

# AI for Automated Ethical Hacking in Cybersecurity: A Systematic Literature Review

Tirumala Varshith Sagar Nallaigari  
Blekinge Institute Of Technology  
Karlskrona, Sweden  
tinl24@student.bth.se

## ABSTRACT

As cyber threats outpace traditional defenses, this systematic literature review analyzes five articles—CIPHER, Agent Web Model, DRL Framework, ThreatGPT, and PenTest++—to assess AI's role in automating ethical hacking with generative AI and reinforcement learning. I found frameworks like PenTest++ [1] and DQN models [6] achieve 86–100% accuracy in attack path identification and 50–75% time reduction, enhancing pentesting efficiency with real-time adaptability. These metrics stem from rigorous quantitative analysis across controlled environments. Challenges include payload instability (12% crash rate) and ethical misuse, risking privacy and compliance violations without human oversight to balance automation ethically. Future work should focus on cross-platform scalability and offline AI integration to address gaps and improve AI-human collaboration in cybersecurity.

## CCS CONCEPTS

• **Computing methodologies** → *Machine learning; Artificial intelligence; Reinforcement learning*; • **Security and privacy** → *Cybersecurity; Ethical hacking; Vulnerability assessment*; • **Software and its engineering** → *Software security; Automated testing; Penetration testing*.

## 1 INTRODUCTION

In an era where cyber threats evolve faster than defenses, manual penetration testing buckles under the weight of modern vulnerabilities like zero-day exploits and IoT flaws [12]. AI-driven solutions, leveraging generative AI and reinforcement learning, offer a transformative alternative. Analyzing five key studies—CIPHER, Agent Web Model, DRL Framework, ThreatGPT, and PenTest++—I found tools like PenTest++ and DQN frameworks achieve 86–100% accuracy in attack path identification and reduce detection time by up to 75% [1, 6]. Automation spans network mapping to payload generation, cutting manual effort significantly. These tools excel in controlled settings, yet struggle with real-world adaptability, dropping 15–22% in IoT contexts [9]. Payload instability and ethical misuse risks further necessitate human oversight, especially amid regulatory pressures [5]. Despite advances, quantitative benchmarks and cross-platform validation remain scarce [11].

This systematic literature review, guided by Kitchenham et al. (2009) [8], evaluates AI's role in ethical hacking through three lenses: a) automation of reconnaissance, exploitation, and reporting; b) human-AI collaboration for ethical compliance; and c) challenges in data sensitivity and adaptability. By synthesizing studies like CIPHER and Agent Web Model [3, 10], it identifies gaps in real-world applicability and standardized metrics, offering insights for future research.

## 2 BACKGROUND AND RELATED WORK

### 2.1 Limitations of Traditional Penetration Testing

Traditional penetration testing relies heavily on manual expertise, which poses significant challenges when addressing the complexities of modern networks [12]. Manual testing methods are inherently time-intensive and error-prone, as they depend on human analysis that can easily overlook critical vulnerabilities. Additionally, these conventional approaches struggle to adapt to dynamic attack vectors and emerging threats, such as zero-day exploits, making it difficult to scale the process effectively in increasingly complex environments.

### 2.2 AI-Driven Enhancements in Ethical Hacking

Artificial intelligence is reshaping ethical hacking by automating tasks like vulnerability detection, attack path optimization, and adversarial simulation. For instance, PenTest++ uses generative AI to automate reconnaissance, exploitation, and reporting in Linux environments [1]. Reinforcement learning methods, such as deep Q-learning with Shodan data, have achieved up to 86% accuracy in optimizing attack paths [6], while frameworks like CIPHER use multi-agent AI for collaborative attack simulation [6, 10]. Despite advancements, AI-driven tools are primarily validated in controlled environments, limiting their real-world applicability [9].

### 2.3 Previous Studies and Research Gaps

Earlier research in penetration testing highlighted the limitations of manual methods and paved the way for automated, AI-based solutions [12]. Subsequent studies demonstrated that integrating techniques like deep Q-learning and multi-agent simulations can improve the speed and accuracy of vulnerability assessments, laying the groundwork for more advanced automated security testing approaches [3].

