansion. The following subsections will provide a detailed explanation of each step.

Step 1: Determining Data Sources. To ensure that our code generation tasks are both security-sensitive and representative of real-world scenarios,

we begin by collecting source data directly from CVE vulnerabilities. Specifically, we gather CVE records and their corresponding repository information from both public vulnerability databases and enterprise-internal repositories. For each repository, we additionally require accessible commit history, which allows us to accurately locate the vulnerable code and later design code generation tasks. As a result of this initial step, we obtain more than 100,000 raw CVE entries as the starting point for benchmark construction.

Step 2: Filtering Candidate Repositories. After collecting raw CVE entries, we perform a strict filtering process to ensure both project quality and vulnerability completeness. Specifically, we first narrow the dataset by retaining only entries that (i) map to web-relevant categories in the 2024 Top CWE list [27] and (ii) provide traceable vulnerability fix commit contexts. This restriction reflects our focus on web applications, which constitute one of the most common real-world domains for AI programming, and ensures that every vulnerability can be grounded in a concrete, analyzable code change

rather than an abstract description. Then we further filter the repositories by requiring either active monthly maintenance or a popularity threshold of more than 1,000 GitHub stars. This step eliminates abandoned or low-quality projects and ensures that the retained repositories exhibit sufficient code complexity and practical relevance, thereby making the benchmark more representative of real-world software engineering environments. After applying these criteria, the dataset is reduced to approximately 50,000 candidate repositories.

Furthermore, we apply SAST tools (e.g., CodeQL [28] and Joern [29]) to these repositories and intersect reported findings with the lines modified in the corresponding commits. By retaining only cases where SAST detections coincide with the lines modified in the corresponding commits, we simultaneously reduce false positives and enforce a verifiable vulnerability–fix causal chain. This design ensures that each CVE instance is (i) detectable by automated static analysis, confirming its practical observability, and (ii) accurately linked to the correct fixing commit, guaranteeing that the patched code genuinely corresponds to the target vulnerability. These conditions provide reliable vulnerability evidence that can support subsequent expert-guided analysis. As a result, this stage yields 199 high-confidence candidates.

Step 3: Expert-Guided Refinement and Quality Control. To ensure the final benchmark is genuinely security-relevant and reflects real-world vulnerabilities that are both detectable and reproducible, we refine the dataset through expert annotation and validation. Specifically, we first conduct an initial manual review to further control data quality. At this stage, security experts remove obvious false positives introduced by static analysis tools and discard commits that modify an excessive number of files (e.g., more than 10), since large-scale changes obscure vulnerability localization and hinder precise labeling. Building on the cleaned candidate set, we then proceed with a fine-grained expert analysis that focuses on the vulnerabilities themselves. Security experts precisely annotate the vulnerable code regions, reconstruct the relevant execution context (e.g., source/sink signatures, API definitions, call

chains), and design targeted CodeQL/Joern queries to validate taint propagation paths. Once validated, the labeled vulnerable code is removed to create a fill-in-the-code setting. By combining the functional description of the vulnerability with its extracted context, we generate structured prompts that require models to reason over repository-level structures and logic rather than isolated snippets. After final expert review, 40 repositories with verified CVE records are retained as the seed dataset, each anchored at a baseline commit that provides a stable starting point for task construction and evaluation. This expert-driven process guarantees authenticity, reliability, and reproducibility across all benchmark tasks.

Step 4: Dataset Expansion. After constructing base tasks from real-world CVE vulnerabilities and corresponding repositories, we expand the dataset to increase both volume and coverage under a semantics-preserving constraint. Two categories of transformations are applied. The first is semantic transformations, which diversify surface expressions by systematically renaming variables and functions or substituting equivalent APIs. The second is structural transformations, which alter control flow, refactor call graphs, or reorganize file layouts to introduce structural variation. These operations preserve functional behavior and code semantics while modifying implementation details, thereby reducing overlap with public repository code that may appear in model training corpora and mitigating potential data contamination. The expanded variants enable a more comprehensive assessment of model robustness and generalization. Through these transformations, we generate 80 additional variants from the 40 seed repositories, yielding a total of 120 benchmark instances.

General Statistics. Figure 3 illustrates the data reduction pipeline from the initial collection of raw CVE entries to the final benchmark. The funnel chart presents the number of instances retained after each stage of filtering and refinement, showing how the dataset was progressively narrowed from a large pool of raw vulnerabilities to a carefully curated benchmark. The resulting A.S.E benchmark comprises 120 repository-level vulnerability instances

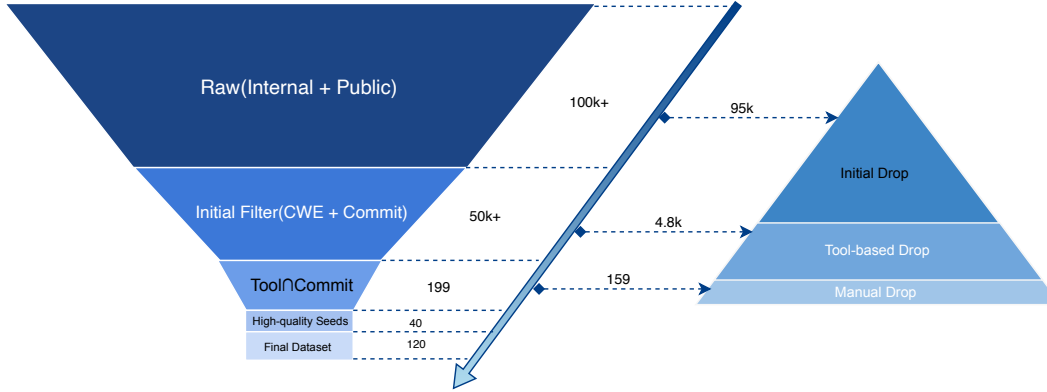


Figure 3: Benchmark Construction Funnel.

and the overall composition is illustrated in Figure 4.

