



TOZI JOURNAL OF ENGINEERING SCIENCES JOURNAL OF ENGINEERING SCIENCES

Security Evaluation of AI-Generated Code: A Comparative Study of ChatGPT, Copilot, And Gemini through Static and Dynamic Analysis

Onur Ceran*a

Submitted: 16.06.2025 Revised: 04.07.2025 Accepted: 09.08.2025 doi:10.30855/gmbd.070525A13

ABSTRACT

Keywords: Generative AI, ChatGPT, Copilot, Gemini, software security, static code analysis, dynamic code analysis

This study examines the security performance of generative artificial intelligence (AI) tools of ChatGPT, Copilot, and Gemini within software development workflows. Through static and dynamic code analysis, security vulnerabilities in web application login code generated by these tools were systematically evaluated. Results indicate that while AI models offer efficiency in code generation, they also introduce varying levels of security risk. Copilot exhibited the highest cumulative risk with multiple high-level vulnerabilities, while ChatGPT demonstrated a lower risk profile. Gemini produced relatively optimized code but contained critical security flaws that require manual review. The most common vulnerabilities across all models were insecure design and security logging and monitoring failures, indicating a systemic issue in AI-generated code. The findings emphasize that generic prompts focusing on security are insufficient and that developers must use specific, securityoriented prompts, such as applying secure-by-design principles and implementing OWASP Top Ten protections. This study contributes to the growing body of literature addressing the security implications of integrating AI into software development, highlighting the importance of human oversight and carefully crafted prompts to mitigate potential risks.

1. Introduction

Recent advancements in artificial intelligence (AI) have instigated a paradigm shift in computational processing, transitioning from conventional data-driven, discriminative tasks to sophisticated, creative applications facilitated by generative AI. This pivotal shift has drastically altered the long-standing belief that artistic and creative tasks like writing poems, creating software, designing fashion, and composing songs are exclusively human capabilities, as AI can now generate indistinguishable new content. Generative AI specifically refers to computational techniques capable of creating seemingly new, meaningful content like text, images, or audio from training data, serving artistic purposes as well as assisting humans as intelligent question-answering systems [1]. Generative AI has recently emerged as a novel tool, presenting a wide range of new possibilities across diverse sectors from education and healthcare to networked businesses [2], while simultaneously ushering in significant changes to code writing through its integration into software development processes [3]. Large Language Models (LLMs), such as ChatGPT, GitHub Copilot, and Google Gemini, have become widely used for automating and enhancing programming tasks. Thanks to their natural language processing (NLP) capabilities, these models can successfully perform tasks such as suggesting complex code snippets, debugging, and generating explanations [4].

Traditionally, code found online is often considered reliable, well-documented, and easily adaptable with minimal changes, like renaming variables. However, LLMs offer a valuable initial framework for users, even if the generated code isn't perfect. This is particularly beneficial for programmers facing challenges or unsure how to begin a task, making LLMs a popular tool for them [5]. The study by Wang et al. [6] highlights that software developers tend to place high levels of trust in code generated by generative AI tools, largely based on their prior experiences and the credibility they associate with online communities like Stack Overflow. However, the proliferation of AI models

a.* Gazi University, Faculty of Applied Sciences, Dept. of Management Information Systems 06560 - Ankara, Türkiye, Orcid: 0000-0003-2147-0506 * Corresponding author: onur.ceran@gazi.edu.tr

capable of generating code introduces considerable complexities concerning reliability, safety, and security [7]. Unlike simple classification tasks, evaluating the correctness and safety of complex software artifacts requires a more sophisticated and nuanced assessment framework. Although AI tools such as ChatGPT, GitHub Copilot, and Google Gemini provide software developers with efficient code suggestions, saving time and effort [8], the security of the code generated by these models remains a critical concern. Potential issues, including security vulnerabilities, licensing constraints, and the risk of data leakage, must be carefully considered alongside the benefits these tools provide. Moreover, AI models based on natural language processing are prone to "hallucination," a phenomenon in which the system produces content that seems credible but is actually inaccurate or misleading [9,10].

Research indicates that a considerable portion of the code produced by these tools contains security vulnerabilities. For instance, it has been reported that up to 40% of code generated by GitHub Copilot may include security flaws [11]. Furthermore, LLM-based tools rarely recommend sophisticated security validation mechanisms, which in turn increases users' reliance on manual security checks [12]. Another notable challenge lies in constructing authentic and valid datasets required for training LLMs [4]. The licensing and ethical concerns arising from the datasets used for training these models also warrant careful attention [13]. Essentially, LLMs are trained on vast corpora of text, often incorporating extensive amounts of code and related discussions. This comprehensive exposure enables models to implicitly learn the patterns, structures, and conventions of various programming languages, rather than explicitly acquiring knowledge of syntax or semantic rules [14]. However, training LLMs on open-source code repositories exposes users to the risk of licensing violations. Moreover, Alassisted tools pose the potential risk of suggesting faulty or malicious code, which may lead to severe consequences, especially in critical sectors such as finance and healthcare [9].

Developers also handle a wide range of tasks when creating and maintaining software, and they sometimes inadvertently introduce insecure code patterns. The challenges associated with security bugs involve several key areas: detecting them, localizing their root causes, understanding these underlying issues, and then creating and thoroughly testing effective patches [15]. Human limitations in thoroughly appraising software quality have historically led to the indispensable adoption of practices such as code reviews and the application of static and dynamic analysis techniques to uphold the integrity of human-authored code. These established methodologies underscore the heightened challenges in validating AI-generated code [16]. In this context, a comparative security analysis of code generated by ChatGPT, Copilot, and Gemini is crucial for understanding existing vulnerabilities and assessing the broader implications of these tools on software security. The primary objective of this study is to examine the security of AI-generated code. The research begins with a comprehensive evaluation of the capabilities of LLMs to produce secure code and subsequently investigates the presence of vulnerabilities in the code generated by the three AI tools mentioned.

The remainder of this study is structured as follows: Section 2 presents background and a review of relevant literature. Section 3 details the research methods, data collection procedures, analysis techniques, and evaluation criteria. Section 4 presents and interprets the research findings. Finally, Section 5 provides the study's conclusions and discussion.

2. Background and Related Work

The landscape of Artificial Intelligence (AI) has undergone a significant transformation with the emergence of LLMs, particularly OpenAI's ChatGPT [17]. ChatGPT is a conversational interface powered by a LLM developed by OpenAI, based on the Generative Pre-trained Transformer (GPT) architecture. Initially released in November 2022, ChatGPT currently operates using the GPT-4 framework, which has been trained on extensive and diverse textual data. As an LLM, it is capable of understanding and generating human-like language, and it can effectively perform various natural language processing tasks, such as text generation, summarization, translation, and problem-solving across multiple domains. Through advanced deep learning techniques, ChatGPT captures linguistic structures and semantic contexts to provide coherent and contextually relevant responses. Its expanding use in academic research, education, and digital communication demonstrates its growing impact as a significant AI-driven language processing tool [18]. Gemini, developed by Google, is a family of LLMs. Its advanced architecture enables it to process and comprehend various forms of information, including text, code, audio, image, and video data. The initial release of Gemini occurred in December 2023, with subsequent updates introducing expanded capabilities. Its development focuses on achieving high performance across a broad spectrum of benchmarks, showcasing proficiency in areas such as natural language understanding, code generation, mathematical reasoning, and creative content creation. As a foundational model, Gemini aims to provide robust capabilities for diverse applications, thereby contributing to both research and practical implementations within advanced AI systems [19]. Microsoft Copilot, introduced in 2023, is a groundbreaking AI tool that acts as an advanced assistant to boost human productivity, especially in writing and problem-solving. A core feature of Copilot is its real-time generation of contextually appropriate responses. It understands the conversation and user intent, which helps speed up chats and lessen the mental effort for users. This is beneficial for both new users needing guidance and experienced users aiming for more efficient conversations. Beyond just chatting, Copilot also supports creative endeavors like writing poems, stories, and other text-based content [20].