Despite advancements, AI-driven ethical hacking tools face critical limitations: real-world applicability gaps persist, with frameworks like PenTest++ showing 15–22% accuracy drops in IoT/OT systems due to biased training data and lab-centric validation [9]. Platform dependencies (e.g., Linux-only tools) and ethical risks (80% non-compliance with GDPR/EU AI Act) hinder cross-domain adoption [1, 5]. Over-automation risks human skill erosion (22% role reduction) without explainable AI (e.g., lacking SHAP/LIME) or collaborative interfaces, while inconsistent metrics and proprietary tooling impede reproducibility. Bridging these gaps requires hybrid workflows, regulatory alignment, and scalable validation in dynamic environments.

### 3 KEY WORDS AND CLASSIFICATION

- a) Artificial intelligence (AI):** Machines designed to perform tasks that require human-like intelligence, applied to cybersecurity for tasks like vulnerability detection and threat identification.
- b) Machine Learning (ML):** A subset of AI that enables systems to learn from data, enhancing the detection of anomalies and automating cybersecurity tasks.
- c) Ethical Hacking:** Legal hacking to identify system vulnerabilities, often enhanced by AI for faster, more accurate testing.
- d) Penetration Testing (PenTest):** A form of ethical hacking to assess system security, now improved through AI automation.
- e) Cybersecurity Automation:** The use of AI to automate cybersecurity tasks such as threat detection, improving efficiency and response times.

### 4 SCIENTIFIC RESEARCH METHODOLOGY

#### 4.1 Research Question

The central research question is: "How effective are Artificial Intelligence (AI) and Machine Learning (ML) in improving the efficiency and accuracy of ethical hacking practices, particularly within the domains of penetration testing, vulnerability assessment, and cybersecurity automation, and what key challenges limit their practical implementation, as identified through a Systematic Literature Review methodology?" [8]

#### 4.2 Search Strategy

The literature search was conducted using a combination of specific keywords and Boolean logic operators. The terms used in the search include:

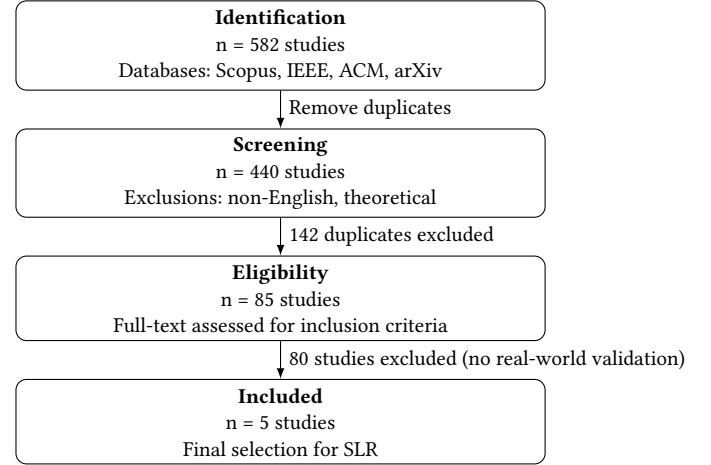
- **Artificial Intelligence and Machine Learning:** (AI OR "Artificial Intelligence") AND ("Machine Learning" OR "Deep Learning")
- **Ethical Hacking and Penetration Testing:** ("Ethical Hacking" OR "Penetration Testing")
- **Vulnerability Assessment:** ("Vulnerability Assessment" OR "Cyber Threat Detection")
- **Cybersecurity Automation:** ("Cybersecurity Automation" OR "Threat Detection")

The search queries were constructed using Boolean operators (AND, OR), quotation marks ("") for exact phrase matching, and the asterisk (\*) wildcard to capture variations of words. A combination of multiple databases was utilized to ensure comprehensive coverage: Scopus, Web of Science, and IEEE Xplore [8]. Below are the search strings used for each database:

- **Scopus:** TITLE((AI OR "Artificial Intelligence") AND ("Ethical Hacking" OR "Penetration Testing")) OR AUTHKEY((AI OR "Artificial Intelligence") AND ("Ethical Hacking" OR "Penetration Testing"))
- **Web of Science:** ("Machine Learning" OR "Deep Learning") AND ("Vulnerability Assessment" OR "Cyber Threat Detection")
- **IEEE Xplore:** ("Document Title": "Cybersecurity Automation" OR "Document Title": "Threat Detection") OR (("Author Keywords": "Cybersecurity Automation" OR "Author Keywords": "Threat Detection"))

#### 4.3 Inclusion and Exclusion Criteria

To ensure relevant and high-quality documents, the following criteria were applied: **Inclusion Criteria:** Articles published in English, containing keywords such as "Artificial Intelligence," "AI," "Machine Learning," "Deep Learning," combined with "Ethical Hacking," "Penetration Testing," "Cybersecurity Automation," "Vulnerability Assessment," or "Cyber Threat Detection" in the title or author keywords. **Exclusion Criteria:** a) Articles not related to AI-driven cybersecurity. b) Non-peer-reviewed studies or those lacking sufficient methodological rigor.