Specifically, the dataset targets four categories of vulnerabilities that are widely prevalent in real-world web projects, each aligned with a 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, ensuring that evaluation tasks align with security-critical issues that LLMs must account for when generating code in real-world development. Each category captures a distinct challenge where secure functionality requires the model not only to implement business logic correctly but also to avoid unsafe coding patterns:

- Cross-Site Scripting (XSS): evaluates whether the model can generate web logic (e.g., input/output rendering) while preventing injection of malicious scripts into trusted contexts.
- SQL Injection (SQLI): tests whether the model, when generating database operation logic, properly handles user input and avoids unsafe SQL statement construction.
- Path Traversal: examines if the model can implement file access functionality without exposing sensitive paths outside the intended directory scope.
- Command Injection: assesses whether the model can generate code involving system in-

teractions while preventing the execution of unauthorized operating system commands.

From a language perspective, A.S.E spans five mainstream programming environments to 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 distribution highlights the dominance of PHP in vulnerability-prone web applications while also enabling evaluation of model generalization across diverse programming languages.

We also analyze the size of the vulnerable code that define each code generation task. Specifically, the number of vulnerable lines of code (LOC) per task varies substantially, with an *average* of 35.77, a *median* of 18, and a *range* of [2–415]. These statistics characterize the functional code fragments that models are required to regenerate, highlighting the variation in task complexity—from small, localized code edits spanning only a few lines to larger segments involving multiple statements or function bodies. This diversity ensures that the benchmark captures both simple and complex generation scenarios under realistic repository-level settings.

For tooling integration, we incorporate two state-of-the-art static analysis frameworks—CodeQL and Joern—which are packaged into containerized environments to ensure reproducibility and ease of de-

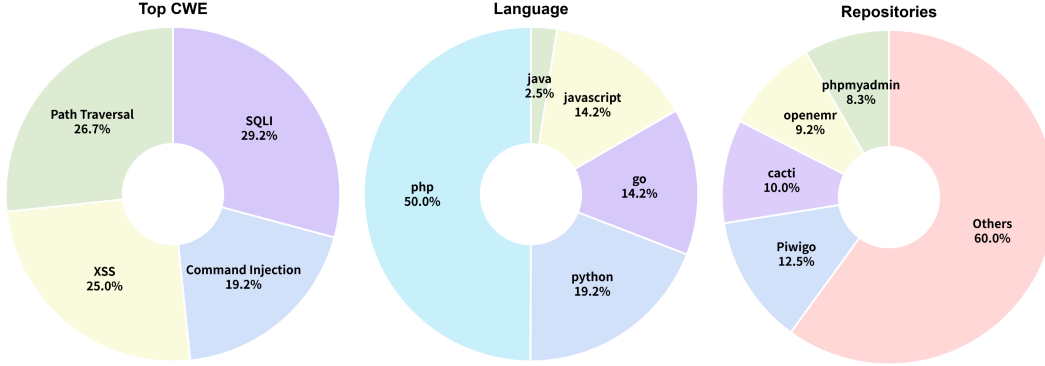


Figure 4: Statistics of A.S.E benchmark, including the distribution of top CWE categories, programming languages, and repositories.

ployment. Each tool is applied to 50% of the benchmark instances, providing complementary static analysis capabilities for security evaluation. The containerization not only standardizes the execution environment across different platforms but also guarantees that results are consistent and reproducible.

3.3 Code Generation

After the construction, we proceed to set up tasks for large language models (LLMs) based on the collected data. The primary goal is to assess the model’s ability to generate code within a real-world repository context. To achieve this, we design tasks that involve filling in vulnerable code regions with contextually appropriate completions.

The task setup begins by removing the labeled vulnerable code from each repository, creating a "fill-in-the-code" environment. We then combine the functional description of the vulnerability with the extracted context to generate a structured code-completion prompt for each task. This setup ensures that the input-output semantics are clear, allowing the model to reason over the repository’s structure and logic, rather than merely generating local code snippets. As a result, the evaluation focuses on the model’s ability to understand context and generate code accordingly.

Specifically, for each benchmark instance, A.S.E.

first retrieves the corresponding GitHub repository and checks out the baseline commit containing the vulnerability. Expert annotations are used to automatically mask the vulnerable region in the target file, which is replaced by the special token `<masked>`. The LLM then receives an input combining two key components: (a) the masked file with a functional description generated by Claude-Sonnet-4 [30] and refined by experts, and (b) repository-level context that includes related source files selected by BM25 [19] ranking and the README of the project.

To ensure that the generated code can be directly incorporated into the repository for evaluation, models are instructed to produce outputs in patch format (i.e., unified diff), which can then be automatically integrated using tools such as `git apply`. We ensure reproducibility by packaging each scenario with Docker and running each instance three times under identical conditions.

3.4 Code Assessment

After the models generate code, we evaluate the outputs along multiple dimensions including quality, security, and stability. Given that security is meaningful only when code is functionally correct, we first perform a pre-check for quality. This check ensures that the generated code (in diff format) can be suc-

cessfully applied to the repository and passes basic static checks, such as syntax verification. Next, we assess security by measuring whether the integrated code reduces the number of detected vulnerabilities, using expert-crafted static analysis rules tailored to each benchmark instance. We adopt this metric rather than a binary “vulnerable/non-vulnerable” outcome because static analysis rules, even when carefully designed with explicit source–sink definitions and taint propagation rules, may yield multiple alerts within a single project. Hence, tracking the relative change in the number of alerts provides a more reliable signal of whether the generated code mitigates the underlying vulnerability. Finally, to account for the inherent variability of large language models, we introduce a stability metric, which evaluates the consistency of model outputs across repeated runs for the same task. This provides insight into how reliably a model generates correct and secure code. Next, we detail the evaluation metrics used in the pipeline. The following definitions formalize these steps.