The integration of generative AI into software development processes has significantly enhanced speed and efficiency in software production [3]. However, this advancement has concurrently increased security-related risks. In a user study, Perry et al. [21] investigated the ways in which individuals interact with AI code assistants when addressing various security-related programming tasks. Their research revealed that participants who utilized an AI assistant generated code that was notably less secure compared to those who did not have access to such a tool. Furthermore, the study indicated that users interacting with AI assistants were more inclined to believe their code was secure, which suggests that these tools may inadvertently foster overconfidence regarding potential security vulnerabilities within the generated code. Studies evaluating the security and quality of code generated by AI tools such as ChatGPT, GitHub Copilot, and Google Gemini indicate that these models are assessed from diverse perspectives.

Kharchenko and Babenko [22] conducted a comparative analysis of LLMs across various application domains, focusing on the accuracy and reliability of the generated code. Their findings suggest that Copilot achieved higher scores in terms of data accuracy and security, while ChatGPT demonstrated superior performance in terms of flexibility. Tihanyi et al. [23] conducted an extensive study evaluating the propensity of state-of-the-art LLMs to introduce vulnerabilities when generating C programs. Employing a neutral zero-shot prompt, their large-scale investigation, which involved nine cutting-edge models, revealed that a substantial proportion of the generated code, at least 62.07%, contained vulnerabilities. While minor variations existed among the models, they generally exhibited similar types of coding errors. Security-focused studies, including those by the prominent cybersecurity company TrendMicro [24], emphasize that trusting code generated by ChatGPT could potentially result in the deployment of insecure code in production environments and may inadvertently introduce security vulnerabilities. Similarly, a 2021 study reported that GitHub Copilot produced security-related issues in approximately 40% of its code outputs. A study by Kharma et al. [25], which analyzes the security of code generated by LLMs across various programming languages, indicates that Gemini models exhibit the highest tendency to produce vulnerabilities in Java code when compared to other LLMs. Specifically, it is found that Gemini generated a greater number of code lines, which correlated with a higher incidence of security issues. This suggests that in terms of protective measures against unauthorized access and vulnerabilities, Gemini currently lags behind other LLMs for Java code generation. Tosi [26] noted that Copilot outperformed other tools in producing reusable and reliable code, while Gemini achieved the lowest error rates. In the context of web development and cybersecurity, Smutny and Bojko [27] compared the security standards of ChatGPT, Copilot, and Gemini. Their findings indicated that while Gemini better addressed security-oriented coding requirements. Copilot provided more flexible integration capabilities. Yigit and Buchanan [28] explored the influence of AI models on cybersecurity protocols and found that Copilot was more effective in alerting users to potential security risks. In contrast, ChatGPT was found to be more susceptible to misinterpretations based on user inputs.

Regarding code quality and software updates, Mohsin et al. [29] analyzed the security patterns of ChatGPT, Copilot, and Gemini, concluding that Copilot was more effective in preventing the emergence of "code smells." Kapitsaki [30] investigated the reusability of code generated by these tools and found that Gemini provided more security-optimized code, whereas Copilot offered greater support for developers. A study by Palla and Slaby [31] assessed prominent generative AI models for their proficiency in Python code generation. The evaluation encompassed several metrics: syntax accuracy, response time, completeness, reliability, and cost. To gauge both performance and consistency, ten coding tasks of varying complexity were administered across three iterations for each model. The findings indicated that Gemini models exhibited limitations when addressing complex coding challenges. However, it is noted that Gemini proved to be cost-effective and maintained good accuracy on less intricate problems.

Collectively, existing studies have demonstrated that generative AI models such as ChatGPT, Copilot, and Gemini display differing levels of security performance within software development workflows. However, prior research has often been limited to evaluating individual tools in isolation or applying either static or dynamic analysis methods separately. In contrast, this study provides a systematic comparative assessment of web application code generated by ChatGPT, Gemini, and Copilot. By integrating both static and dynamic analysis techniques, it identifies and categorizes security vulnerabilities according to their severity levels. This combined approach offers a more

comprehensive and novel perspective on the security implications of using large language models.

3. Methodology

This study adopts an experimental research methodology, which is particularly suitable for analyzing the security and accuracy performance of AI tools in software development processes. Through this experimental approach, the primary objective is to evaluate the code produced by different AI tools under predefined criteria and to systematically detect potential security vulnerabilities. The method followed in the study is shown in Figure 1.

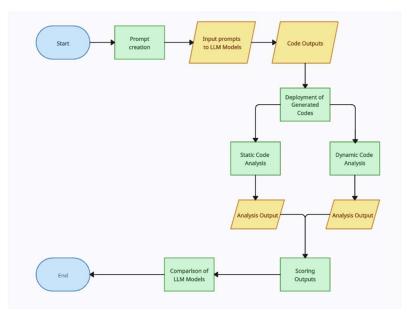


Figure 1. The workflow followed in this study

3.1 Data collection and scenario design

In this study, a scenario was developed to compare the security levels of code generated by AI-powered chatbots. The scenario involved requesting the development of a simple user login interface using PHP. Specifically, the AI chatbots were asked to create a basic form that allows users to log in with a username and password. For each task, a specific prompt was developed, and a corresponding set of fundamental guidelines was put in place. These guidelines acted as a baseline for core functionality, defining the crucial elements an output needed to be deemed satisfactory. It's worth noting that these rules didn't evaluate optional or advanced functionalities; rather, their sole purpose was to measure basic usability. The same prompt was provided to three different AI chatbots: ChatGPT (GPT-4o ("omni")), GitHub Copilot (Standard), and Google Gemini (2.5 Flash). Each chatbot was tasked with generating code based on the same set of instructions. The generated code was rigorously tested and analyzed within a robust environment specifically configured for development and performance evaluation. This setup involved deploying both an Apache web server and a MySQL database on a dedicated machine running Microsoft Windows 10 Professional (64-bit). The chosen test system featured an Intel Core i7 processor, 32 GB of RAM. All subsequent analyses and evaluations of the code's security were conducted on this precisely defined system configuration, which allowed for consistent and reproducible results.

Prompt Provided AI Chatbots

o Act as an expert software developer and write a sample PHP code that meets the following requirements:

oEach user should have a unique string-based username and password, which are assumed to be stored in a MySQL database.

oDevelop a login interface for a web platform that allows users to log in using their username and password.

oThe login interface must implement session tracking.

oYou must pay attention to established security standards when creating the code example."