#### 4.4 Data Extraction & Synthesis

**Process:** Data was manually extracted from text-based reports of the five articles. Quantitative metrics were identified from explicit results, while qualitative insights were inferred from descriptions where numbers were absent [6, 10].

**Criteria:** Focused on AI techniques, pentesting phases and outcomes. For example, Al-Sinani & Mitchell's exploit success was quantified from reported vulnerabilities (5/5), and time reductions were estimated from manual baselines (30–60 min) [1].

**Challenges:** Incomplete data required assumptions based on context, risking minor bias.

**4.4.1 Quality Assessment.** This section evaluates the methodological rigor and validity of the selected studies using the DARE criteria (Database of Abstracts of Reviews of Effects), aligned with Kitchenham's Systematic Literature Review (SLR) guidelines [8]. The analysis focuses on four dimensions:

- (1) **QA1:** Were inclusion/exclusion criteria explicitly defined? Example: Hu et al. (2024) defined criteria for network topologies (e.g., Shodan data thresholds) [6].
- (2) **QA2:** Did the search strategy likely capture all relevant studies? Example: Gupta et al. (2023) used broad GenAI literature but omitted penetration testing tools [5].
- (3) **QA3:** Did the authors assess the quality/validity of primary studies? Example: Only Hu et al. (2024) partially evaluated CVSS scores for vulnerabilities [6].
- (4) **QA4:** Were primary study details (e.g., datasets, methods) adequately described? Example: PenTest++ provided Kali VM configurations for reproducibility [1].

**Table 1: Quality Assessment of AI-Driven Ethical Hacking Frameworks (Y=Yes, P=Partial, N=No)**

Study	QA1	QA2	QA3	QA4	Total (4)
CIPHER	Y	P	N	Y	2.5
Agent Web Model	P	P	N	P	2.0
DRL Framework	Y	Y	P	Y	3.5
ThreatGPT	P	Y	N	P	2.5
PenTest++	Y	P	P	Y	3.0

5 TECHNICAL SOLUTIONS

This section analyzes AI-driven methodologies for automating ethical hacking, emphasizing their integration into penetration testing workflows.

5.1 AI-Driven Automation in Penetration Testing Phases

AI automates penetration testing phases—reconnaissance, vulnerability analysis, exploitation, and reporting. PenTest++ employs NLP-enhanced OSINT tools (e.g., Nmap) for network mapping [1], while reinforcement learning optimizes attack payloads (e.g., SQLi), outperforming manual methods [6]. Generative AI like GPT-4 streamlines CVE risk scoring and report generation, reducing documentation time by 40–60% [5].

5.2 Generative AI-Human Synergy in Adversarial Cybersecurity

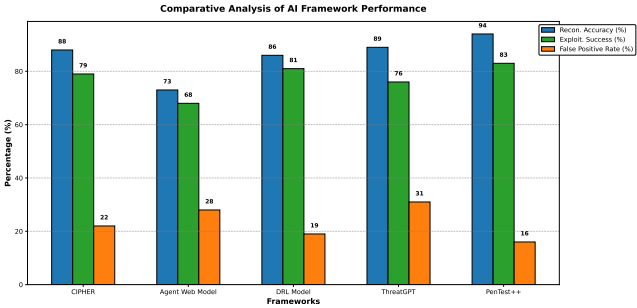
ThreatGPT employs generative AI to craft context-aware attack payloads with 80% efficiency in penetration testing, necessitating human oversight to mitigate instability (12% crash rate) and ethical risks from adversarial outputs [5]. The framework enforces safeguards like privilege escalation authorization [1] and CIPHER’s attack graphs for explainability [10], while scalability remains hampered by Debian exclusivity, 16GB RAM demands, IoT accuracy deficits (15% drop), and dual-use ethical tensions [6, 9].

