A Survey on Web Application Testing: A Decade of Evolution

TAO LI, School of Computer Science and Engineering, Macau University of Science and Technology, China RUBING HUANG, School of Computer Science and Engineering, Macau University of Science and Technology, China and Macau University of Science and Technology Zhuhai MUST Science and Technology Research Institute, China

CHENHUI CUI, School of Computer Science and Engineering, Macau University of Science and Technology, China

DAVE TOWEY, School of Computer Science, University of Nottingham Ningbo China, China

LEI MA, School of Science, University of Tokyo, Japan

YUAN-FANG LI, Faculty of Information Technology, Monash University, Australia

WEN XIA, School of Computer Science and Technology, Harbin Institute of Technology Shenzhen, China

As one of the most popular software applications, a web application is a program, accessible through the web, to dynamically generate content based on user interactions or contextual data, for example, online shopping platforms, social networking sites, and financial services. Web applications operate in diverse environments and leverage web technologies such as HTML, CSS, JavaScript, and Ajax, often incorporating features like asynchronous operations to enhance user experience. Due to the increasing user and popularity of web applications, approaches to their quality have become increasingly important. Web Application Testing (WAT) plays a vital role in ensuring web applications' functionality, security, and reliability. Given the speed with which web technologies are evolving, WAT is especially important. Over the last decade, various WAT approaches have been developed. The diversity of approaches reflects the many aspects of web applications, such as dynamic content, asynchronous operations, and diverse user environments. This paper provides a comprehensive overview of the main achievements during the past decade: It examines the main steps involved in WAT, including test-case generation and execution, and evaluation and assessment. The currently available tools for WAT are also examined. The paper also discusses some open research challenges and potential future WAT work.

CCS Concepts: • Software and its engineering \rightarrow Software verification and validation; • Information systems \rightarrow Web applications; • General and reference \rightarrow Surveys and overviews.

Additional Key Words and Phrases: Software testing, web application testing, survey

Authors' addresses: Tao Li, 3220007015@student.must.edu.mo, School of Computer Science and Engineering, Macau University of Science and Technology, Taipa, Macau, China, 999078; Rubing Huang, rbhuang@must.edu.mo, School of Computer Science and Engineering, Macau University of Science and Technology, Taipa, Macau, China, 999078 and Macau University of Science and Technology Research Institute, Zhuhai, Guangdong Province, China, 519099; Chenhui Cui, 3230002105@student.must.edu.mo, School of Computer Science and Engineering, Macau University of Science and Technology, Taipa, Macau, China, 999078; Dave Towey, dave.towey@nottingham.edu.cn, School of Computer Science, University of Nottingham Ningbo China, Ningbo, Zhejiang, China, 999078; Lei Ma, ma.lei@acm.org, School of Science, University of Tokyo, Tokyo, Japan, 113-0033; Yuan-Fang Li, yuanfang.li@monash.edu, Faculty of Information Technology, Monash University, Melbourne, Victoria, Australia; Wen Xia, xiawen@hit.edu.cn, School of Computer Science and Technology, Harbin Institute of Technology Shenzhen, Shenzhen, China, 518055.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

0004-5411/2018/8-ART111 \$15.00

ACM Reference Format:

1 INTRODUCTION

The Software Development Life Cycle (SDLC) [161] is composed of several essential phases, including requirements specification, architecture design, coding, testing, and maintenance [76]. Each phase involves specific processes and activities that contribute to the overall development of the software. Among these phases, testing is particularly critical, often accounting for a significant proportion of time and cost for developing a software application. Recent studies [29, 107] indicate that the testing process can consume a range of 10% to 80% of total development time, depending on the complexity of the project. The primary objective of testing is to identify errors that manifest during program execution. An effective testing methodology improves the probability of detecting such errors, for example, by systematically running and examining a wide range of execution paths, including edge cases and fault-prone areas [104, 111].

Since the adoption of the World Wide Web in the early 1990s [26], web applications have become integral to various sectors, including finance, healthcare, education, and commerce [50]. These applications have introduced innovative features and functionalities that facilitate consistent and efficient access to information across diverse sources, significantly transforming how society interacts with digital content [84]. Meanwhile, web applications rely on a set of core technologies, including HTML to define the structure of web pages, CSS to style, JavaScript to enable interactivity, and Ajax for asynchronous data exchanges between the client and the server. These technologies are the foundation of modern web applications and are critical in supporting their dynamic and distributed nature [19, 64, 158].

As web applications grow increasingly complex and technologies continue to evolve, the requirement for robust and effective testing methodologies has become more crucial atkinson2002towards: It is necessary to continuously evolve and advance testing strategies to ensure the reliability, security, and overall quality of web applications [104]. Web Application Testing (WAT) [50] has become an essential process in the SDLC of web applications, which is used to evaluate and ensure the quality of Applications Under Test (AUTs) [86], making it a cornerstone of modern software quality assurance. WAT aims to ensure that web applications function as intended, meet specified requirements, and are free from vulnerabilities. The scope of WAT goes beyond bug detection and fixing, and includes enhancing user experience, ensuring data security, and maintaining the overall reliability of the web systems. As web technologies evolve, testing methodologies have advanced to address web applications' increasing complexities and demands.

Web applications have had an important impact on society [63], reshaping various aspects of daily life and industry [138]. Despite their significance, however, WAT reviews have tended to be narrowly focused [145], addressing specific aspects without providing a comprehensive overview [91]. One WAT survey paper, by Prazina et al. [145], only concentrates on the testing and analysis methods for automated web layouts. In contrast, significant attention has been devoted to the security aspects of web applications [9, 21, 41, 70, 93, 112, 126, 149, 150, 181, 193]. Additionally, Gupta et al. [72] focused on the automated generation of regression test cases for web applications.

Early reviews of WAT, such as those by Li et al. [104] and Dougan et al. [50], established foundational insights into testing dynamic web environments, focusing on early challenges and techniques. More recently, Balsam and Mishra [23] have provided a review to discuss some challenges and opportunities in WAT. However, this paper focused on only 72 primary studies, which means that

many significant contributions and broader advancements, over the past decade, have been overlooked. To address this gap, we aim to provide a comprehensive examination of the latest progress in WAT, detailing core processes, evaluating current tools, and exploring emerging research directions to offer a thorough overview of the field's evolution over the past decade.

This paper provides a comprehensive examination and discussion of WAT across its entire life cycle. We explore the research status and WAT advances from the perspectives of test case generation and execution, failure diagnosis, evaluation and assessment, regression testing, and available tools. We also examine some key open research challenges and discuss potential directions for future WAT work.

The rest of this paper is organized as follows: Section 2 provides background information on the formal definition of WAT and the fundamental testing techniques employed in the field. Section 3 outlines the survey methodology, including presenting the six research questions, the literature retrieval process, and the statistical analysis of the retrieved studies. Sections 4 to 9 address each of the research questions (RQ1 to RQ6), respectively. Finally, Section 10 concludes the paper by summarizing key findings and proposing directions for future research.

2 BACKGROUND

This section presents some preliminary WAT concepts and provides an overview of web applications and WAT, covering some core definitions and relevant testing techniques.

2.1 Web Applications

A web application [84] is a software system that facilitates user interaction through a combination of front-end and back-end components, over a network, typically accessed through a web browser. The front-end, or client-side, renders the user interface and handles user inputs. It consists of web pages constructed using technologies such as HTML, CSS, and JavaScript. This layer is designed to provide a responsive and interactive user experience, enabling seamless communication between the user and the application. The back-end, or server-side, encompasses the core functionalities of the application: These include data processing, business logic execution, and database management. The back-end is the key to the application: It processes requests from the front-end and ensures that the correct data and functionalities are provided to the user. Communication between the front-end and back-end typically relies on protocols such as HTTP and HTTPS, which support smooth data exchange and maintain the application's state [92, 130]. Together, these components provide an integrated system that supports a wide range of user interactions, from simple content browsing to complex transaction processing.