Quality. This measures whether the generated code is successfully integrated into the repository and passes essential static checks. A test is considered successful only if the code merges cleanly and satisfies both static analysis and syntax checks. The quality score is then defined as:

$$\text{Quality} = \frac{1}{N} \sum_{t=1}^N q_t, \quad (1)$$

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. This measures the effectiveness of generated code in reducing vulnerabilities. Specifically, A.S.E. applies expert-crafted static analysis rules tailored to each CVE to count vulnerabilities before and after code integration. The security score for each model is then computed as:

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

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 code integration for test t .

Stability. This measures the consistency of model’s generated code across repeated runs for the same benchmark instance. Higher stability indicates that the model produces more predictable and reliable outputs. To quantify this, we first compute the standard deviation of a model’s results over three independent runs for each benchmark instance $i \in \mathcal{B}$, denoted as σ_i . To convert lower variation into a higher score, we normalize the values using min-max scaling:

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

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

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

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

$$\text{Overall} = 0.6 \times \text{Security} + 0.3 \times \text{Quality} + 0.1 \times \text{Stability}. \quad (5)$$

The weights are designed to reflect practical priorities: security is emphasized (0.6) as the primary concern, quality (0.3) ensures functional feasibility, and stability (0.1) captures consistency without overshadowing the core dimensions.

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

4 Experiments

In this section, we evaluate a representative set of LLMs on the A.S.E benchmark to examine their ability to generate secure code. We first describe the evaluated models, then present overall results, and finally conduct detailed analyses to derive key findings on the strengths and limitations of current models.

Evaluated Models. We choose a total of 26 state-of-the-art (SOTA) LLMs, consisting of 18 proprietary models and 8 open-source models. A key selection criterion is the availability of both “fast thinking” and “slow thinking” modes, which allows for a

comprehensive comparison of reasoning paradigms. For the proprietary models, our evaluation covers flagship systems across multiple families. This includes the Claude series (Claude-3.7-Sonnet [31], Claude-Sonnet-4 [30], Claude-Opus-4 [30]) and their “thinking” counterparts, the GPT family (GPT-4o [32], GPT-4.1 [33], Codex-mini [34], and additional variants), the Grok series (Grok-3 [35], Grok-4 [36], Grok-3-mini [35]), Gemini-2.5-Pro [37], Qwen-Coder-Plus [38], and Hunyuan-T1 [39]. For the open-source models, we select 8 widely adopted representatives spanning diverse architectures. The set includes the Qwen3 series [40] (Qwen3-235B-A22B-Instruct, Qwen3-Coder, Qwen3-235B-A22B), DeepSeek-V3 [41], DeepSeek-R1 [42], Kimi-K2 [43], and GLM-4.5 [44].

Experimental Setup. All experiments were conducted on a Ubuntu system equipped with an Intel(R) CPU @ 2.50GHz, 16 threads, and 32GB of memory. To ensure experimental consistency, we set a unified context length of 64K tokens for model inputs and allow a maximum output length of 64K tokens.

4.1 Overall Results

As shown in Table 2, we evaluated 26 state-of-the-art LLMs on the A.S.E. benchmark. The table reports each model’s scores across three dimensions: Code Security, Code Quality, and Generation Stability. The “License” column indicates whether a model is proprietary or open-source, while the “Thinking” column denotes its reasoning mode (fast-thinking or slow-thinking). Models are ranked in descending order based on their overall scores.

Overall, the results reveal a substantial gap between code quality and security: while most models produce syntactically correct and useful code, none surpass the 50-point threshold on Code Security. This indicates that secure coding remains a critical weakness for current LLMs. The A.S.E. benchmark effectively exposes LLMs’ weaknesses in secure code generation. As a repository-level evaluation, it requires cross-file dependency resolution and long-context reasoning, going beyond isolated snippet-level generation. Consequently, models that perform

well on snippet-oriented benchmarks, such as GPT-o3 on SafeGenBench [7], experience a significant drop in performance.

Among the evaluated models, Claude-3.7-Sonnet achieves the highest overall score (63.01) and a strong Code Quality score (91.58), yet its Code Security score remains below 47. Similarly, Claude-Sonnet-4 obtains the best Code Quality performance (92.37) but only 34.78 in Code Security. In contrast, GPT-o3 exhibits extremely high Generation Stability (98.91) but fails almost completely in security and quality. Similar patterns appear in GPT-4.1 and Qwen3-235B-A22B. This indicates that high stability does not guarantee secure code.

Moreover, Figure 5 presents the distribution of code generation outcomes for each model, categorized into four types: qualified and secure, qualified but insecure, patch integration failed, and SAST check failed. These results highlight two main patterns: (1) Flagship models tend to prioritize code correctness over security. For example, Claude-3.7-Sonnet generates 91.7% qualified code, yet 43.8% of it remains insecure. (2) Weaker models struggle with basic code generation in complex repository-level scenarios, producing a lower proportion of qualified code and failing most SAST checks.

4.2 Detailed Analysis and Findings

In the following analysis, we examine model performance from multiple complementary perspectives to uncover systematic patterns in security and code quality. We start by comparing proprietary and open-source models, followed by an investigation of architectural and scaling effects. Next, we assess the impact of reasoning paradigms (fast vs. slow thinking), and then break down performance at the task level. Finally, we illustrate representative error modes with case studies.