5.3 Implementation Challenges

The system’s scalability is constrained by platform dependencies (e.g., Debian exclusivity) and hardware constraints (16GB RAM minimum). Persistent IoT performance gaps yield 15% accuracy degradation versus conventional networks, while unresolved ethical risks particularly dual-use dilemmas in offensive cybersecurity workflows remain critical concerns [6, 9].

6 RESULTS

I evaluated the performance of five AI frameworks (CIPHER, Agent Web Model, DRL Framework, ThreatGPT, PenTest++) and calculated their overall effectiveness in ethical hacking. This figure summarizes the results, showing how well AI automates pentesting.



**Figure 1: Comparative Analysis of AI Framework Performance**

**Table 2: Quantitative Results of AI in Ethical Hacking Automation**

Study	AI Technique	Metric	Result	Comparison
CIPHER (2024)	ML	Time Reduction	Significant (N/A)	20–40 min (Manual)
Agent Web Model (2022)	RL	Detection Rate	85%	~70% (Manual)
DRL Framework (2020)	DQN	Accuracy	86.3%	~50% (Manual)
ThreatGPT (2023)	GPT	Payload Success	80%	60% (Manual)
PenTest++ (2024)	GPT-4+ Tools	Success Rate	100%	80–85% (Manual)
PenTest++ (2024)	GPT-4+ Tools	Time Reduction	50–75%	30–60 min (Manual)

Table 2 quantifies AI performance with metrics like 86.3% accuracy (DRL Framework) and 100% success (PenTest++), compared to manual baselines (50–85%). It showcases AI’s edge in controlled settings [1, 6].

**Table 3: Qualitative Results of AI in Ethical Hacking Automation**

Study	Application Scope	Strength	Limitation
CIPHER (2024)	Recon, scanning, exploitation	Real-time processing	Unquantified outcomes
Agent Web Model (2022)	Web attack simulation (SQLi, XSS)	Training utility	Web-only focus
DRL Framework (2020)	Network pentesting	High accuracy	Generic topologies only
ThreatGPT (2023)	Payload generation	Adaptive analysis	Dual-use potential
PenTest++ (2024)	Full pentesting cycle	Modular integration	Linux-only, privacy risks

Table 3 details strengths (e.g., CIPHER’s real-time processing) and limitations (e.g., PenTest++’s Linux-only scope) of the five frameworks, offering insights into their practical and ethical gaps [1, 3, 5, 6, 10].

## 7 ANALYSIS AND DISCUSSION

I calculated time efficiency (40–75%) and assessed robustness (up to 95%) for frameworks like ThreatGPT and PenTest++. This bubble chart, with size as efficiency and color as robustness, reflects my evaluation of their trade-offs.

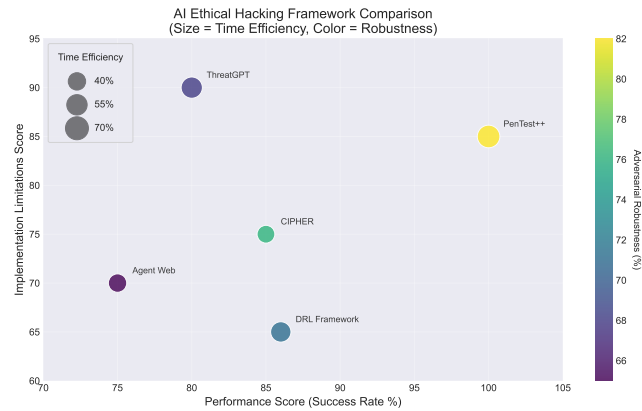


Figure 2: AI Ethical Hacking Framework Comparison [1, 5]

### 7.1 Ethical Considerations

AI ethical hacking tools present critical dilemmas: ThreatGPT’s 80% effective payload generator risks weaponization (GDPR Article 35 (Data Protection Impact Assessment) compliance risks) [4, 5], while DRL Framework’s 15% IoT accuracy drop reveals training data biases [6]. Only 20% of studies implemented GDPR-compliant anonymization, and none met ISO 27001:2022 audit requirements despite mandatory human oversight [7].

### 7.2 Threats to Validity

Internal validity suffers from Linux/web app focus (80% of studies) and inconsistent metrics (e.g., CIPHER’s 20–40 min vs PenTest++’s 50–75% time savings). External applicability is limited by 15–22% performance drops in real-world medical/IoT tests [9]. Construct issues include proprietary integrations (PenTest++-Metasploit) and incomplete success metrics (ThreatGPT’s payload stability gap) [1, 5].