Web applications can be either static or dynamic. Static web applications have fixed content that remains unchanged regardless of user interactions, making them suitable for delivering straightforward information, such as in digital brochures or basic company websites [176]. Dynamic web applications, in contrast, provide a more interactive experience, with content that adapts to user inputs, interactions, or real-time data [158]. The dynamic nature of these applications allows for greater flexibility and user engagement, making them well-suited for complex platforms like social media networks, e-commerce sites, and online service portals that require continuous interaction and real-time responsiveness [59].

2.2 Web Application Testing

Web applications have become a fundamental part of our daily lives, making thorough testing a vital process to ensure their functionality and dependability. Effective testing ensures that these applications not only meet the required functionality but also have sufficient security, and good performance and reliability. Thorough testing of web applications can help to identify and address

potential vulnerabilities, performance bottlenecks, and functional issues: This enhances the user experience and helps to safeguard sensitive data, which is vital for maintaining trust in web-based services. Comprehensive testing ultimately contributes to the development of robust, dependable web applications that can adapt to the dynamic demands of users and the evolving technological landscape.

WAT involves several key steps:

- Analysis of the functional and non-functional requirements of the web application to guide the development of comprehensive test cases.
- Generation of test data according to the testing methodology (such as model-based or security testing), and appropriate design of test cases, to ensure a broad coverage of potential scenarios.
- Execution of the test cases by sending HTTP requests to the web application, and comparison of the responses with the expected outcomes to identify inconsistencies.
- Documentation of the test results, analysis of any issues, and application of regression testing to ensure that recent changes have not introduced new defects.
- Implementation of continuous testing and monitoring to maintain the application's stability and security as it evolves.

3 METHODOLOGY

In this section, we present the six research questions (RQs) related to WAT that structured our study, and this paper. Our WAT survey involved a systematic and structured methodology, guided by the published frameworks of Huang et al. [80] and Webster et al. [187]. Our findings' presentation draws from recent surveys on related topics, including the testing of RESTful APIs [66], mutation analysis [85], and metamorphic testing [38, 160]. The following sections detail key aspects of the review process and results.

3.1 Research Questions

The goal of this survey paper was to identify and categorize the available WAT information. To achieve this, the following six RQs guided the study:

- RQ1: How has WAT evolved, and what WAT topics have been examined in published studies?
- **RQ2**: What methods and approaches have been used for WAT test case generation?
- **RQ3:** How are test cases executed in WAT?
- **RQ4:** What metrics are used to evaluate WAT? How are they used?
- **RQ5**: What tools are available to support WAT?
- **RQ6:** What are the challenges and future work related to WAT?

An answer to RQ1 should provide an understanding of the WAT literature and the distribution of its key topics. RQ2 prompted an examination of the various methods and approaches used for WAT test case generation, revealing the application of a diversity of techniques. RQ3 focuses on the execution of WAT test cases, aiming to clarify the processes and practices involved. RQ4 should reveal the metrics and evaluation frameworks used in WAT, enabling an understanding of how the WAT effectiveness and efficiency are assessed. RQ5 should reveal the currently available WAT tools, offering a comprehensive overview of the technological resources for WAT. Finally, based on the insights gathered from the other RQs, RQ6 shall identify unresolved challenges and propose potential directions for future WAT research. Figure 1 presents an overview of the structure of this paper.

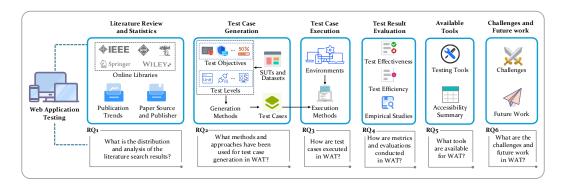


Fig. 1. Structure of this survey paper.

Table 1. Selected Digital Libraries, with Search Queries

No.	Name	Search Query
1	ACM	[[Title: "web"] OR [Title: "browser based"]] AND [[Title: "test"] OR [Title: "testing"] OR [Title: "measure"] OR [Title: "measurement"] OR [Title: "measuring"] OR [Title: "check"] OR [Title: "checking"] OR [Title: "detection"]]
2	Elsevier	(("web") OR ("browser based")) AND (("test") OR ("testing") OR ("measure") OR ("check") OR
		("checking") OR ("detect") OR ("detecting"))
3	IEEE	("Document Title": "web" OR "Document Title": "browser based") AND ("Document Title": "test" OR
		"Document Title":"testing" OR "Document Title":"measure" OR "Document Title":"measurement" OR
		"Document Title":"measuring" OR "Document Title":"check" OR "Document Title":"checking" OR
		"Document Title": "detect" OR "Document Title": "detecting" OR "Document Title": "detection")
4	Springer	"web" AND ("test" OR "testing" OR "measure" OR "measurement" OR "measuring" OR "check" OR
		"checking" OR "detect" OR "detecting" OR "detection")
	Wiley	"("web") OR ("browser based")" in Title and "("test") OR ("testing") OR ("measure") OR
5		("measurement") OR ("measuring") OR ("check") OR ("checking") OR ("detect") OR ("detecting")
		OR ("detection")" in Title

3.2 Literature Search and Selection

Similar to previous survey studies [15, 66, 79, 80, 187], we selected the following online repositories of technical research literature:

- ACM Digital Library (ACM)
- Elsevier Science Direct (Elsevier)
- IEEE Xplore Digital Library (IEEE)
- Springer Online Library (Springer)
- Wiley Online Library (Wiley)

These repositories were chosen for their extensive collections of conference, symposium, and workshop papers, as well as for their access to multiple journals that are considered highly relevant to WAT [80]. Our survey examined computer science papers published between January 1, 2014, and December 31, 2023.

Once the literature repositories were identified, search strategies were developed and tailored for each of them, using WAT-specific terminology and formatting. Given the unique characteristics and limitations of the advanced search functions in each repository (such as Elsevier's restriction on using a maximum of eight Boolean connectors per field), the search methods were designed to strictly adhere to the repository rules and constraints. To ensure comprehensive coverage of all relevant literature, multi-dimensional and multi-layered filtering keywords were used to maximize the retrieval of the WAT materials. Table 1 illustrates the specific search queries for five selected digital libraries.

No.	Name	No. of studies from the search keywords-based results	No. of studies excluded based on the selection criteria	No. of studies after the selection criteria
1	ACM	232	138	94
2	Elsevier	27	19	8
3	IEEE	343	171	172
4	Springer	65	33	32
5	Wiley	19	11	8
Total		686	372	314

Table 2. Selection Results of Primary Studies

The survey included studies written in English that were focused on WAT, but excluded academic theses (e.g., master's and doctoral dissertations), books, posters, and previous survey studies. In addition, some processes relevant to WAT, such as failure diagnosis and regression testing for web applications, were excluded from the survey. The inclusion criteria were:

- The paper was written in English.
- The paper was related to WAT.
- The paper was not a thesis.
- The paper was not a survey or a systematic literature review.
- The paper was freely available in open access.

After duplicates were removed and the exclusion criteria applied, the initial set of candidate studies was significantly refined. A snowballing approach [66], which involved examining the references lists of in the selected studies, led to the identification and inclusion of several additional papers. Table 2 summarized the details of the search and filtering process. Eventually, 314 papers were included for the preliminary review and statistical analysis.