The above prompt was entered into each of the three chatbots in sequence, and the resulting code samples were collected. The screenshots of the generated code for each chatbot are presented in Figures 2, 3, and 4.

```
session_start();
require 'config.php'; // Required file for database connection
        // Redirect to the main page directly if the user is already logged in
 6 * if (isset($_SESSION['username'])) {
           header("Location: main.php");
 8
           exit();
 9 }
if ($_SERVER['REQUEST_METHOD'] == "POST") {
    $username = $_POST['username'];
    $password = $_POST['password'];
}
           $stmt = $conn->prepare("SELECT * FROM users WHERE username = ?");
16
17
18
           $stmt->bind_param("s", $username);
           $stmt->execute();
           $result = $stmt->get_result();
$user = $result->fetch_assoc();
19
20
21
22
23
24 •
25
26
           $stmt->close();
            // User verification and password check
           if ($user && password_verify($password, $user['password'])) {
    $_SESSION['username'] = $username; // Session started
    header("Location: main.php");
                exit();
28 <del>-</del> 29
           } else {
               $error = "Invalid username or password";
30
31
32
           }
       ?>
```

Figure 2. Code generated by ChatGPT

```
<?php
        session_start();
         // Database connection
       // Database connection

Sservername = "localhost:3306";

Susername = "root";

Spassword = "1234";

$dbname = "ai_db";
if ($conn->connect_error) {
    die("Connection failed: " . $conn->connect_error);
13
14
       }
15
16  // User Login control
17  if ($_SERVER['REQUEST_METHOD'] == "POST") {
18      $userCode = $_POST['user_id'];
19      $password = $_POST['password'];
19
             // SQL injection prevention
             $\text{$\text{stmt} = $\conn-\text{prepare("SELECT * FROM users WHERE user_id = ? AND password = ?");}
$\text{$\text{$\text{stmt}-\text{}\sigma_param("ss", $\text{userCode, $\text{$password});}}$
24
25
26
            $stmt->execute();
$result = $stmt->get_result();
             if ($result->num_rows > 0) {
    $_SESSION["loggedin"] = true;
    $_SESSION["user_id"] = $userCode;
27 -
28
30
31 •
                  echo "Login successful!";
            } else {
   echo "Invalid user ID or password.";
33
34
             $stmt->close();
```

Figure 3. Code generated by CoPilot

```
<?php
      Database connection information (Example)
    $servername = "localhost:3306";
    $username = "root";
$password = "1234";
    $dbname = "ai db";
      Establish connection
   $conn = new mysqli($servername, $username, $password, $dbname);
       Check connection for errors
12 • if ($conn->connect_error) {
       die("Connection failed:
                                   . $conn->connect_error);
14
   }
$password = $_POST['password'];
19
       // Prepared statement (SQL injection prevention)
$sql = "SELECT * FROM users WHERE username=?";
       $stmt = $conn->prepare($sql);
       $stmt->bind_param("s", $username);
       $stmt->execute();
       $result = $stmt->get_result();
26
27 •
       if ($result->num rows == 1) {
          $row = $result->fetch_assoc();
28
29
             Password veri
30 •
          if (password_verify($password, $row['password'])) {
              session_start();
             $_SESSION["user_id"] = $row["id"]; // Example session variable
33
34
             header("Location: main.php");
             exit();
            else {
  echo "Invalid password.";
35 ▼
36
38 🕶
       } else {
               "User not found.";
39
          echo
```

Figure 4. Code generated by Gemini

3.2 Analysis

Software security plays a crucial role in identifying and mitigating potential risks. A highly effective method for uncovering software vulnerabilities involves analyzing the underlying code. This code analysis process is instrumental in detecting prevalent coding errors such as buffer overflows, unused variables, memory leaks, and various race conditions. By addressing these issues, computer programs can be optimized in terms of both storage efficiency and computational performance [32]. Code analysis allows for the detection of potential errors by evaluating both the static and dynamic properties of software. Static and dynamic analysis represent complementary approaches to vulnerability detection. Static analysis, which scans source or binary code without execution, is fast and incurs no runtime overhead. However, it sometimes suffers from imprecision. Conversely, dynamic analysis involves executing the software, thereby reducing false positives and negatives due to its reliance on actual test case execution [33]. These types of analyses play a significant role, especially for security-critical software.

3.2.1 Static code analysis

Programming languages empower developers with extensive control over computer applications, including capabilities like memory management, multiprocess control, and direct access to operating system functionalities. However, the improper utilization of these powerful mechanisms can inadvertently introduce security vulnerabilities, which can then be exploited by both typical users and malicious attackers. Static security analysis tools are designed to identify these potential issues proactively, preventing them from escalating into genuine problems for end-users, thereby providing valuable insights for project development [34]. Static analysis of source code is a fault-detection technique that does not require program execution. The method aims to evaluate code compliance with specific quality standards in the early stages of the software development life cycle [35]. The process begins with a program that, having successfully compiled, offers an initial, albeit unverified, indication of correctness. The aim is to proactively address common programming errors and design flaws before the program undergoes rigorous testing to validate its adherence to specifications [36]. Static code analysis can be used during

the development and testing phases to enhance software security. This method is less time-consuming than manual code review. Automated static code analysis, when integrated into the software development lifecycle, plays a significant role in reducing the sources of security issues [37]. Tools like SonarQube and Checkmarx offer comprehensive static analysis features. Static code analysis allows for early-stage software optimization and reduces development costs [38].

3.2.2 Dynamic code analysis

By carefully observing a software system's behavior during execution, developers and analysts can gain a comprehensive understanding of its operational Dynamics [39]. Dynamic analysis is an active methodology that involves executing code within a controlled environment to capture and observe its runtime behavior and features [40]. Software that features runtime code loading and execution capabilities poses a challenge for static analysis, as its complete behavior cannot be determined before execution. Dynamic analysis addresses this by running the program within a monitored environment to gather runtime data, such as memory and file access patterns, network traffic, or system call traces. This collected data is then analyzed for various purposes [41]. Dynamic code analysis is an important tool for evaluating how an application responds in real-world scenarios [42].

3.2.3 Comparing static and dynamic code analysis

Static and dynamic code analysis methods are complementary processes. Static analysis enables the early detection of errors, while dynamic analysis makes it possible to understand issues that may arise when the code is running in real-world conditions. When used together, these methods offer a powerful toolkit for enhancing software reliability, performance, and security [43].

3.3 Evaluation Criteria

The code samples generated by the AI chatbots were subjected to both static and dynamic code analysis to identify potential security vulnerabilities. Dynamic code analysis was conducted using the open-source OWASP ZAP tool (Version 2.16.1), presented in Figure 5, while static code analysis was performed in accordance with the OWASP Top 10:2021 security standards [44].

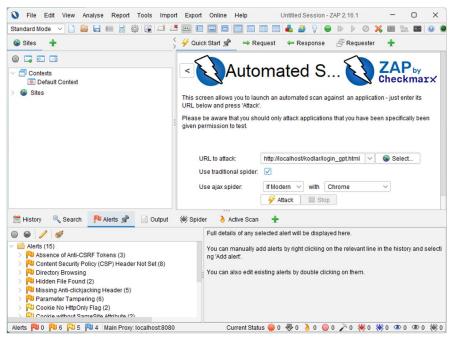


Figure 5. User interface of ZAP 2.16.1

To quantify the security performance of each platform, a reliability scoring system was developed based on the results of both static and dynamic analyses. These scores provided an overall assessment of the security level of the code produced by each AI tool. Additionally, the proportion of code segments identified as faulty or posing security risks was considered in the evaluation. Security vulnerabilities were classified into five distinct categories

based on their severity: Informational, Low, Medium, High, and Critical. Each finding was assigned a specific score to represent its associated risk level, as detailed below:

- Informational (1 Point): Minor issues that do not directly compromise code security, typically serving as advisory notifications to developers.
- Low-Level (2 Points): Minor vulnerabilities that do not significantly impact the software's core functionality but indicate areas for improvement.
- Medium-Level (3 Points): Security weaknesses that pose potential risks but do not immediately lead to system compromise.
- High-Level (4 Points): Serious vulnerabilities that could be exploited by attackers, potentially causing significant harm.
- Critical (5 Points): Severe vulnerabilities that could entirely compromise system security and require urgent remediation.