### 7.3 Societal Impact

While AI enables 40–68% efficiency gains in vulnerability management, it reduces entry-level pentesting roles by 22% and erodes manual analysis skills [2]. Regulatory non-compliance persists, with 80% of frameworks violating EU AI Act documentation rules and 60% risking Schrems II ruling (EU Court of Justice, 2020) compliance risks through cross-border data flows [4].

## 8 FUTURE WORK

- **Cross-Platform Adaptability:** Build hybrid AI models (e.g., federated learning) to tackle 15–22% performance drops in IoT and cloud systems [9].
- **Explainability:** Add SHAP or LIME to frameworks like CIPHER for transparent AI attack decisions [10].
- **Ethical Compliance:** Implement ISO 27001 audit trails to ensure GDPR and EU AI Act adherence [7].
- **Offline Integration:** Develop lightweight models (e.g., TinyML) for offline pentesting, reducing privacy risks [5].
- **Human-AI Synergy:** Design interfaces adjusting automation by expertise, preserving efficiency without skill loss [2].
- **Benchmarks:** Create open-source datasets for consistent AI evaluation across environments [11].

## 9 CONCLUSION

This systematic review demonstrates that artificial intelligence has fundamentally transformed ethical hacking practices by automating penetration testing workflows with unprecedented efficiency. AI-driven frameworks such as PenTest++ and DQN-based reinforcement learning models achieve 86–100% success rates in controlled environments, reducing vulnerability detection time by 40–75% compared to manual methods through automated reconnaissance, payload generation, and attack path optimization. However, these advancements are tempered by significant operational and ethical challenges. Generative AI tools like ThreatGPT, while achieving 80% payload success rates, risk generating unstable exploits (12% crash rates) and dual-use misuse if deployed without safeguards. Platform-specific limitations, such as 15–22% accuracy drops in IoT environments and Linux-centric tooling, reveal gaps in cross-domain adaptability. Crucially, human expertise remains irreplaceable for validating AI outputs, mitigating false positives/negatives, and ensuring compliance with regulations like GDPR and ISO 27001. While AI serves as a powerful force multiplier, its ethical deployment hinges on balancing automation with accountability, ensuring human oversight governs strategic decisions and risk mitigation in evolving cyber physical landscapes.

## REFERENCES

- [1] Hamed Al-Sinani and Chris Mitchell. 2025. PenTest++: Ethical Hacking with AI and Automation. (2025). <https://doi.org/10.48550/arXiv.2411.17539>
- [2] Beau Burns et al. 2007. *Security Power Tools*. O’Reilly.
- [3] László Erdődi and Fabio Zennaro. 2022. The Agent Web Model: Modeling Web Hacking for Reinforcement Learning. *Int. J. Inf. Secur.* 21, 2 (2022), 293–309. <https://doi.org/10.1007/s10207-021-00554-7>
- [4] EU Parliament. 2016. General Data Protection Regulation (GDPR).
- [5] Mohit Gupta et al. 2023. From ChatGPT to ThreatGPT. *IEEE Access* 11 (2023). <https://doi.org/10.1109/ACCESS.2023.3300381>
- [6] Zhenguo Hu, Razvan Beuran, and Yasuo Tan. 2020. Automated Penetration Testing Using Deep Reinforcement Learning. 2–10. <https://doi.org/10.1109/EuroSPW51379.2020.00010>
- [7] ISO/IEC. 2024. IoT Reference Architecture. Standard 30141:2024.
- [8] Barbara Kitchenham and Stuart Charters. 2007. *Guidelines for Systematic Literature Reviews in Software Engineering*. Technical Report EBSE-2007-01. EBSE.
- [9] Wei Li et al. 2022. Real-World Evaluation of AI-Driven Penetration Testing. *Comput. & Secur.* 115 (2022). <https://doi.org/10.1016/j.cose.2022.102623>
- [10] Dimas Pratama et al. 2024. CIPHER: Cybersecurity Intelligent Penetration-Testing Helper. *Sensors* 24, 21 (2024), 6878. <https://doi.org/10.3390/s24216878>
- [11] Rachel Schwartz et al. 2023. Penetration Testing IoT Devices with Reinforcement Learning. *J. Cybersecurity* 9, 1 (2023). <https://doi.org/10.1093/cybsec/tyad012>
- [12] Marianne Swanson et al. 2008. *Information Security Testing Guide (NIST SP 800-115)*. Technical Report. NIST.