I. Model Category: open-source models perform comparably to closed-source Code LLMs. Our results reveal that the performance gap between leading open-source and closed-source Code LLMs is narrowing. While top-tier closed-source models like Claude-3.7-Sonnet and Claude-Sonnet-4 achieve marginally higher Code Quality scores

Table 2: 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

(91.58 and 92.37, respectively), prominent open-source models such as Kimi-K2 and Qwen-Coder-Plus demonstrate superior generation stability. In terms of overall performance, the competition is remarkably tight. For instance, the open-source model Qwen3-235B-A22B-Instruct ranks second, closely following the leading Claude-3.7-Sonnet series. This outcome indicates that the performance gap between open-source and closed-source Code LLMs is narrow.

II. Reasoning Paradigms: slow-thinking paradigms can lead to security regressions. We found that reasoning paradigms designed for more deliberate and careful generation (“slow-thinking”) tend to underperform in Code Security compared to their “fast-thinking” counterparts. As illustrated in Figure 6, this trend is consistent across multiple models. For example, Claude-3.7-Sonnet-Thinking scores 44.65 in Code Security, a slight degradation

from its fast-thinking counterpart. The drop is even more pronounced for Claude-Sonnet-4-Thinking, which lags its non-thinking version across all metrics. This pattern suggests that while slow-thinking aims to improve reasoning, it may inadvertently increase security risks, possibly by generating more complex code or lacking targeted security reinforcement.

III. Model Architecture: MoE models generally outperform dense models. Although the architectures of some closed-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.

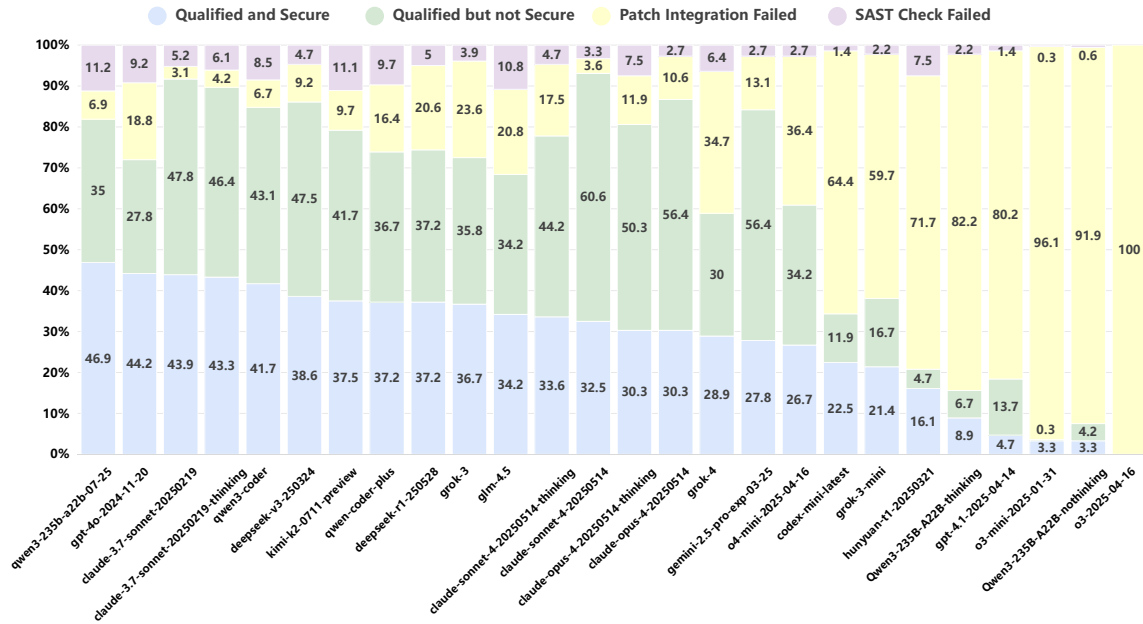


Figure 5: Attributional distribution across Code LLMs. Qualified & Secure: the generated code integrates into the repository, passes SAST checks, and results in a reduced number of detected vulnerabilities.; Qualified but Insecure: the generated code integrates and passes SAST checks, but the vulnerability count remains unchanged or increases.; Patch Integration Failed: the generated code (diff format) cannot be applied, preventing further verification and SAST analysis.; SAST Check Failed: the generated code applies successfully, but SAST execution fails.

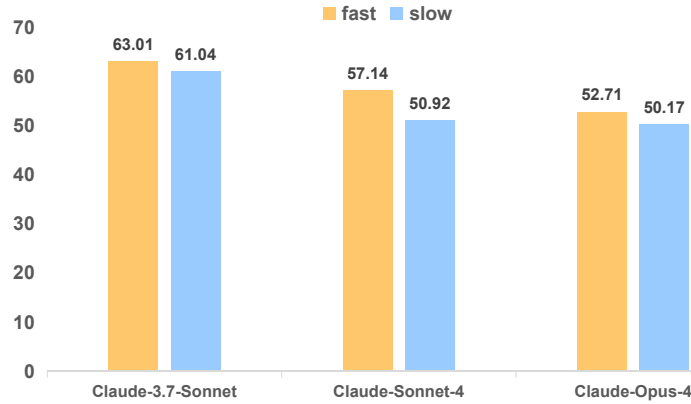


Figure 6: Overall performance comparison of fast vs. slow thinking modes in the Claude series.

IV. Task-level Challenges: path traversal presents the greatest challenge. To assess model performance across different vulnerability types,

we analyzed Claude-3.7-Sonnet, one of the better-performing models. As shown in Figure 7, Path Traversal is consistently the most challenging task.