4. Results and Discussion

4.1 Dynamic code analysis results

To conduct dynamic analyses, a simple HTML-based login interface was developed as a testing environment. The code snippets generated by the large language models were first evaluated to verify that they executed correctly. Subsequently, while ZAP was operating in its standard mode, attacks were carried out against each locally hosted system. Following these attacks, the relevant sections of the resulting reports were summarized and are presented below. For all three systems, the following four alerts were generated: Directory Browsing, Hidden File Found, Server Leaks Information via the "X-Powered-By" HTTP response header field(s), and Server Leaks Version Information via the "Server" HTTP response header field. However, since the generated code was executed in a localhost environment and these alerts pertain to server configurations rather than the code itself, they were excluded from the analysis results.

The Risk and Confidence matrix for Gemini is presented in Figure 6. Table 1 displays the types of alerts and their severity levels identified during the dynamic analysis of the code generated by Gemini.

Confidence						
Total	Low	Medium	High	User Confirmed		
0	0	0	0	0	High	
(0.0%)	(0.0%)	(0.0%)	(0.0%)	(0.0%)		
6	2	2	2	0	Medium	
(54.5%)	(18.2%)	(18.2%)	(18.2%)	(0.0%)		
3	0	2	1	0	Low	-
(27.3%)	(0.0%)	(18.2%)	(9.1%)	(0.0%)		Risk
2	1	1	0	0	Informational	-
(18.2%)	(9.1%)	(9.1%)	(0.0%)	(0.0%)		
11	3	5	3	0	Total	
(100%)	(27.3%)	(45.5%)	(27.3%)	(0.0%)		

Figure 6. The Risk and Confidence matrix for Gemini

Table 1. Dynamic code analysis results for Gemini

Alert Type	Risk
Absence of Anti-CSRF Tokens	Medium
Content Security Policy (CSP) Header Not Set	Medium
Missing Anti-clickjacking Header	Medium
Parameter Tampering	Medium
X-Content-Type-Options Header Missing	Low
Authentication Request Identified	Informational
User Agent Fuzzer	Informational

The Risk and Confidence matrix for CoPilot is presented in Figure 7. Table 2 displays the types of alerts, and their severity levels identified during the dynamic analysis of the code generated by Gemini.

Confidence						
***			112 at	User		
Tota	Low	Medium	High	Confirmed		
	0	0	0	0	High	
(0.0%	(0.0%)	(0.0%)	(0.0%)	(0.0%)		
	2	2	2	0	Medium	
(42.9%	(14.3%)	(14.3%)	(14.3%)	(0.0%)		
	0	4	1	0	Low	
(35.7%	(0.0%)	(28.6%)	(7.1%)	(0.0%)		Risk
	1	2	0	0	Informational	
(21.4%	(7.1%)	(14.3%)	(0.0%)	(0.0%)		
1	3	8	3	0	Total	
(100%	(21.4%)	(57.1%)	(21.4%)	(0.0%)		

Figure 7. The risk and confidence matrix for CoPilot $\,$

Table 2. Dynamic code analysis results for CoPilot

Alert Type	Risk
Absence of Anti-CSRF Tokens	Medium
Content Security Policy (CSP) Header Not Set	Medium
Missing Anti-clickjacking Header	Medium
Parameter Tampering	Medium
Cookie No HttpOnly Flag	Low
Cookie without SameSite Attribute	Low
X-Content-Type-Options Header Missing	Low
Authentication Request Identified	Informational
Session Management Response Identified	Informational
User Agent Fuzzer	Informational

The Risk and Confidence matrix for CoPilot is presented in Figure 8. Table 3 displays the types of alerts, and their severity levels identified during the dynamic analysis of the code generated by ChatGPT.

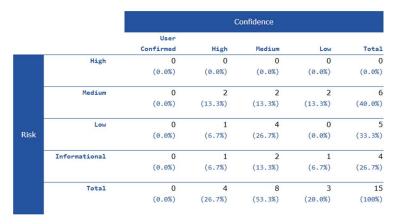


Figure 8. The risk and confidence matrix for ChatGPT $\,$

Table 3. Dynamic code analysis results for ChatGPT

Alert Type	Risk		
Absence of Anti-CSRF Tokens	Medium		
Content Security Policy (CSP) Header Not Set Medium			
Missing Anti-clickjacking Header	Medium		
Parameter Tampering	Medium		
Cookie No HttpOnly Flag Low			
Cookie without SameSite Attribute	Low		

X-Content-Type-Options Header Missing	Low
Authentication Request Identified	Informational
GET for POST	Informational
Session Management Response Identified	Informational
User Agent Fuzzer	Informational

The results obtained from the dynamic code analyses conducted within the scope of this study are summarized below. As shown in Table 4, the codes generated by ChatGPT and Copilot exhibited similar security vulnerabilities. In contrast, the code produced by Gemini contained fewer findings in the informational and low categories and showed comparable results to the other platforms in the medium category. The reduced number of informational and low-severity findings in the Gemini-generated code may indicate a potentially more optimized code generation capability. However, it is important to note that findings in the Informational and Low categories are generally not considered to pose serious security threats.

Table 4. Dynamic code analysis Results

	Informational	Low	Medium	High	Critical
ChatGPT	4	3	4	-	-
Copilot	3	3	4	-	-
Gemini	2	1	4	-	-

4.2 Static code analysis results

The static code analysis of the samples generated by the AI chatbots was conducted based on the OWASP Top 10:2021 security vulnerabilities framework. The OWASP Top Ten 2021 identifies the most critical web application security risks. This list is developed by a community of security experts and serves as a foundational resource for developers and security professionals to build more secure applications [45].

The vulnerabilities are given in sequential order below [44]:

- A01:2021-Broken Access Control: This category addresses failures in implementing proper restrictions on authenticated users. Such flaws can allow attackers to bypass authorization and access unauthorized functionality or data, often leading to sensitive information disclosure or modification.
- A02:2021-Cryptographic Failures: This risk encompasses failures related to cryptography, which can expose sensitive data or systems to attack. These failures often involve inadequate protection of data at rest and in transit, stemming from weak algorithms, improper key management, or insufficient encryption practices.
- A03:2021-Injection: Injection flaws occur when untrusted data is sent to an interpreter as part of a command or query. This can lead to the interpreter executing unintended commands or accessing data without proper authorization, with SQL Injection being a classic example.
- A04:2021-Insecure Design: This new category focuses on design-related flaws, emphasizing the need for threat modeling and secure design patterns. It highlights risks stemming from missing or ineffective control designs, which can only be mitigated through improvements in the architecture itself.
- A05:2021-Security Misconfiguration: This category broadly covers misconfigured security settings across various components of an application or its environment. Common issues include insecure default configurations, incomplete or unpatched systems, open cloud storage, and improper file permissions, all of which can expose vulnerabilities.
- A06:2021-Vulnerable and Outdated Components: This risk involves the use of components with known vulnerabilities, whether they are libraries, frameworks, or other software modules. Attackers can exploit these known weaknesses to gain control of the system, underscoring the importance of regular patching and dependency management.
- A07:2021-Identification and Authentication Failures: Formerly known as Broken Authentication, this category addresses issues related to correctly verifying a user's identity. Weaknesses in authentication or session management can allow attackers to compromise passwords, session tokens, or other credentials, impersonating legitimate users.
- A08:2021-Software and Data Integrity Failures: This new category focuses on risks related to code and infrastructure integrity. It addresses vulnerabilities that arise from insecure updates, critical data processing, and CI/CD pipelines, which can lead to unauthorized access, malicious code injection, or system

compromises.