Although some papers may have been omitted due to the focus on a subset of reputable publishers, we are confident that the overall trends we report on are accurate and provide a fair representation of the current state of the art in WAT.

3.3 Data Extraction and Collection

Key information was gathered from each study regarding the research motivation, contributions to the field, empirical-evaluation details, common misconceptions, and remaining challenges in WAT. These were reviewed and verified by co-authors of the survey, minimizing the risk of overlooking critical information, and reducing potential errors in the analyses.

We organized the RQ data collection methods as follows:

- **RQ1:** Fundamental information for each paper such as publication year, type of paper, and paper source.
- **RQ2:** Testing AUT, testing objective, and test case generation techniques.
- **RQ3:** Execution methods, tools and platforms, environments, tested items, dataset URL, and open-source status.
- **RQ4:** Evaluation metrics, methods and approaches, and tools and techniques.
- RQ5: Tool name, type, main features and capabilities, URL, and open-source status.
- **RQ6**: Overview of the remaining challenges.

4 ANSWER TO RQ1: PUBLICATION EVOLUTION AND DISTRIBUTION

This section addresses RQ1 through a detailed analysis of publication trends and their distribution across the literature sources. Our methodology provides a comprehensive analysis of the evolution and distribution of WAT research over the past decade.

4.1 Publication Trends

We collected publication year data from 314 papers (as illustrated in Figure 2) to highlight trends in WAT research from 2014 to 2023.

Figure 2(a) shows the annual number of publications, and Figure 2(b) shows the cumulative total. Figure 2(a) shows the fluctuating (but relatively steady) nature of WAT-related publications from 2014 to 2023. The annual publication counts exhibit some variation over the decade, with 32 publications in 2014, followed by slight fluctuations: 35 in 2015, 25 in 2016, 27 in 2017, 29 in 2018, and 34 in 2019. In 2020, the count was 31, which increased to 38 in 2021. This was followed by a decrease to 28 in 2022 and a subsequent rise to 35 in 2023. This pattern reflects a generally stable trend in WAT research output over the years.

The cumulative publication data presented in Figure 2(b) reflects a consistent growth in research contributions. The trend in Figure 2(b) suggests that WAT has maintained a steady pace of development, driven by ongoing technological advancements and sustained research interest in the field.

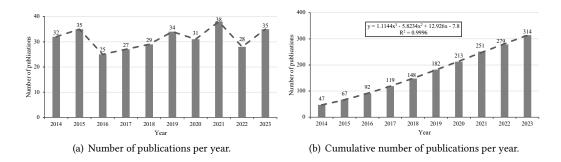


Fig. 2. WAT papers published between January 1, 2014, and December 31, 2023.

4.2 Paper Source and Publisher

The papers analyzed in this study were from various journals and conferences. As shown in, the majority of the publications were from conferences (71%), with a smaller proportion being from journals (29%). The data presented in Figure 3(b) reveals distinct patterns in publication volumes over the years. For instance, conference papers exhibited fluctuations, starting at 25 in 2014, rising slightly to 27 in 2015, and then decreasing to 17 in 2018 before recovering to 26 in 2023. In contrast, journal publications, while also fluctuating, reached their peak of 15 papers in 2019, demonstrating a steadier but less pronounced growth compared to conference papers.

5 ANSWER TO RQ2: TEST CASE GENERATION

This section addresses RQ2 by systematically examining key aspects of WAT test case generation, including the *Applications Under Test* (AUTs) and datasets, testing objectives, and test case generation methods.

5.1 AUTs and Datasets

5.1.1 AUTs. This section reviews the types of AUTs used in WAT, referencing the selected research papers. In WAT, AUTs refer to the web platforms or systems being evaluated (for functionality, performance, security, and usability). The diversity of WAT AUTs emphasizes the need for tailored testing approaches to address specific challenges across different domains: E-commerce platforms,

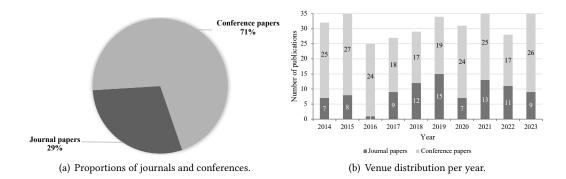


Fig. 3. Venue distribution of surveyed papers.

for example, require stringent security measures, whereas government portals may focus on accessibility.

- *E-commerce Platforms:* E-commerce applications such as Magento, Shopify, and custom-built online stores are popular WAT AUTs. This may be due to their complex transactional workflows, which involve user authentication, shopping carts, and payment gateways [144]. Testing often focuses on functional correctness and security, with particular attention to vulnerabilities to potential attacks like *SQL Injection* (SQLi) and *Cross-Site Scripting* (XSS) [30]. Performance testing is also essential to ensure that the platforms can handle high traffic volumes during peak times [180].
- Content Management Systems (CMSs): CMS platforms such as WordPress, Drupal, and Joomla are frequently chosen as AUTs because of their modular and extensible nature, which can potentially make them vulnerable, through third-party plugins [18, 61]. Testing typically emphasizes plugin vulnerabilities, core stability, and validating performance after updates [102].
- Government and Public Service Portals: Government web applications, such as university and public service portals, are often tested for accessibility, usability, and security [137]. These platforms often need to comply with accessibility standards such as the Web Content Accessibility Guidelines (WCAG) [17], making accessibility testing critical [177]. Security testing emphasizes the prevention of unauthorized access and data breaches, as illustrated by frameworks targeting GDPR compliance in web portals [11].
- Social Media Platforms: Social media platforms, like Facebook and Twitter, pose unique challenges for WAT, due to their real-time interaction and dynamic content. Testing for these applications often examines performance under heavy user loads, cross-browser compatibility, and session management vulnerabilities [16, 55].
- Web Application Firewalls (WAFs): WAFs play a key role in defending web applications from attacks like SQLi, XSS, and Cross-Site Request Forgery (CSRF) [142, 196]. Testing WAFs involves simulating real-world attack scenarios to evaluate how effectively they distinguish between legitimate and malicious traffic [37, 197].
- 5.1.2 Datasets. Datasets are essential in WAT, providing the inputs needed to simulate real-world usage and evaluate application behavior across key areas such as security, functionality, and performance. The datasets are often comprised of user data, predefined test cases, or automatically

Table 3. An overview of datasets.