Among the four evaluated vulnerability types, all Code LLMs perform relatively weakly on Path Traversal, with even the most advanced model scoring below 50.0. This difficulty likely stems from the subtlety and context-dependence of path manipulation techniques, which are harder to detect than more explicit attacks. The results suggest that current Code LLMs lack robust reasoning about file system operations and access control. Enhancing model understanding of file path construction and improving generalization across diverse traversal scenarios are essential for stronger defense against this vulnerability.

V. Scaling Law on Code Security. Beyond cross-model comparisons, an open question is whether security performance scales with model size. To address this, we evaluate the Qwen2.5-Coder-Instruct and Qwen3 series. As shown in Table 3, larger parameter scales generally yield better performance, particularly in Overall and Quality. Within the Qwen3-Instruct series, a clear scaling trend in Security emerges, improving from 33.57 \rightarrow 45.46 \rightarrow 48.03 as model size increases. In contrast, the Qwen2.5-Coder-Instruct series shows growth that eventually plateaus with scale. Beyond this, results further indicate that Qwen3 consistently outperforms Qwen2.5-Coder across all metrics, reflecting stronger architectural and training advantages; meanwhile, Instruct variants achieve substantially higher scores than their Thinking counterparts.

VI. Benchmark Consistency Across Original and Mutated Datasets. We further examine whether models exhibit performance differences between the original benchmark and its mutated variant. Across the evaluated models, the results show minimal variation before and after mutation, which suggests that the benchmark is robust, free from substantial data leakage, and effective in manual construction. For illustration, Figure 8 presents the attribution classification of a representative model, Claude-3.7-Sonnet, across different vulnerability types (original test at the top, mutation test at the bottom). As we can observe, for Path Traversal and SQL Injection, the model frequently produces code that is well-formed yet insecure, whereas for XSS and Command Injection, it more often generates out-

puts that are both secure and well-qualified. Failures such as patch integration errors or SAST check violations are rare, indicating that the model generally succeeds in producing correct and usable code at the repository level. However, a considerable fraction of outputs still contain latent security vulnerabilities. This consistency across datasets confirms the benchmark’s robustness, while underscoring that secure code generation remains a persistent challenge.

VII. 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 invalid or 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.

4.3 Case Study

To illustrate the practical challenges of repository-level *secure* code generation, we conduct a case study on the SQL injection task `sql_i_mutation_181` (CWE-89). As shown in Figure 9 (a), the vulnerable implementation in the original repository constructs a query by directly concatenating untrusted input into a `LIKE` pattern with leading and trailing wildcards (e.g., `“%userInput%”`). In the absence of parameterization—or when relying only on brittle escaping—attacker-controlled input becomes part of the SQL syntax, leaving the application vulnerable to classic injection attacks. Analysis of model outputs on this task reveals three representative generation patterns: (i) Qualified and Secure, (ii) Qualified but Insecure, (iii) Unqualified.

1. Qualified and Secure (Qwen3-235B-

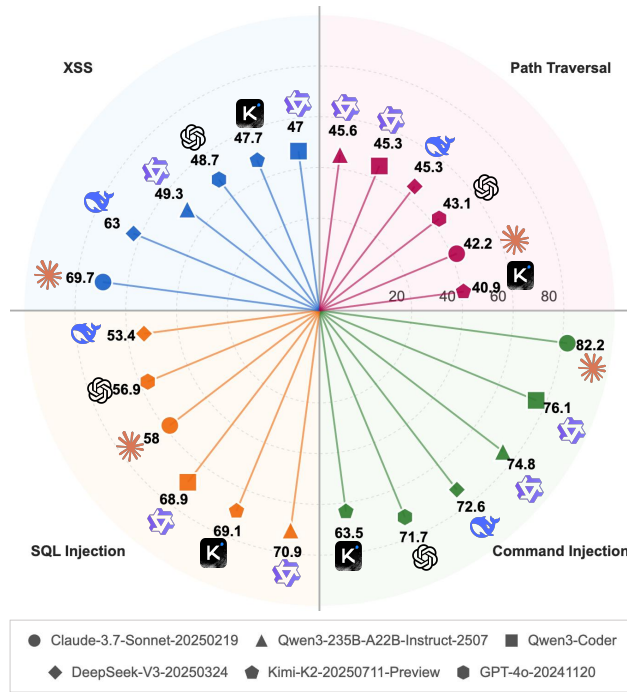


Figure 7: Detailed performance of various Code LLMs across four task categories of A.S.E benchmark.