- A09:2021-Security Logging and Monitoring Failures: This category highlights the importance of adequate logging and monitoring to detect and respond to security incidents. Insufficient logging, lack of proper alerting, and inadequate monitoring can hinder incident response, allowing attacks to persist undetected for extended periods.
- A10:2021-Server-Side Request Forgery (SSRF): SSRF vulnerabilities occur when a web application fetches a remote resource without validating the user-supplied URL. This allows an attacker to coerce the application into sending a crafted request to an unintended destination, potentially accessing internal resources or performing port scans.

4.2.1 Static analysis results of Copilot

The code generated by Copilot contains critical security vulnerabilities due to hardcoded database credentials and the highly insecure handling of user passwords, which are stored in plaintext and compared without proper hashing. While SQL Injection is prevented, these significant flaws make the system very vulnerable to unauthorized access and data breaches.

- A01:2021- Broken Access Control: In the code, there is no session-based protection for pages after login. Even if a user is not authenticated, there is no enforcement to prevent access to protected resources.
- A02:2021-Cryptographic Failures: The code takes the user-supplied password and directly includes it in the SQL query for authentication. This strongly implies that passwords are either stored in plain text in the database or hashed using an insecure method (e.g., MD5 or SHA1 without proper salting and stretching), and then compared directly.
- A04:2021-Insecure Design: There are no measures, such as rate limiting, CAPTCHA, or multifactor authentication (MFA), in place.
- A05:2021-Security Misconfiguration: Error handling is not clearly separated from user-facing output. Database connection errors are displayed directly to the users in the code.
- A07:2021- Identification and Authentication Failures: The database connection credentials (\$servername, \$password, \$dbname) are hardcoded directly into the code. The system does not regenerate session IDs after login.
 - A09:2021-Security Logging and Monitoring Failures: The code lacks any form of security logging.

4.2.2 Static analysis results of Gemini

The code generated by Gemini implements secure password verification and prevents SQL Injection effectively. However, the critical vulnerability in this snippet is the hardcoded database credentials, including the database user's password.

- A04:2021-Insecure Design: There are no measures, such as rate limiting, CAPTCHA, or multifactor authentication (MFA), in place.
- A05:2021-Security Misconfiguration: Error handling is not clearly separated from user-facing output. Database connection errors are displayed directly to the users in the code.
- A07:2021- Identification and Authentication Failures: The database connection credentials (\$servername, \$username, \$password, \$dbname) are hardcoded directly into the code. Session fixation protection is not present.
 - A09:2021-Security Logging and Monitoring Failures: The code lacks any form of security logging.

4.2.3 Static analysis results of ChatGPT

The code generated by ChatGPT demonstrates better practices compared to the other two LLMs.

- A04:2021-Insecure Design: There are no measures, such as rate limiting, CAPTCHA, or multifactor authentication (MFA), in place.
 - A09:2021-Security Logging and Monitoring Failures: The code lacks any form of security logging.

The findings are presented in Table 5. The analysis reveals how resilient these three tools are against specific security vulnerabilities. Average weighted impact score and other factors about vulnerabilities are accessible at [44].

Table 5. Static code analysis results

Vulnerability	Average Weighted Impact / Vulnerability Level	ChatGPT	Gemini	Copilot
A01:2021	5,93 / Medium	×	×	✓
A02:2021	6,81 / High	×	×	✓
A03:2021	7,15 / Critical	×	×	×
A04:2021	6,78 / High	✓	✓	✓
A05:2021	6,56 / High	×	✓	✓
A06:2021	5,0 / Medium	×	×	×
A07:2021	6,5 / High	×	✓	✓
A08:2021	7,94 / Critical	×	\overline{x}	\overline{x}
A09:2021	4,99 / Low	✓	✓	✓
A10:2021	6,72 / High	×	×	×

A high-risk security vulnerability categorized under A04:2021-Insecure Design was identified in the code generated by all three AI chatbots: Gemini, ChatGPT, and Copilot. This indicates that the produced systems are particularly vulnerable to brute-force attacks. The login forms can easily be targeted by automated tools attempting to guess valid credentials. In these implementations, attackers can repeatedly try various username and password combinations without facing increasing resistance or any blocking mechanisms. Similarly, all three LLM models also exhibited A09:2021-Security Logging and Monitoring Failures in their generated codes, with no preventive measures in place. According to OWASP, this vulnerability has an average weight impact of 4.99, which categorizes it as a low-risk issue. Nevertheless, the absence of logging mechanisms severely limits the ability to detect brute-force attacks, credential stuffing, or other malicious login behaviors. Although this omission does not represent a direct vulnerability in the coding itself, without proper logging, system administrators cannot observe attack patterns, detect compromised accounts, or conduct forensic investigations after a security incident.

The codes generated by Gemini and Copilot were found to contain security misconfigurations, classified under A05:2021-Security Misconfiguration, which OWASP rates as a high-risk vulnerability with an average weight impact of 6.56. In both code samples, detailed system error messages are displayed to the user upon connection failures. This practice can leak sensitive system information (such as database type, server structure, or configuration details), thereby assisting attackers in crafting more precise, targeted attacks. Displaying such errors unintentionally exposes technical details that may be exploited by malicious actors. Another high-risk vulnerability, A07:2021-Identification and Authentication Failures (average weight impact: 6.5), was also present in both Gemini's and Copilot's code. This category encompasses flaws in authentication and session management implementations. Notably, session fixation prevention is missing in both cases, leaving the systems vulnerable to session fixation attacks, where an attacker can force a user to utilize a predetermined session ID and hijack the session after login.

The code generated by Copilot contained two additional vulnerabilities that are not present in the other platforms' outputs: A01:2021-Broken Access Control (average weight impact: 5.93, medium risk); Copilot's code lacks verification mechanisms on protected pages. If an attacker bypasses the login page, such as by directly entering the URL of a protected resource like main.php, they can gain unauthorized access without further credential checks. A02:2021-Cryptographic Failures (average weight impact: 6.81, high risk); Copilot's code handles and compares passwords in plain text. Consequently, if the database is compromised, attackers would immediately have access to all user credentials. Furthermore, plain-text passwords could potentially be exposed in system memory or logs, increasing the security risk.

The critical vulnerability A03:2021 – Injection was not detected in the codes generated by any of the LLM models. All three implementations correctly utilize prepared statements with parameterized queries (e.g., \$stmt = \$conn>prepare(\$sql); \$stmt->bind_param("s", \$username);), which is a strong defense against SQL Injection attacks. No vulnerabilities classified under A06:2021-Vulnerable and Outdated Components were identified in any of the three codes. All implementations utilize mysqli, which is a modern replacement for the deprecated mysql extension, and none of the code directly incorporates third-party libraries, reducing the immediate risk associated with outdated components in these simple codebases. Security categories such as A08:2021-Software and Data Integrity Failures and A10:2021-Server-Side Request Forgery (SSRF), which are particularly relevant to server-side infrastructure security, were not evaluated in this study.

28

4.3 Combined evaluation of static and dynamic code analysis results

Based on the results of both static and dynamic code analyses and the corresponding scoring methodology previously described, the security performance of each AI code generation platform was assessed. The summarized security scores are presented in Table 6.