Dataset Name	Overview	Resource Link	Reference
Software Assurance Reference Dataset (SARD)	It is a dataset consisting of PHP test cases, where each test case represents a vulnerable or non-vulnerable SQL-related code snippet. The test cases include examples of user input propagating through the application code until a sensitive function (e.g., a database query) is executed.	https://samate.nist.gov/SARD/index.php	[57]
Near-Duplicate Study DataSet	It is a dataset of about 100,000 annotated pairs of web pages from open-source web applications and real websites. Each pair is labeled as distinct, near-duplicate, or clone to facilitate the evaluation of near-duplicate detection methods.	https://zenodo.org/records/ 3376730	[40]
Website Screen- shots Dataset from Roboflow	It is a dataset of 1,206 screenshots from popular websites with 54,215 UI element annotations such as buttons, text, images, links, titles, fields, iframes, and labels. This dataset is used to train object detection models and evaluate the visual consistency of web applications.	https://public.roboflow.com/ object-detection/website- screenshots	[1]
PhishTank dataset	It is a widely used dataset containing known phishing URLs.		[135]
Web injection publicly available datasets	The datasets consist of various HTTP request samples labeled according to the type of attack they represent, including SQL injection, XSS, command injection, and path traversal.	SQL Injection: (1) https://www.kaggle.com/syedsaqlainhussain/sql-injection-dataset, (2) https://github.com/mhamouei/rat_datasets, XSS Injection:https://www.kaggle.com/syedsaqlainhussain/cross-site-scripting-xss-dataset-for-deeplearning, HTTP Parameters:https://github.com/Morzeux/HttpParamsDataset	[13, 154]
Mitch Framework HTTP Request Dataset	It is a dataset based on the usage of the Mitch framework, which contains 5,828 HTTP requests but is reduced to 1901 records after balancing.	https://github.com/alviser/mitch	[24]
NASA 1995 Apache Web Log	It is a dataset containing log files from NASA web servers, recording HTTP requests from real users during July 1995.	https://ita.ee.lbl.gov/html/contrib/ NASA-HTTP.html	[141]
DBpedia	It is a structured knowledge base that extracts information from Wikipedia. It contains millions of entities represented as RDF triples.	https://www.dbpedia.org/	[113]
fuzzdb	It is an open-source database containing a large number of attack patterns.	https://github.com/fuzzdb- project/fuzzdb	[194]
ATO-data	It is a collection of keystroke dynamics data collected un- der controlled conditions in a lab and anonymized data from 187,332 sessions of financial web pages, which supports ana- lyzing the changes in different user behaviors and is used to evaluate the effectiveness of imposter detection models.	https://github.com/mluckner/ ATO-data.git	[67]
Web Accessibility Evaluation Data	It is an automated web accessibility assessment of nearly 3 million web pages obtained between March and September 2021 using Qual Web.	https://zenodo.org/records/ 7494722	[115]
PHP-Webshell- Dataset	It is a dataset containing 2,917 samples from 17 webshell collection projects.	https://github.com/Cyc1e183/ PHP-Webshell-Dataset	[65]
PHP Vulnerability Test Suite	It is a dataset containing 9408 PHP source code samples, in- cluding 5600 safe samples and 3808 unsafe samples, which is widely used for the detection of XSS vulnerabilities and the training of prediction models.	https://github.com/stivalet/PHP- Vulnerability-test-suite	[71]
Datasets of three common web attack payloads (SQLi, XSS, RCE)	It is a dataset generated by attack grammar, including SQL injection, cross-site scripting, and remote command execution attack payloads, used to train and test the performance of GPTFuzzer on WAF.	https://github.com/ hongliangliang/gptfuzzer	[105]
Joomla! and Man- tisBT test suite	It is a data set that contains multiple representative use cases, covering functions such as article management and user management. It has 47 manually written test cases and 453 test steps.	https://zenodo.org/record/ 4973219	[90]
GHTraffic	It is a dataset of significant size comprising HTTP transac- tions extracted from GitHub data (i.e., from 04 August 2015 GHTorrent issues snapshot) and augmented with synthetic transaction data.	https://zenodo.org/record/ 1034573	[28]

generated data. They support a range of testing scenarios, including edge cases, aiming to identify potential vulnerabilities and performance bottlenecks.

Each WAT dataset was designed to meet specific testing objectives, allowing for targeted assessments of various aspects. For example, the *Software Assurance Reference Dataset* (SARD) [57] can help to identify security vulnerabilities, such as SQLi, thus contributing to more secure software development practices. The PhishTank dataset [135], which contains real-world phishing URLs, can support evaluation of phishing detection systems, enhancing web application security.

Datasets like the Website Screenshots Dataset [1] facilitate the assessment of visual consistency and UI accuracy, helping to ensure an optimal user experience. The NASA web log data [141] provides real-world HTTP requests, offering valuable insights into traffic patterns and server performance under various loads. DBpedia [113] supports the validation of semantic web applications: It provides structured data that can be used to help verify data integration and usage.

These datasets enable testers to adopt focused testing strategies that ensure a thorough evaluation of web applications across multiple dimensions. Table 3 summarizes key datasets used in WAT, each with a distinct focus, ranging from security testing.

5.2 Test Objectives

WAT test objectives include detecting bugs, ensuring security, evaluating performance, assessing the functionality of *Graphical User Interfaces* (GUIs), and identifying vulnerabilities (such as for XSS and SQLi).

Figure 4 shows the proportion of the different objectives across the surveyed literature. Web vulnerability detection, including XSS and SQL Injection detection, accounted for the largest proportion (29.94%), emphasizing the necessity of addressing critical security risks. Security testing represented 27.71%, highlighting its importance in protecting applications from security threats. GUI testing accounted for 19.43%. Performance evaluation constituted 12.42%, and bug/defect detection represented 10.51%. The following sections provide a detailed analysis of each testing objective.

5.2.1 Bug/Defect Detection. Detecting and correcting defects is essential for maintaining web application integrity. Techniques such as combining unsupervised learning with embedding layers effectively optimize test case prioritization while reducing redundancy [47]. Modern web applications often involve intricate interdependencies between components, requiring advanced techniques

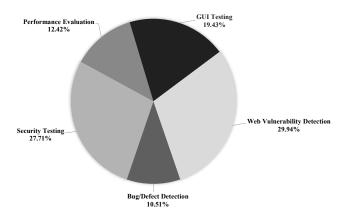


Fig. 4. Proportions of test objectives in surveyed papers.

to detect defects across diverse states and modules. This complexity necessitates strategies such as dynamic test case generation and static analysis to ensure robust coverage [116].

Static code analysis can complement regression testing, and may identify issues early in development. Tools that use static analysis (such as taint analysis) detect potential errors in the source code without execution [116, 136]. For example, the Q-learning-based Selenium-Java tool dynamically explores web GUI elements and learns optimal test paths, significantly improving testing coverage for complex systems [55].

5.2.2 Security Testing. The goal of security testing is to ensure that web applications remain protected against unauthorized access, data breaches, and other security threats. Security testing provides a systematic framework for evaluating the robustness of applications through the identification and mitigation of vulnerabilities before and during runtime [118].

To achieve comprehensive protection, security testing employs a combination of methodologies. Static analysis examines source code for vulnerabilities, such as unhandled exceptions and unsanitized inputs, enabling early detection of issues during development [197]. Dynamic analysis evaluates applications during runtime, simulating potential attack scenarios to uncover vulnerabilities that static analysis might overlook [143]. Hybrid analysis combines static and dynamic methods to address their individual limitations, offering enhanced coverage and reducing false positives [68]. Automated tools, such as Burp Suite, enhance the efficiency of dynamic analysis by enabling real-time evaluation of application defenses [61].

5.2.3 Performance Evaluation. Performance evaluation aims to ensure that an application remains efficient and stable under varying load conditions [142]. Essential to this are load testing and stress testing, which simulate diverse user traffic scenarios to help developers identify performance bottlenecks, especially in environments with high concurrency [108, 142].

Advanced tools such as JMeter [88, 174] and LoadRunner [32] can automate measurement of some important metrics, including response times, throughput, and system resource utilization (e.g., CPU, memory, disk I/O) [55]. They can do this across a spectrum of load levels [55]. These tools can provide comprehensive insights into system behavior under stress, helping the optimization of the system architecture, and ensuring reliability during peak load conditions [95].

5.2.4 GUI Testing. The aim of GUI testing is to ensure that the application's user interface remains consistent, in terms of both functionality and user experience, across different devices and browser environments [195]. Automated tools can simulate user interactions, exploring the GUI's responsiveness under diverse conditions [106]: Cross-browser testing with tools like Selenium Grid [89], for example, can help to ensure compatibility across different platforms [142]. Meanwhile, Scout uses gamified features [59], including progress tracking and dynamic error injection, to encourage systematic exploration and improve test coverage.