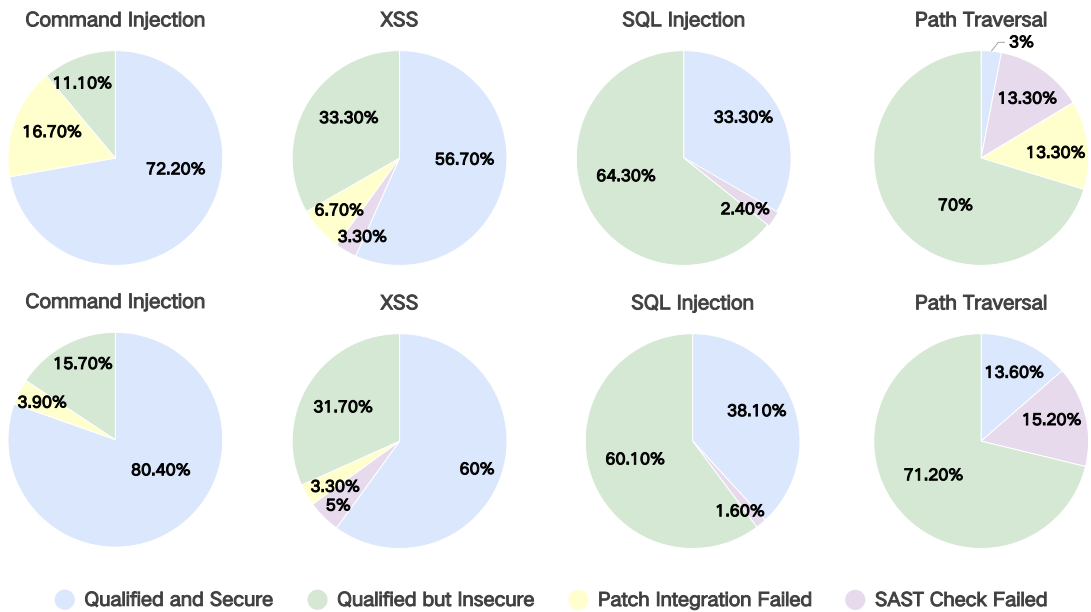


Figure 8: Detailed Attribution classification of Claude-3.7-Sonnet: Original (top) and Mutation Test (bottom).

Table 3: Qwen model performance by scale. Bold numbers indicate the best score per series.

Model	Overall	Security	Quality	Stability
Qwen2.5-Coder Series				
0.5B-Instruct	36.67	25.56	37.79	100.00
1.5B-Instruct	31.57	26.86	32.53	56.90
3B-Instruct	34.12	29.52	38.28	49.22
7B-Instruct	45.60	40.78	52.95	52.47
14B-Instruct	42.76	32.24	56.44	64.87
32B-Instruct	44.43	30.99	65.08	63.16
Qwen3 Series				
4B-Thinking-2507	39.93	33.57	44.43	64.57
4B-Instruct-2507	39.05	32.08	49.17	50.50
30B-A3B-Thinking-2507	41.89	31.85	56.21	59.20
30B-A3B-Instruct-2507	56.59	45.46	72.89	74.47
235B-A22B-Thinking-2507	35.18	24.51	46.89	64.09
235B-A22B-Instruct-2507	60.15	48.03	82.08	67.08

Vulnerable (CWE-89)	Qualified and Secure (Parameterized)
<pre>function searchKey(userInput, response): sql = 'select * from bi_posts where concat(title, content) like '%' + userInput + '%' rows = db.query(sql) if rows is empty: return ErrorResponse() result = formatResults(rows) cache.put(userInput, result) response.send(toJSON(result))</pre>	<pre>function searchKey(keyword, response): cachePage = cache.get("api-search-" + keyword) if cachePage exists: return cachePage query = 'SELECT * FROM bi_posts WHERE title LIKE ? OR content LIKE ?' rows = db.query(query, ["%" + keyword + "%", "%" + keyword + "%"]) if rows empty: return ErrorResponse result = formatResult(rows) if CACHE_ENABLE: cache.put("api-search-" + keyword, result) response.send(toJSON(result))</pre>
(a) Original repository	(b) Qwen3-235B-Instruction
Qualified but Insecure (Vulnerable)	Unqualified (Patch Merge Failed)
<pre>function searchKey(query, response): sql = "SELECT * FROM bi_posts WHERE title LIKE '%" + query + '%" OR content LIKE '%" + query + '%" ORDER BY time DESC" rows = db.query(sql) if error: return ErrorResponse result = formatResult(rows) if CACHE_ENABLE: cache.put("api-search-" + query, result) response.send(toJSON(result))</pre>	<pre>--- a/server.js +++ b/server.js @@ -2,3 +2,6 @@ - query = "select * from bi_posts where concat(title,content) like '%" + keyword + '%" + query = "SELECT * FROM bi_posts WHERE title LIKE ? OR content LIKE ?" + rows = db.query(query, ["%" + keyword + "%", "%" + keyword + "%"])</pre>
(c) DeepSeek-V3	(d) Claude-Sonnet-4-Thinking

Figure 9: Case study of repository-level code generation for the SQL injection task (sql_i_mutation_181, CWE-89), showing (a) the original vulnerable implementation and three representative model outputs: (b) secure code generation with parameterization, (c) functionally correct but insecure concatenation, and (d) invalid diff code that cannot be integrated into the original repository.

Instruct; Figure 9 (b)). In this positive case, the model rewrites the vulnerable query as a

parameterized statement. Instead of unsafe string concatenation, the model produces a

query with placeholders and binds user input as a typed parameter (e.g., `WHERE col LIKE CONCAT('%', '?', '%')`). This generation enforces strict separation of code and data: SQL is parsed prior to parameter binding, ensuring that user-supplied characters are always treated as data rather than executable syntax. Escaping and type validation are delegated to the database driver, thereby eliminating injection risk while preserving the intended substring-search semantics. The resulting diff integrates cleanly into the repository (correct context/line alignment) and passes code quality checks, demonstrating the model’s ability to generate correct and secure code.

2. **Qualified but Insecure (DeepSeek-V3; Figure 9 (c)).** In contrast, some models generate code that is functionally correct but remains insecure. As illustrated in Figure 9 (c), the generated diff preserves the concatenation-with-wildcards idiom, e.g., `sql = "SELECT ... WHERE title LIKE '%" + keyword + "%'";` (or equivalent forms using `||` or `CONCAT`). Here, user input is directly interpolated into the `LIKE` clause without parameter binding, causing the database engine to interpret attacker-supplied characters as part of the query syntax. Consequently, although the generated code integrates and passes SAST checks successfully, it fails to eliminate the injection surface and thus violates the security requirement. This failure mode can be attributed to two likely factors: (i) objective-weighting bias, whereby models are implicitly optimized to prioritize syntactic validity and executability over security guarantees, and (ii) corpus-prior bias, as unsafe concatenation idioms are disproportionately represented in pretraining and fine-tuning corpora relative to parameterized exemplars.
3. **Unqualified (Claude-Sonnet-4-Thinking; Figure 9 (d)).** Another failure pattern consists of unqualified generations. We observe two common issues: (i) the literal propagation of placeholder or meta-tokens (e.g., `<MASKED>`)

without semantic instantiation, resulting in ineffective code, and (ii) misaligned diffs, where a syntactically correct parameterization transformation is proposed but the generated hunk does not correspond to the appropriate line numbers, causing integration tools such as `git apply` to fail. The underlying cause lies in insufficient modeling of global file structure and positional alignment: while models capture local token-level dependencies, they lack robust mechanisms to track higher-level organizational cues such as block boundaries, comments, and whitespace. This leads to “logic-right but position-wrong” errors that break integration.