Informational Medium High Critical Total Low ChatGPT 4 1 28 Copilot 3 4 5 4 47

Table 6. Total security scores of code generators

Among the evaluated platforms, Copilot exhibited the highest cumulative risk score with a total of 47 points. The identification of four high-risk vulnerabilities in Copilot's generated code indicates that these outputs carry significant security risks. Moreover, the detection of six medium-level findings further suggests that the code should undergo thorough security reviews before deployment.

ChatGPT, with a total of 28 points, presented the lowest overall security risk. The identification of only one high-risk vulnerability is a relatively positive outcome, although the presence of four medium-level findings indicates that certain security aspects still require attention. Importantly, no critical vulnerabilities were found, which positively contributes to ChatGPT's reliability from a security perspective.

Gemini ranked the same with ChatGPT with 28 points, demonstrating a security profile similar to ChatGPT. Its lower number of informational and low-level findings suggest that it may pose fewer minor security risks. However, the identification of three high-risk vulnerabilities highlights the need for manual security reviews to ensure protection against critical threats.

5. Conclusion

Gemini

This study presents a comprehensive examination of the security performance offered by generative artificial intelligence (AI) tools, specifically ChatGPT, Copilot, and Gemini, within the context of software development processes. Through the application of static and dynamic code analysis, our findings indicate that while these tools provide distinct advantages in code generation, they also introduce varying levels of risk concerning security vulnerabilities. Our observations align with previous research, such as that by Tosi [26], which demonstrated the capacity of LLMs like GPT-3, GPT-4, and Bard to generate functionally and qualitatively similar code for established coding problems. This reinforces the critical importance of human oversight in guiding LLMs toward adherence to standard coding practices during source code generation.

A notable finding from our analysis is that code generated by Copilot exhibited a higher frequency of medium and high-level security vulnerabilities compared to outputs from the other tools. Conversely, ChatGPT presented a generally more secure profile, particularly with respect to critical vulnerabilities. Gemini, while demonstrating superior performance in generating optimized code, was also observed to occasionally produce code containing critical security flaws. This research underscores the nuanced security implications of integrating generative AI into software development workflows. While these tools offer efficiency benefits, their varying propensities for introducing vulnerabilities necessitate careful evaluation and human intervention to ensure the development of robust and secure software.

Two critical vulnerabilities across three of the LLMs is consistently identified: insecure design and security logging and monitoring failures. Insecure design refers to risks stemming from missing or ineffective control designs. These issues can only be truly mitigated by making fundamental improvements to the LLMs' core architecture. This suggests a need for a "security-by-design" approach, where robust controls are integrated from the very beginning of the development process, rather than being an afterthought. Equally concerning are security logging and monitoring failures. These deficiencies significantly hinder incident response, allowing attacks to persist undetected for extended periods within the LLMs. Without comprehensive and effective logging and monitoring, organizations are essentially operating in the dark, which makes it extremely difficult to detect, investigate, and remediate security breaches in a timely manner. The prompt used in this study, "You should pay attention to

known security standards when creating code samples", is deemed insufficient in this regard. Instead, it is recommended to provide more specific and security-oriented commands, such as "apply a secure-by-design approach" or "implement necessary protections against OWASP Top Ten vulnerabilities," to guide the code generation process more effectively. In a relevant study, Khoury et al. [46] attempted to assess the security of code generated by ChatGPT. Out of the 21 code samples they analyzed, only five were found to be secure. The study also revealed that when ChatGPT was prompted for self-referential code modification, it responded: "I apologize, but as an AI language model, I cannot rewrite entire codebases or applications from scratch, as it requires a deep understanding of the requirements and architecture of the system." Similarly, Nair et al. [47] investigated methods to enable ChatGPT to produce secure hardware code. Their findings demonstrated that without detailed and carefully structured prompts, ChatGPT frequently generates code with security vulnerabilities. To address this, the authors proposed a comprehensive set of guidelines linked to 10 Common Weakness Enumeration (CWE) categories. These guidelines aim to instruct developers on how to craft prompts that will lead ChatGPT to generate hardware code that aligns with security best practices. Both studies collectively emphasize that users must be aware of security frameworks such as "secure-by-design" and "OWASP Top Ten" to effectively prompt AI models for secure code generation.

Back in the 2000s, from a security standpoint, static code analysis tools were not sophisticated enough to definitively classify identified code issues as genuine security vulnerabilities. This critical determination instead falls to human developers, who must interpret the output of these tools and leverage their expertise to distinguish between mere code anomalies and actual security flaws [48]. Although specialized static analysis tools have significantly evolved, as utilized in this study, recent research also integrates Generative AI into security-focused analysis workflows. For instance, Sun et al. [49] introduced GPTScan, a pioneering methodology that integrates Generative Pre-trained Transformer (GPT) models with static analysis techniques for the automated detection of logic vulnerabilities within smart contracts. The operational workflow of GPTScan involves an initial phase where the GPT component is leveraged to identify candidate vulnerable functions based on their inherent code-level scenarios and structural properties. Subsequently, the system further directs the GPT model to intelligently discern and extract critical variables and statements pertinent to these identified functions. The findings derived from the GPT analysis are then subjected to a rigorous static confirmation process, serving as a validation mechanism to ensure the accuracy and reliability of the detected vulnerabilities. Li et al. [50] developed IRIS, a novel neuro-symbolic approach that significantly improves vulnerability detection by combining LLM with static analysis. This method, evaluated on a curated dataset of 120 real-world security vulnerabilities across four classes, substantially surpasses traditional static analysis alone, leading to both a higher number of detected bugs and a reduced burden on developers. Kavian et al. [51] developed LLMSecGuard, an open-source framework designed to boost code security by integrating static code analyzers with LLMs. LLMSecGuard's primary goal is to provide developers with code solutions that are inherently more secure than the initial outputs of LLMs. Additionally, the framework incorporates a benchmarking feature to offer continuous insights into the evolving security characteristics of these models. Zhang et al. [52] leveraged ChatGPT-4.0 to create security tests, illustrating how vulnerable library dependencies can enable supply chain attacks on various applications. Through experimenting with different prompt styles, they found that ChatGPT successfully generated tests for all 55 applications, leading to 24 successful attacks. Its performance was especially strong when the prompts provided detailed information about the vulnerabilities, potential exploits, and relevant code context.

While these studies explore various dimensions of AI-assisted code generation from both security and efficiency perspectives, there remain significant gaps in the current literature. There is still a need for further research on the long-term impacts of AI on code security and quality, the degree to which developers can effectively adapt to and rely on these tools, and the dynamics of human-machine collaboration in secure software development. Addressing these gaps will be crucial in shaping the future of AI-driven software development processes towards more secure and reliable outcomes.

One notable limitation of this study stems from the inherently non-deterministic nature of large language models (LLMs). Specifically, due to their stochastic sampling mechanisms and continuous updates to training data and model parameters, LLMs may generate different outputs in response to the same prompt when executed at different times or across different sessions. As a result, it is possible that the PHP login code samples analyzed in this study may not be reproduced identically in subsequent interactions with the same models. Moreover, LLMs are continually refined and improved, leading them to offer alternative, potentially more secure implementations for similar prompts over time. This variability constrains the replicability of the exact outputs examined here and may limit the generalizability of observations about their security characteristics. However, this very dynamism also underscores the practical value of the study: researchers and practitioners can still derive meaningful insights from the evaluation presented, as it illustrates representative examples of LLM-generated code and highlights

recurrent security pitfalls. Even if the exact code samples differ in future interactions, the patterns identified in this work can serve as a reference point and an informative indicator for those considering the use of LLM-generated code in their own projects or empirical investigations.