Image comparison and computer vision technologies are also extensively used to detect interface changes, and to identify potential UI errors [95]. These tools capture and compare screenshots from different stages of testing, ensuring that updates do not introduce unintended changes to the user interface [197]. The recent integration of AI in testing can further enhance efficiency by learning from user interactions and predicting areas that may be likely to fail [37, 55].

5.2.5 Web Vulnerability Detection. Web vulnerability detection focuses on identifying and mitigating specific threats in web applications, such as Cross-Site Request Forgery (CSRF), directory traversal, Cross-Site Scripting (XSS), and SQL Injection (SQLi). These vulnerabilities, if left unaddressed, can compromise application functionality, data security, and user privacy [25, 196].

To detect these vulnerabilities, static analysis examines source code to identify input validation flaws and hardcoded credentials, enabling early detection of coding issues [108, 197]. Dynamic

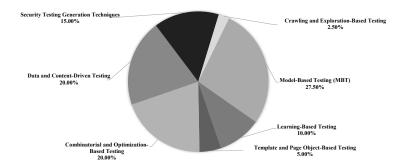


Fig. 5. Proportions of generation methods in surveyed papers.

analysis simulates attack scenarios during runtime, evaluating the application's behavior under unexpected interactions [25]. Hybrid analysis combines the strengths of static and dynamic approaches, providing enhanced detection accuracy and reducing false positives [68].

Modern advancements in vulnerability detection integrate machine learning and deep learning techniques. Methods based on Recurrent Neural Networks (RNNs) analyze input and output patterns to detect malicious behaviors, demonstrating high effectiveness in identifying XSS and SQLi attacks [68]. Additionally, OWASP ZAP is widely used to automate the identification of vulnerabilities, offering strong support for detecting SQLi and XSS [196].

5.3 Generation Methods

WAT can make use of many test case generation methods, including: Model-Based [49], Combinatorial and Optimization-Based [134], Learning-Based [5], Data and Content-Driven [52], Crawling and Exploration-Based [119], Template and Page Object-Based [169], and Security Test Case Generation techniques [56]. These methods offer a broad selection of strategies that can address diverse testing needs.

Figure 5 presents the proportions of some key WAT test generation methods in the surveyed literature. Model-Based Testing (MBT), which uses structural models (e.g., UML, statecharts) that reflect system behavior, is the most prominent, appearing in 27.50% of the papers. Combinatorial and Optimization-Based Testing (20.00%) enhances coverage through input-variable combinations. Data and Content-Driven Testing (20.00%) make use of input validation across diverse scenarios. Security Testing was found in 15.00% of the papers. Learning-Based Testing (10.00%), Template and Page Object-Based Testing (5.00%), and Crawling and Exploration-Based Testing (2.50%) were less commonly reported.

5.3.1 MBT. Model-Based Testing (MBT) automatically generates test cases, enhancing both efficiency and coverage. UML diagrams, statecharts, ER diagrams, and BPMN models are commonly used to represent the system, and to generate the concrete test cases. Panthi & Mohapatra [139] and Suhag & Bhatia [127] used UML models; Indumathi & Begum [81] used ER diagrams and state transition models for database testing. Akpinar et al. [3] showed how MBT could enhance coverage of complex interactions by integrating GraphWalker, Selenium, and JUnit. Moura et al. [45] used BPMN models to generate automated test scripts that covered functional test paths in Business Process Management (BPM) applications. Wang et al. [183] proposed an approach based on UML statecharts for generating Web link security test scenarios. Thummala & Offutt [173] used Petri nets to detect concurrency behaviors. Fard et al. [120] and Waheed et al. [178] showed how model inference and script automation could enhance the automation and test coverage.

- 5.3.2 Combinatorial and Optimization-Based Testing. Techniques based on combinatorial or optimization can generate efficient test cases, which can significantly reduce the number of tests while improving coverage. Qi et al. [148] introduced an automated approach based on combinatorial strategies: This enhanced the coverage of dynamic pages. Bozic et al. [31] reported on the effectiveness of IPOG and IPOG-F combinatorial testing algorithms for detecting XSS vulnerabilities. Wang et al. [186] and Wang et al. [184] applied evolutionary algorithms and combinatorial testing methods to generate test cases covering complex paths in Web applications: Their work demonstrated the potential for genetic algorithms and simulated annealing in Web testing. Wang et al. [185] used combinatorial testing to construct navigation maps for dynamic Web applications, enhancing the test coverage and effectiveness.
- 5.3.3 Al-based Testing. Testing methods can also use machine learning and artificial intelligence (AI) techniques to analyze and generate Web application behaviors and inputs, automating the creation of effective test cases. Chang et al. [37] proposed the WebQT tool based on reinforcement learning, significantly improving testing efficiency and coverage through intelligent path selection. Pavanetto & Brambilla [141] used Generative Adversarial Networks (GANs) to generate highly realistic Web navigation paths, effectively simulating real user operations. Chaleshtari et al. [36] demonstrated the potential of AI techniques in security testing by using mutation testing to detect security vulnerabilities in Web systems. Shahbaz et al. [162] introduced a technique for generating test data based on regular expressions and Web searches, producing valid and invalid WAT data. Li et al. [98] developed the SymJS framework, which combines symbolic execution and dynamic feedback to improve the coverage and accuracy of JavaScript WAT.
- 5.3.4 Data and Content-Driven Testing. Data and content-driven testing methods aim to cover all possible input combinations, to comprehensively verify application functionality and security. Hanna & Jaber [77] proposed a method based on analyzing client-side user input fields, effectively covering various input scenarios. Hanna & Munro [78] used semantically invalid input data to ensure that Web applications reject inputs that do not meet semantic constraints, thereby enhancing security. Shahbaz et al. [162] combined regular expressions with Web searches to automatically generate test cases, verifying string validation logic in Web applications. Nagowah & Kora-Ramiah [128] developed the CRaTCP tool, which automatically extracts control combinations and generates test cases, achieving comprehensive WAT coverage. Mariani et al. [113] utilized the data from applications to automatically generate valid input data (which satisfies the syntactical and semantical requirements) to test complex web applications and ensure the coverage of input combinations.
- 5.3.5 Crawling and Exploration-Based Testing. Methods based on crawling and exploration are particularly suited for testing the complex navigation and state changes of dynamic Web applications. Fard et al. [120] combined automated crawling with manually written test cases to enhance coverage and explore unexplored paths. Santiago Balera and Junior [22] proposed a multi-perspective crawling approach using generative hyper-heuristic algorithms to capture diverse user behaviors and generate comprehensive test cases. Yousaf et al. [191] integrated model-driven techniques with crawling strategies to navigate interfaces and improve test case quality systematically. Jesus et al. [44] leveraged task-driven methods to generate navigation paths aligned with functional workflows. Nguyen et al. [131] applied Q-learning to explore dynamic state spaces, effectively detecting vulnerabilities like XSS. Pavanetto and Brambilla [141] utilized deep learning models, combining RNNs and GANs to simulate realistic navigation paths, enhancing test coverage and interaction fidelity.

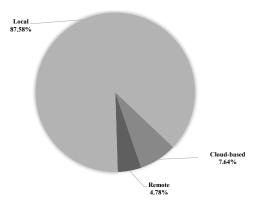


Fig. 6. Proportions of execution environments in surveyed papers.