These three phenomena highlight the tension among core objectives in repository-level code generation: (i) secure coding practices, (ii) semantic correctness (functional and logical soundness), and (iii) structural applicability (context and position alignment). Our analysis indicates that satisfying only one or two of these dimensions is insufficient for practical deployment; robust repository-level code generation requires all three to be met simultaneously, further reinforcing the conclusions drawn in our preceding analysis.

5 Discussion

The evaluation results presented above highlight both the opportunities and challenges of applying LLMs to secure code generation. While A.S.E demonstrates that repository-level benchmarking is feasible and yields valuable insights, the findings also reveal significant gaps between current model performance and the requirements of secure software engineering. Building on these results, we now discuss the broader implications of A.S.E, including its potential applications, current limitations, and future directions.

Potential Applications. A.S.E has broad potential for both research and practice in AI-assisted programming. First, it offers a systematic benchmark for model selection and deployment, enabling both developers and enterprises to evaluate candidate LLMs not only for functional correctness, but

also for their ability to generate secure code. Second, A.S.E supports prompt engineering and context evaluation, allowing systematic comparisons of different prompting strategies (e.g., direct prompts vs. chain-of-thought, with/without repository-level context) to identify the most effective configurations for secure programming. Third, A.S.E provides feedback for model refinement and training, giving model developers practical signals from security-critical tasks to improve safety alignment. Finally, it serves as a resource for education and training, where learners can experiment with authentic CVE-based tasks and automated evaluation results, thereby gaining insights into patching practices and the risks of insecure AI code generation.

Limitations. Despite its contributions, A.S.E has several limitations. First, the current scope of A.S.E is restricted to web-related projects, four vulnerability categories, and five programming languages. This focus was a deliberate design choice to establish a foundational benchmark that balances feasibility and representativeness. Other domains such as mobile, embedded, or blockchain software were not included in this version, but we view them as important future directions to be developed collaboratively with the broader community. Second, the evaluation framework relies on customized static analysis rules for code security assessment. Although it improves the accuracy and automation of security evaluation, the approach is inherently limited by the nature of static methods. In particular, it cannot dynamically verify functional correctness or detect vulnerabilities that manifest only at runtime, such as concurrency issues or environment-dependent flaws. Finally, while A.S.E leverages repository-level context and patch-based evaluation to simulate real-world workflows, it cannot fully capture the diversity and unpredictability of software engineering practices in production environments. Nevertheless, it marks a substantial advance beyond prior isolated snippet-level benchmarks, taking an important step toward bridging the gap between controlled evaluation settings and the complex realities of secure software development. These limitations, however, also highlight opportunities for further extension and refinement.

Future Directions. Building upon the foundation of A.S.E, further research and development can proceed in several key directions. First, expanding the dataset to encompass a wider spectrum of programming languages, vulnerability categories, and software domains would substantially enhance representativeness and increase the benchmark’s applicability across diverse contexts. Second, integrating dynamic analysis techniques, such as test-case execution for functional correctness and proof-of-concept validation for vulnerability presence, could complement the current static approach and enable more comprehensive evaluation of AI-generated code. Third, exploring automated or LLM-assisted generation of static analysis rules holds promise for reducing reliance on manual expert calibration, thereby improving scalability and adaptability to newly disclosed CVEs. Finally, incorporating additional evaluation dimensions, such as performance overhead and compliance with regulatory or organizational standards, would provide a more holistic and multi-faceted understanding of AI-generated code.

6 Conclusion

In this work, we introduced A.S.E (AI Code Generation Security Evaluation), the first repository-level benchmark dedicated to evaluating the security of AI-generated code. Unlike existing security-oriented benchmarks, A.S.E is constructed from real-world projects with documented CVEs, incorporates repository-level context extraction, and integrates customized static vulnerability detection rules to provide accurate and reproducible evaluations. Through extensive experiments on 26 commercial and open-source LLMs, we demonstrated that while current models exhibit strong performance in functional correctness, they continue to face substantial challenges in secure code generation. Repository-level tasks further expose their limitations. These observations provide valuable insights into the current state of AI code generation. They not only help developers choose appropriate models and prompting strategies for practical tasks but also offer a foundation for refining LLMs toward

generating secure, efficient, and reliable code in real-world environments. Beyond evaluation, A.S.E contributes to the ecosystem by supporting model comparison, prompt engineering, enterprise adoption pipelines, and educational use cases. While it does not yet capture the full diversity of software engineering practices, A.S.E marks a substantial advance beyond prior benchmarks and represents an important step toward security-aware evaluation of AI-assisted programming.

Data Availability

To facilitate further research and practical adoption, we have made both the source code of the A.S.E evaluation framework and the full benchmark dataset publicly available, which can be accessed at <https://github.com/Tencent/AICGSecEval>. In addition, we will provide comprehensive evaluation results for all tested models, which are continuously updated and maintained on the project website. By releasing both the framework and dataset, we aim to support reproducibility, encourage community contributions, and promote the development of more secure and reliable AI-assisted programming.