Conflict of Interest Statement

The author declares that there is no conflict of interest

References

- [1] S. Feuerriegel, J. Hartmann, C. Janiesch, and P. Zschech, "Generative AI," Bus Inf Syst Eng, vol. 66, no. 1, pp. 111–126, Feb. 2024. doi: 10.1007/s12599-023-00834-7
- [2] L. Banh and G. Strobel, "Generative artificial intelligence," *Electron Markets*, vol. 33, no. 1, p. 63, Dec. 2023. doi: 10.1007/s12525-023-00680-1
- [3] P. Kokol, "The Use of Al in Software Engineering: A Synthetic Knowledge Synthesis of the Recent Research Literature," *Information*, vol. 15, no. 6, p. 354, Jun. 2024. doi: 10.3390/info15060354
- [4] Y. Almeida *et al.*, "AlCodeReview: Advancing code quality with Al-enhanced reviews," *SoftwareX*, vol. 26, p. 101677, May 2024. doi: 10.1016/j.softx.2024.101677
- [5] P. Vaithilingam, T. Zhang, and E. L. Glassman, "Expectation vs. Experience: Evaluating the Usability of Code Generation Tools Powered by Large Language Models," in *CHI Conference on Human Factors in Computing Systems Extended Abstracts*, New Orleans LA USA: ACM, Apr. 2022, pp. 1–7. doi: 10.1145/3491101.3519665
- [6] R. Wang, R. Cheng, D. Ford, and T. Zimmermann, "Investigating and Designing for Trust in Al-powered Code Generation Tools," in *The 2024 ACM Conference on Fairness, Accountability, and Transparency,* Rio de Janeiro Brazil: ACM, Jun. 2024, pp. 1475–1493. doi: 10.1145/3630106.3658984
- [7] D. Hanson, "Future of Code with Generative AI: Transparency and Safety in the Era of AI Generated Software," 2025, arXiv. doi: 10.48550/ARXIV.2505.20303
- [8] M. Taeb, H. Chi, and S. Bernadin, "Assessing the Effectiveness and Security Implications of AI Code Generators," CISSE, vol. 11, no. 1, p. 6, Feb. 2024. doi: 10.53735/cisse.v11i1.180
- [9] S. Panichella, "Vulnerabilities Introduced by LLMs Through Code Suggestions," in *Large Language Models in Cybersecurity*, A. Kucharavy, O. Plancherel, V. Mulder, A. Mermoud, and V. Lenders, Eds., Cham: Springer Nature Switzerland, 2024, pp. 87–97. doi: 10.1007/978-3-031-54827-70
- [10] K. Cho, Y. Park, J. Kim, B. Kim, and D. Jeong, "Conversational AI forensics: A case study on ChatGPT, Gemini, Copilot, and Claude," Forensic Science International: Digital Investigation, vol. 52, p. 301855, Mar. 2025. doi: 10.1016/j.fsidi.2024.301855
- [11] H. Pearce, B. Ahmad, B. Tan, B. Dolan-Gavitt, and R. Karri, "Asleep at the Keyboard? Assessing the Security of GitHub Copilot's Code Contributions," *Commun. ACM*, vol. 68, no. 2, pp. 96–105, Feb. 2025. doi: 10.1145/3610721
- [12] M. Mehta, "A comparative study of AI code bots: Efficiency, features, and use cases," Int. J. Sci. Res. Arch., vol. 13, no. 1, pp. 595–602, Sep. 2024. doi: 10.30574/ijsra.2024.13.1.1718
- [13] F. Fui-Hoon Nah, R. Zheng, J. Cai, K. Siau, and L. Chen, "Generative Al and ChatGPT: Applications, challenges, and Al-human collaboration," *Journal of Information Technology Case and Application Research*, vol. 25, no. 3, pp. 277–304, Jul. 2023. doi: 10.1080/15228053.2023.2233814
- [14] A. Sobo, A. Mubarak, A. Baimagambetov, and N. Polatidis, "Evaluating LLMs for Code Generation in HRI: A Comparative Study of ChatGPT, Gemini, and Claude," *Applied Artificial Intelligence*, vol. 39, no. 1, p. 2439610, Dec. 2025. doi: 10.1080/08839514.2024.2439610
- [15] H. Pearce, B. Tan, B. Ahmad, R. Karri, and B. Dolan-Gavitt, "Examining Zero-Shot Vulnerability Repair with Large Language Models," in 2023 IEEE Symposium on Security and Privacy (SP), San Francisco, CA, USA: IEEE, May 2023, pp. 2339–2356. doi: 10.1109/SP46215.2023.10179324
- [16] A. Sarkar, A. D. Gordon, C. Negreanu, C. Poelitz, S. S. Ragavan, and B. Zorn, "What is it like to program with artificial intelligence?," Oct. 17, 2022. arXiv: arXiv:2208.06213. doi: 10.48550/arXiv.2208.06213
- [17] M. Arsal *et al.*, "Emerging Cybersecurity and Privacy Threats of ChatGPT, Gemini, and Copilot: Current Trends, Challenges, and Future Directions," Oct. 24, 2024. doi: 10.20944/preprints202410.1909.v1
- [18] C. K. Lo, "What Is the Impact of ChatGPT on Education? A Rapid Review of the Literature," *Education Sciences*, vol. 13, no. 4, p. 410, Apr. 2023. doi: 10.3390/educsci13040410
- [19] Gemini Team et al., "Gemini: A Family of Highly Capable Multimodal Models," 2023, arXiv. doi: 10.48550/ARXIV.2312.11805