5.3.6 Template and Page Object-Based Testing. Robust and efficient test cases can be generated using templates or the Page Object Model. Yu et al. [192] used an incremental testing approach based on the Page Object Model that automatically generates page objects to represent Web pages. Neves et al. [43] developed the Morpheus Web Testing tool, which extracts UI component information from JSF and Primefaces frameworks to generate system-level functional test cases, effectively covering complex user interaction scenarios. Leotta et al. [97] utilized meta-heuristic techniques, including greedy and genetic algorithms, to automatically generate robust XPath locators, which reduces the maintenance overhead caused by DOM evolution.

5.3.7 Security Testing Generation Techniques. There are also techniques for generating test cases to detect security vulnerabilities in Web applications. Common vulnerabilities include for XSS and SQLi attacks. Akrout et al. [4] proposed an automated black-box Web vulnerability identification and attack scenario generation method, enhancing vulnerability detection and accuracy. Awang et al. [20] used Cartesian product algorithms [2] to generate SQLi test data, effectively detecting vulnerabilities. Wang et al. [183] used UML statecharts to generate Web link security test scenarios, enhancing robustness in complex security environments. Nguyen et al. [131] used Q-learning algorithms [83], detecting security vulnerabilities, and demonstrating the potential for intelligent algorithms in Web security testing. Wang et al. [184] used evolutionary testing, combining client-side and server-side testing scenarios to coverage and efficiency, and demonstrating the potential for combinatorial optimization algorithms in Web security testing.

6 ANSWER TO RQ3: TEST CASE EXECUTION

This section addresses RQ3, examining the execution of test cases within the context of WAT. The section explores the various methods, tools, platforms, environments, and key considerations of WAT testing: Each aspect is analyzed with a focus on its practical application, and its contribution to test-execution effectiveness.

6.1 Environments

The selection of testing environments is critical for WAT effectiveness, directly influencing the accuracy and scope of performance assessments. According to the literature, testing environments can be categorized into local [186], remote [34], and cloud-based [188].

Figure 6 shows that local environments were the most frequently used, accounting for 87.58%. The control and immediate feedback potential of local environments make them particularly suitable

for iterative testing in early development [53]. Remote environments (4.78%) support testing under various network conditions: This is essential for applications with geographically dispersed users. Cloud-based environments (7.64%) offer scalability and flexibility to meet complex testing demands and support continuous integration in large-scale applications.

6.1.1 Local Environments. The initial stages of development and testing often take place in local environments, which provide a controlled and isolated setting that facilitates rapid iteration and debugging. In these environments, because the web applications are deployed on the same machine that the tests are being executed on, the speed of iteration and bug fixing can be increased. Tools like Apache JMeter are increasingly adopted for testing web applications in local environments [174]. JMeter excels in performance and functional testing and provides a reliable framework to validate isolated components during unit and integration testing, especially in controlled local settings.

A challenge for local environments, however, is that they cannot easily some simulate real-world scenarios like network latencies, distributed system behaviors, and large-scale performance [146]. To address these limitations, tools like Docker have been widely used to create containerized environments that maintain consistency across different stages of development and production [75]. Docker ensures that testing environments are uniform, thus reducing discrepancies among development, testing, and production setups [114].

6.1.2 Remote Environments. Remote environments extend local capabilities by enabling the execution of test cases across multiple physical or virtual machines. Tools like Selenium Grid [22] play an important role in remote environments, enabling cross-browser and cross-platform testing [46]. Selenium Grid helps ensure that web applications function consistently across different environments and configurations [82].

Despite their advantages, remote environments also face challenges, including network latency and the need to synchronize test execution across different systems [198]. Docker, in conjunction with Selenium Grid, helps alleviate these issues by ensuring that the test environments are consistent across remote nodes [75].

6.1.3 Cloud-Based Environments. Cloud-based environments offer the most flexible and scalable solution for large-scale WAT. Platforms such as Amazon Web Services (AWS), Google Cloud Platform (GCP), and Microsoft Azure offer dynamic resource allocation, allowing tests to be executed ondemand across multiple geographical locations [109, 170]. A cloud environment can simulate real-world scenarios (including high traffic and various network conditions), and can support global users [18]. AWS Lambda, for example, is an important tool for cloud environments, providing a serverless architecture that allows testing to run without manual infrastructure management [125]. It can also dynamically expand resources based on test execution requirements, ensuring efficient test runs.

6.2 Execution Methods

Test-case execution has a critical role in the verification of web application functionality, security, and performance. Various methods can be used to execute test cases, each with certain advantages, depending on the nature and complexity of the application. This section provides a detailed exploration of the three primary approaches to test-case execution: manual, automated, and hybrid. Each approach is analyzed for its practical application, highlighting any contributions to the overall effectiveness of WAT test execution.

Figure 7 presents an overview of the s of different approaches to test-case execution in the surveyed literature. Automated execution dominates, appearing in 91.08% of the literature. This is as expected, given its efficiency and scalability for repetitive tasks and regression testing. Automated

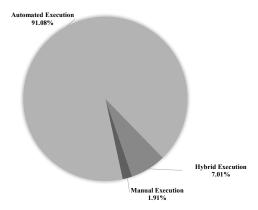


Fig. 7. Proportions of execution methods in surveyed papers.

execution can support tasks requiring consistency across environments. Manual execution, although only in 1.91% of the surveyed papers, remains essential for tasks like exploratory testing and scenarios where human intuition is required. Hybrid execution, accounting for 7.01%, combines both approaches, optimizing test coverage and flexibility by automating routine tasks while leveraging manual testing for subjective analysis. This can be particularly valuable in large-scale applications.

6.2.1 Manual Execution. Manual execution involves human testers directly interacting with web applications to verify their behavior. This can be particularly effective for exploratory testing, as testers need flexibility to dynamically respond to application behavior and detect unexpected issues during testing. Exploratory testing can help to identify user-experience issues, or unexpected errors that are difficult to predict when designing automated scripts [55].

Manual testing is also often used to evaluate new features or complex user interfaces, which can require human judgment to assess usability, aesthetics, and overall user satisfaction [95]. However, there can be significant limitations to manual testing, including that it can be time-consuming, labor-intensive, and prone to human errors. This can be especially the case when dealing with large applications that require the execution of hundreds of test cases [60]. Therefore, although manual testing remains essential for specific tasks, it may be most suitable for early development stages or areas that require subjective validation [73].

6.2.2 Automated Execution. Automated execution is a more scalable and efficient approach, especially when repeating test cases in different environments or during regression testing. Typically, scripts and tools are used to execute test cases, capture results, and compare them with predefined expected outcomes [95]. This is often done without need for further human intervention. Selenium, for example, is widely used to automate functional testing across various browsers, ensuring consistent behavior and compatibility [16].

Automated execution can significantly improve efficiency by reducing human involvement and allowing tests to run continuously in different environments. However, automated execution can require considerable upfront investment in scripting and tool configuration [12]. Furthermore, it may not be suitable when subjective analysis is required, such as for visual-interface testing [60, 108].

6.2.3 Hybrid Execution. Hybrid execution combines both manual and automated techniques, drawing on the strengths of both approaches. In hybrid models, manual testing is typically used for exploratory tasks or validation of new features, while automated scripts handle repetitive tasks and regression testing [142]. This approach offers flexibility by allowing testers to manually

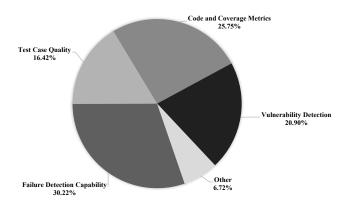


Fig. 8. Proportions of test effectiveness metrics in surveyed papers.

address complex, subjective scenarios while taking advantage of automation for the repetitive, time-consuming processes [60, 68].