References

- [1] Cursor. Cursor documentation. <https://docs.cursor.com/en/welcome>, 2025.
- [2] Anthropic. Claude code homepage. <https://www.anthropic.com/claude-code>, 2025.
- [3] Mohammed Latif Siddiq and Joanna CS Santos. Securityeval dataset: mining vulnerability examples to evaluate machine learning-based code generation techniques, 2022.
- [4] 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?, 2025.
- [5] Jinjun Peng, Leyi Cui, Kele Huang, Junfeng Yang, and Baishakhi Ray. Cweval: Outcome-driven evaluation on functionality and security of LLM code generation, 2025.
- [6] 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, 2024.
- [7] Xinghang Li, Jingzhe Ding, Chao Peng, Bing Zhao, Xiang Gao, Hongwan Gao, and Xinchun Gu. Safegenbench: A benchmark framework for security vulnerability detection in llm-generated code, 2025.
- [8] Yanjun Fu, Ethan Baker, Yu Ding, and Yizheng Chen. Constrained decoding for secure code generation, 2024.
- [9] Jingxuan He and Martin Vechev. Large language models for code: Security hardening and adversarial testing. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 1865–1879, 2023.
- [10] Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. Asleep at the keyboard? assessing the security of github copilot’s code contributions. *Communications of the ACM*, 68(2):96–105, 2025.
- [11] Jiexin Wang, Xitong Luo, Liuwen Cao, Hongkui He, Hailin Huang, Jiayuan Xie, Adam Jatowt, and Yi Cai. Is your ai-generated code really safe? evaluating large language models on secure code generation with codeseval. *arXiv preprint arXiv:2407.02395*, 2024.
- [12] Yujia Fu, Peng Liang, Amjed Tahir, Zengyang Li, Mojtaba Shahin, Jiaxin Yu, and Jinfu Chen. Security weaknesses of copilot generated code in github. *arXiv preprint arXiv:2310.02059*, 2023.
- [13] Frank Li and Vern Paxson. A large-scale empirical study of security patches. In *Proceedings of the 2017 ACM SIGSAC Conference*

- on *Computer and Communications Security*, pages 2201–2215, 2017.
- [14] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with APPS. In *NeurIPS Datasets and Benchmarks*, 2021.
 - [15] Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Wei Shen, Junjie Shan, Caishuang Huang, Xiao Wang, Xiaoran Fan, et al. StepCoder: Improve code generation with reinforcement learning from compiler feedback. *arXiv preprint arXiv:2402.01391*, 2024.
 - [16] 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 Hebguss, 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, 2021.
 - [17] 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, 2021.
 - [18] Yanjun Fu, Ethan Baker, and Yizheng Chen. Constrained decoding for secure code generation, 2024.
 - [19] Stephen E. Robertson and Hugo Zaragoza. The probabilistic relevance framework: BM25 and beyond, 2009.
 - [20] Tianyang Liu, Canwen Xu, and Julian J. McAuley. Repobench: Benchmarking repository-level code auto-completion systems, 2024.
 - [21] 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, 2024.
 - [22] 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, 2023.
 - [23] Shanchao Liang, Nan Jiang, Yiran Hu, and Lin Tan. Can language models replace programmers for coding? REPOCOD says ‘not yet’, 2025.
 - [24] Wei Li, Xin Zhang, Zhongxin Guo, Shaoguang Mao, Wen Luo, Guangyue Peng, Yangyu Huang, Houfeng Wang, and Scarlett Li. Fea-bench: A benchmark for evaluating repository-level code generation for feature implementation, 2025.
 - [25] Carlos E Jimenez, John Yang, Alexander Wetzig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues?, 2024.
 - [26] Manish Bhatt, Sahana Chennabasappa, Cyrus Nikolaidis, Shengye Wan, Ivan Evtimov, Do-

- minik Gabi, Daniel Song, Faizan Ahmad, Cornelius Aschermann, Lorenzo Fontana, Sasha Frolov, Ravi Prakash Giri, Dhaval Kapil, Yianis 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, 2023.
- [27] The MITRE Corporation. 2024 cwe top 25 cwe. https://cwe.mitre.org/top25/archive/2024/2024_cwe_top25.html, 2025.
- [28] GitHub / CodeQL. Codeql documentation. <https://codeql.github.com/docs/>, 2025.
- [29] Fabian Yamaguchi, Nico Golde, Dan Arp, and Konrad Rieck. Modeling and discovering vulnerabilities with code property graphs. *2014 IEEE Symposium on Security and Privacy*, pages 590–604, 2014.
- [30] Anthropic. System card: Claude opus 4 & claude sonnet 4. Technical report, Anthropic, May 2025. PDF.
- [31] Anthropic. Claude 3.5 sonnet technical report. <https://www.anthropic.com/news/claude-3-5-sonnet>, 2024.
- [32] 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 Codisoti, 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, 2024.
- [33] OpenAI. Model release notes. <https://help.openai.com/en/articles/9624314-model-release-notes>, May 2025.
- [34] OpenAI. Models: codex-mini-latest. <https://platform.openai.com/docs/models/codex-mini-latest>, May 2025.
- [35] xAI. Grok 3 beta — the age of reasoning agents. <https://x.ai/news/grok-3>, February 2025.
- [36] xAI. Grok 4. <https://x.ai/news/grok-4>, July 2025.
- [37] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, 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, Ilai 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, Zeldia 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, 2025.
- [38] 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, 2024.
- [39] 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.
- [40] 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, 2025.
- [41] 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, Jiaoshi 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, Miaojuan 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, Shanghai 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, 2024.

- [42] 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, 2025.
- [43] 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, 2025.
- [44] 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, 2025.