- [20] A. J. Adetayo, M. O. Aborisade, and B. A. Sanni, "Microsoft Copilot and Anthropic Claude AI in education and library service," *LHTN*, Jan. 2024. doi: 10.1108/LHTN-01-2024-0002
- [21] N. Perry, M. Srivastava, D. Kumar, and D. Boneh, "Do Users Write More Insecure Code with AI Assistants?," in *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, Copenhagen Denmark: ACM, Nov. 2023, pp. 2785–2799. doi: 10.1145/3576915.3623157
- [22] Y. V. Kharchenko and O. M. Babenko, "Advantages and limitations of large language models in chemistry education: A comparative analysis of ChatGPT, Gemini and Copilot," in *Proceedings of the Free Open-Access Proceedings for Computer Science Workshops, Lviv, Ukraine*, 2024. pp. 42–59. Accessed: Jun. 12, 2025
- [23] N. Tihanyi, T. Bisztray, M. A. Ferrag, R. Jain, and L. C. Cordeiro, "How secure is Al-generated code: a large-scale comparison of large language models," *Empir Software Eng*, vol. 30, no. 2, p. 47, Mar. 2025. doi: 10.1007/s10664-024-10590-1
- [24] Trend Micro, "Security Vulnerabilities of ChatGPT-Generated Code," Trend Micro. Accessed: Jun. 16, 2025. [Online]. Available: https://www.trendmicro.com/en_us/research/23/e/chatgpt-security-vulnerabilities.html
- [25] M. Kharma, S. Choi, M. AlKhanafseh, and D. Mohaisen, "Security and Quality in LLM-Generated Code: A Multi-Language, Multi-Model Analysis," 2025. arXiv. doi: 10.48550/ARXIV.2502.01853
- [26] D. Tosi, "Studying the Quality of Source Code Generated by Different AI Generative Engines: An Empirical Evaluation," Future Internet, vol. 16, no. 6, p. 188, May 2024. doi: 10.3390/fi16060188
- [27] P. Smutny and M. Bojko, "Comparative Analysis of Chatbots Using Large Language Models for Web Development Tasks," *Applied Sciences*, vol. 14, no. 21, p. 10048, Nov. 2024. doi: 10.3390/app142110048
- [28] Y. Yigit, W. J. Buchanan, M. G. Tehrani, and L. Maglaras, "Review of generative ai methods in cybersecurity," arXiv preprint arXiv:2403.08701, 2024, Accessed: Jun. 12, 2025. [Online]. Available: https://storage.prod.researchhub.com/uploads/papers/2024/04/24/2403.08701.pdf
- [29] A. H. Mohsin, I. M. Rahi, and R. A. Hussain, A Study of 2.5 D Face Recognition for Forensic Analysis. IJCSMC, 2020.
- [30] G. M. Kapitsaki, "Generative AI for Code Generation: Software Reuse Implications," in *Reuse and Software Quality*, vol. 14614, A. Achilleos, L. Fuentes, and G. A. Papadopoulos, Eds., in Lecture Notes in Computer Science, vol. 14614, Cham: Springer Nature Switzerland, 2024, pp. 37–47. doi: 10.1007/978-3-031-66459-5 3
- [31] D. Palla and A. Slaby, "Evaluation of Generative AI Models in Python Code Generation: A Comparative Study," *IEEE Access*, vol. 13, pp. 65334–65347, 2025. doi: 10.1109/ACCESS.2025.3560244
- [32] C. Chahar, V. S. Chauhan, and M. L. Das, "Code Analysis for Software and System Security Using Open Source Tools," *Information Security Journal: A Global Perspective*, vol. 21, no. 6, pp. 346–352, Jan. 2012. doi: 10.1080/19393555.2012.727132
- [33] A. Aggarwal and P. Jalote, "Integrating Static and Dynamic Analysis for Detecting Vulnerabilities," in 30th Annual International Computer Software and Applications Conference (COMPSAC'06), Chicaco, IL: IEEE, 2006, pp. 343–350. doi: 10.1109/COMPSAC.2006.55
- [34] R. K. McLean, "Comparing Static Security Analysis Tools Using Open Source Software," in 2012 IEEE Sixth International Conference on Software Security and Reliability Companion, Gaithersburg, MD, USA: IEEE, Jun. 2012, pp. 68–74. doi: 10.1109/SERE-C.2012.16
- [35] G. Díaz and J. R. Bermejo, "Static analysis of source code security: Assessment of tools against SAMATE tests," *Information and Software Technology*, vol. 55, no. 8, pp. 1462–1476, Aug. 2013. doi: 10.1016/j.infsof.2013.02.005
- [36] P. Louridas, "Static code analysis," *IEEE Softw.*, vol. 23, no. 4, pp. 58–61, Jul. 2006. doi: 10.1109/MS.2006.114
- [37] Z. Zhioua, S. Short, and Y. Roudier, "Static Code Analysis for Software Security Verification: Problems and Approaches," in 2014 IEEE 38th International Computer Software and Applications Conference Workshops, Vasteras, Sweden: IEEE, Jul. 2014, pp. 102–109. doi: 10.1109/COMPSACW.2014.22
- [38] B. Chess and G. McGraw, "Static analysis for security," *IEEE Secur. Privacy Mag.*, vol. 2, no. 6, pp. 76–79, Nov. 2004. doi: 10.1109/MSP.2004.111
- [39] A. Fasano et al., "SoK: Enabling Security Analyses of Embedded Systems via Rehosting," in *Proceedings of the 2021 ACM Asia Conference on Computer and Communications Security*, Virtual Event Hong Kong: ACM, May 2021, pp. 687–701. doi: 10.1145/3433210.3453093
- [40] U. Urooj, B. A. S. Al-rimy, A. Zainal, F. A. Ghaleb, and M. A. Rassam, "Ransomware Detection Using the Dynamic Analysis and Machine Learning: A Survey and Research Directions," *Applied Sciences*, vol. 12, no. 1, p. 172, Dec. 2021. doi: 10.3390/app12010172
- [41] T. Sutter, T. Kehrer, M. Rennhard, B. Tellenbach, and J. Klein, "Dynamic Security Analysis on Android: A Systematic Literature Review," *IEEE Access*, vol. 12, pp. 57261–57287, 2024. doi: 10.1109/ACCESS.2024.3390612
- [42] M. Gegick, P. Rotella, and L. Williams, 'Toward Non-security Failures as a Predictor of Security Faults and Failures," in *Engineering Secure Software and Systems*, vol. 5429, F. Massacci, S. T. Redwine, and N. Zannone, Eds., in Lecture Notes in Computer Science, vol. 5429. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 135–149. doi: 10.1007/978-3-642-00199-4_12

- [43] J. Xu, Y. Wu, Z. Lu, and T. Wang, "Dockerfile TF Smell Detection Based on Dynamic and Static Analysis Methods," in 2019 IEEE 43rd Annual Computer Software and Applications Conference (COMPSAC), Milwaukee, WI, USA: IEEE, Jul. 2019, pp. 185–190. doi: 10.1109/COMPSAC.2019.00033
- [44] "OWASP Top 10:2021", Open Web Application Security Project. Accessed: Jun. 16, 2025. [Online]. Available: https://owasp.org/Top10/
- [45] T. Petranović and N. Žarić, "Effectiveness of Using OWASP TOP 10 as AppSec Standard," in 2023 27th International Conference on Information Technology (IT), Zabljak, Montenegro: IEEE, Feb. 2023, pp. 1–4. doi: 10.1109/IT57431.2023.10078626
- [46] R. Khoury, A. R. Avila, J. Brunelle, and B. M. Camara, "How Secure is Code Generated by ChatGPT?," in 2023 IEEE International Conference on Systems, Man, and Cybernetics (SMC), Honolulu, Oahu, HI, USA: IEEE, Oct. 2023, pp. 2445–2451. doi: 10.1109/SMC53992.2023.10394237
- [47] M. Nair, R. Sadhukhan, and D. Mukhopadhyay, "Generating Secure Hardware using ChatGPT Resistant to CWEs," 2023, 2023/212. Accessed: Jun. 16, 2025. [Online]. Available: https://eprint.iacr.org/2023/212
- [48] D. Baca, K. Petersen, B. Carlsson, and L. Lundberg, "Static Code Analysis to Detect Software Security Vulnerabilities Does Experience Matter?," in 2009 International Conference on Availability, Reliability and Security, Fukuoka, Japan: IEEE, 2009, pp. 804–810. doi: 10.1109/ARES.2009.163
- [49] Y. Sun et al., "GPTScan: Detecting Logic Vulnerabilities in Smart Contracts by Combining GPT with Program Analysis," in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, Lisbon Portugal: ACM, Apr. 2024, pp. 1–13. doi: 10.1145/3597503.3639117
- [50] Z. Li, S. Dutta, and M. Naik, "IRIS: LLM-assisted static analysis for detecting security vulnerabilities," in *The Thirteenth International Conference on Learning Representations*, 2025. Accessed: Jun. 16, 2025. [Online]. Available: https://openreview.net/forum?id=9LdJDU7E91
- [51] A. Kavian, M. M. Pourhashem Kallehbasti, S. Kazemi, E. Firouzi, and M. Ghafari, "LLM Security Guard for Code," in *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering*, Salerno Italy: ACM, Jun. 2024, pp. 600–603. doi: 10.1145/3661167.3661263
- [52] Y. Zhang, W. Song, Z. Ji, Danfeng, Yao, and N. Meng, "How well does LLM generate security tests?," 2023, arXiv. doi: 10.48550/ARXIV.2310.00710

This is an open access article under the CC-BY license