It has been reported that combining UML activity diagrams with manual and automated test cases could improve overall test coverage and shorten execution time [197]. Similarly, a hybrid approach can be adopted in security testing, where potential vulnerabilities (such as SQLi points) can be detected automatically, and manual testers can then validate the findings [165]. The combination of the thoroughness of manual testing and the efficiency of automation makes hybrid execution often the most practical solution for complex and large-scale web applications [73].

7 ANSWER TO RQ4: EVALUATIONS AND METRICS

This section addresses RQ4 by exploring WAT metrics and evaluation methods. We first discuss the primary effectiveness metrics, such as test case quality, coverage, failure diagnosis, and vulnerability detection. This is followed by an analysis of efficiency metrics, such the time required for generation, execution, and error detection. Finally, empirical studies are reviewed to explain how these metrics have been applied to evaluate and improve WAT effectiveness and efficiency.

7.1 Test Effectiveness

Test effectiveness evaluates the ability of WAT methodologies to identify faults, vulnerabilities, and failures. This involves the use of specific metrics, including test-case quality, code coverage, failure-detection capability, and vulnerability-detection effectiveness. These metrics provide insights into the thoroughness and accuracy of WAT testing processes.

Figure 8 presents an overview of the different test effectiveness metrics used in the surveyed literature. Failure-detection capability is the most frequently used metric, accounting for 30.22% of the papers. Code and coverage metrics (25.75%) measure the breadth of testing. Vulnerability detection comprises 20.90%, while test-case quality accounts for 16.42%. A small proportion of the literature, 6.72%, involved other metrics: For example, Kumar et al. [157], and Tiwari et al. [174], for example, refer to performance-based metrics such as response time, error rate, and throughput.

7.1.1 Test-Case Quality. The quality of test cases directly affects the ability to detect defects in web applications. Dynamic testing approaches, such as Q-learning-based frameworks, have demonstrated substantial improvements in test-case adaptability. These methods optimize element interaction paths, increase branch and element coverage, and effectively identify critical defects in

complex application environments [55]. These tools provide a scalable solution for testing modern web applications by enhancing exploratory capabilities.

Prioritizing test cases based on risk and modification frequency has enhanced defect detection efficiency. These strategies uncover critical defects earlier by focusing on high-risk and frequently modified components, which streamlines testing processes and improves efficiency and effectiveness [7].

Furthermore, fast evaluation techniques that leverage clustering-based sampling can improve test-case efficiency in large-scale applications. These methods group similar application paths and ensure representative coverage, which significantly reduces redundant testing efforts while maintaining reliability [195].

7.1.2 Code Coverage Metrics. Code coverage metrics — including line, branch, and function coverage — can help show how thoroughly an application's codebase is exercised during testing [132, 185]. Coverage analysis tools like Selenium allow testers to capture interactions across multiple layers of the web application, from the front-end UI to the back-end logic [140].

However, high coverage alone does not guarantee comprehensive fault detection. Research has shown that combining coverage metrics with fault injection methods enhances the ability to uncover defects [36]. Injecting controlled faults into the system, and observing whether or not the test cases can detect them, provides a more rigorous assessment of a test suite's fault-detection capabilities [73]. Mutation testing complements code coverage metrics by introducing controlled changes (i.e., "mutants") to evaluate the fault-detection capabilities of test cases [74, 96]. While coverage metrics measure the breadth of testing, mutation testing ensures that covered paths are rigorously tested, enhancing the reliability of the testing process.

7.1.3 Failure-Detection Capability. Failure-detection capability refers to the accurate identification and localization of faults during testing. Techniques such as dynamic trace-based analysis compare the execution traces of different versions of the application to identify discrepancies [175]. This approach has proven effective in failure detection, particularly when combined with automated tools for analyzing execution logs and detecting potential issues, such as security vulnerabilities, in large-scale applications. This can help testers to pinpoint the specific lines of code or configuration changes responsible for faults, significantly reducing the time required for debugging and resolution [75, 179].

Automated diagnostic tools further improve this process by analyzing execution logs and error traces to identify the exact points of failure [172]. These tools are particularly useful for large-scale web applications, where manual fault-localization could be time-consuming and error-prone [55].

7.1.4 Vulnerability Detection. Vulnerability detection includes the identification of potential security issues such as SQLi and XSS that could be exploited by attackers [20, 122]. Automated vulnerability detection tools simulate real-world attack scenarios, testing the robustness of the application's security defenses [129].

Recent advances in machine learning have enhanced vulnerability-detection capabilities. These systems analyze traffic and behavioral patterns to identify anomalies that may indicate vulnerabilities [82]. As machine learning models learn from new data, they continuously improve their ability to detect both known and unknown threats, providing a dynamic defense against emerging security risks [33].

7.2 Test Efficiency

Test efficiency is a core WAT metric, directly affecting the cost effectiveness of the testing process. Efficiency metrics assess the resources required for comprehensive test coverage and fast fault

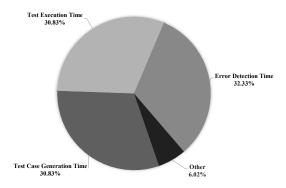


Fig. 9. Proportions of test efficiency metrics in survey papers.

detection. Some of the main metrics include the time taken for test case generation [73] and execution [142], and error detection [55].

Figure 9 presents an overview of the different test efficiency metrics used in the surveyed literature. Test execution time and test case generation time are the most frequently measured metrics, each accounting for 30.83% of the literature. Error detection time is slightly more widely used, appearing in 32.33% of the surveyed papers. A smaller proportion, 6.02%, includes other metrics: For instance, Sung et al. [171] used static analysis methods to reduce redundant tests, thereby improving testing speed and efficiency.

- 7.2.1 Test Case Generation Time. The test-case generation time measures how efficiently test cases are created and is a major part of the overall testing efficiency. Several techniques have been developed to streamline this process. The Page Object Pattern (POP) method has also been able to reduce the test-case generation time as the test suite grows due to the reuse of defined objects. This approach improves efficiency as the same components are used across multiple test cases [95]. Studies involving WebQT, which employs machine learning algorithms, have shown that test-case generation can be significantly accelerated through optimized exploration and test-sequence generation [73]. Tools like CRaTCP can automatically generate and execute all possible test case combinations to further reduce the manual efforts [128].
- 7.2.2 Test-Execution Time. The test-execution time refers to the time taken to execute the generated test cases. Automated testing frameworks like Selenium and Katalon Studio can dramatically reduce execution time (compared to manual testing) [142]. Research has confirmed their effectiveness in optimizing test execution times, particularly in large-scale web applications [142]. Genetic algorithms, for example, have been used to reduce execution time by optimizing test paths, leading to earlier detection of faults [167]. Concurrent execution techniques, where multiple test cases are run in parallel, can also further reduce execution times [60]. Prioritization techniques that execute high-risk test cases first can lead to more efficient use of resources by focusing on the most critical components.
- 7.2.3 Error-Detection Time. The error-detection time reflects how quickly errors or vulnerabilities are identified during testing. Hybrid methods that integrate static and dynamic analysis improve detection efficiency by focusing on critical code paths and minimizing redundant analysis [197]. Techniques leveraging anomaly detection through multi-model approaches accelerate fault detection by identifying deviations in real-time and ensuring low false positive rates [194]. By minimizing manual intervention, automated test suites enable real-time fault identification and significantly

accelerate the detection process [140]. Moreover, combining diverse static analysis tools has significantly improved detection speed and accuracy, particularly in scenarios with limited resources or time constraints [136].

7.3 Empirical Studies

This section examines the use of empirical studies in WAT, focusing on effectiveness and efficiency metrics. It explains how these studies can enhance testing performance, optimize testing processes, and offer insights for further improvements in testing tools.

7.3.1 Application of Effectiveness and Efficiency Metrics in WAT. Coverage has been widely used to measure testing effectiveness in empirical WAT studies [62, 117]. MBT can generate high-quality test cases, with good code coverage and fault-detection capabilities [127, 132]. Compared to manual testing, MBT provides a more systematic approach to cover application requirements, especially for complex systems, and achieves better testing efficiency [184]. MBT is particularly effective at handling the complex interactions between multiple modules or interfaces within large-scale web applications. Empirical studies have shown that MBT reduces test redundancy and optimizes test case generation [78].

Automated tools like Model-Based Testing Leveraged for Web Tests (MoLeWe) [117] have demonstrated excellent testing efficiency [103]: MoLeWe can generate high-quality test cases, increasing test coverage while reducing execution time [62]. It can significantly reduce manual testing time, especially in applications with frequent version updates. MoLeWe has also demonstrated flexibility in resource allocation, optimizing system resource utilization in multi-threaded testing.

Fault Detection Density (FDD) and Fault Detection Effectiveness (FDE) are core effectiveness metrics used to evaluate the fault-detection capabilities of testing tools [74]. The MUTANDIS mutation-testing tool can generate non-equivalent mutants, significantly improving the fault-detection rate in JavaScript applications [121]. Mutation testing can covers a wide range of fault scenarios using diverse mutants: This can improve the comprehensiveness of the test cases. MUTANDIS has also been shown to enhance fault-detection efficiency while reducing resource waste caused by redundant testing paths.

Efficiency metrics are not only used for coverage and fault-detection capabilities, but also in execution time and resource utilization [6, 147]. Empirical studies have shown that tools like JMeter can optimize performance testing through distributed architectures [88, 157]. By parallelizing the execution of test cases, JMeter can significantly reduce execution time while ensuring system stability under high-load conditions. This architecture has been shown to perform well in real-world scenarios, especially in large-scale multi-user concurrent environments [69].

7.3.2 Improvements in Testing Performance through Empirical Studies. Empirical studies have confirmed the effectiveness of reinforcement learning techniques in the prioritization of test cases [37]: High-risk areas can be prioritized, thereby reducing overall testing time and improving fault-detection efficiency.

Automated tools can also improve testing efficiency. MoLeWe can automatically generate test cases and scripts, minimizing manual intervention and significantly enhancing automation [117]: This not only reduces the error rate from manual operations but also enables the rapid adaptation of test suites to accommodate new features or updates.

8 ANSWER TO RQ5: AVAILABLE TOOLS

This section addresses RQ5, identifying and cataloging the WAT tools that are currently available. This should provide an overview of the technological resources available to support this field. We

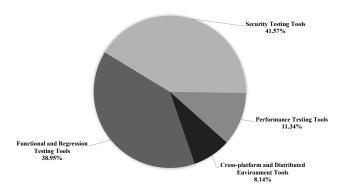


Fig. 10. Proportions of testing tools in surveyed papers.

first list the tools, and then classify them based on their types and features. The section concludes with some details on tool accessibility, including links for further exploration.

8.1 Testing Tools

Various WAT tools are available to streamline and automate the testing process, supporting efficiency, scalability, and consistency across diverse environments. These tools can support the automation of various different types of testing (including functional, regression, security, and performance testing), helping to examine the reliability and robustness of web applications.

Figure 10 presents an overview of the proportions of testing tools found in the surveyed literature. Security testing tools account for 41.57%, highlighting the increasing focus on mitigating cybersecurity risks. Functional and regression testing tools represent 38.95%, emphasizing their critical role in validating application functionality and maintaining system stability following changes. Performance testing tools comprise 11.34%, and cross-platform/distributed environment tools make up 8.14%, reflecting their specialized contributions to optimizing system performance and ensuring consistency across platforms. This distribution reveals varying priorities in testing, with security and functional testing receiving the most attention, and performance and cross-platform testing appearing to be applied only as required.

- 8.1.1 Functional and Regression Testing Tools. Functional and regression testing tools, like Selenium [142] and Katalon Studio [142], have been widely used for cross-browser and cross-platform testing: By automating test execution across various browsers and operating systems, they help to verify that web applications function correctly after updates and code changes
- 8.1.2 Security Testing Tools. Tools like Burp Suite [129] and OWASP ZAP [10] can help to identify security vulnerabilities, including SQLi, XSS, and other common web application threats. These tools are essential for protecting web applications from attacks, and ensuring compliance with security standards.
- 8.1.3 Performance Testing Tools. JMeter [55] and LoadRunner [87] have been used to simulate high-traffic loads, helping to evaluate how web applications perform under stress. They assess critical aspects of performance, including response times and throughput, evaluating whether or not the application remains stable under heavy load conditions.
- 8.1.4 Cross-platform and Distributed Environment Tools. Tools such as Docker [114] and Selenium Grid [89] enable test cases to be executed in distributed environments and containerized platforms.

They help to ensure consistency in execution and provide scalability for larger systems, supporting the efficient management of complex testing scenarios.

8.2 Tool Accessibility

The ease of access to WAT tools is a crucial factor in their adoption, and eventual effectiveness in the software development and testing community. Table 4 lists some popular WAT tools¹, categorized by type, core functionalities, and access links: This structured presentation facilitates the selection of tools based on specific testing requirements, ranging from performance evaluation and security vulnerability detection to automated UI testing and code-quality analysis. Due to page limit constraints, we provide a comprehensive summary of nearly all relevant tools, available in the online supplementary material [101].

Table 4. A comprehensive overview of tools.

Tool Name	Tool Type	Overview	Resource Link	Reference
Selenium	Open-source	It is open source, cross-platform, and supports multiple languages and browsers.	https://www.selenium. dev/ https://github.com/ SeleniumHQ	[1]
ZAP	Open-source	It is an open-source web application scanner that excels at finding common web application vulnerabilities.	https://www.zaproxy.org/ https://github.com/ zaproxy	[10]
JMeter	Open-source	It is an open-source tool for load testing and performance measurement, especially for web applications. It allows the simulation of a large number of users and provides detailed reports.	https://jmeter.apache.org/ https://github.com/apache/ jmeter.	[32]
Arachni	Open-source	It is a web vulnerability scanner designed to automatically detect and assess security vulnerabilities in web applications.	https://github.com/ Arachni/arachni	[155]
Crawljax	Open-source	It is a tool for automatically crawling and testing web applications, which automatically explores the GUI of a web application, creates state flow diagrams, and gener- ates test cases based on interactions.	https://github.com/ crawljax/crawljax	[172]
Cucumber	Open-source	It is a tool that automatically runs business case-driven tests, ensuring that all scenarios are tested as per the acceptance criteria.	https://github.com/ cucumber/cucumber	[152]
Nikto	Open-source	It is an open-source command line tool that focuses on quickly scanning web server and website configuration errors, known security vulnerabilities, and missing XSS protection headers.	https://cirt.net/Nikto2	[48]
Puppeteer	Open-source	It automates web interactions to perform testing.	https://github.com/ puppeteer/puppeteer	[39]
Cypress	Open-source	It is a JavaScript-based tool for automated end-to-end testing of modern web applications.	https://github.com/ cypress-io/cypress	[123]
Axiom	Open-source	It is a dynamic infrastructure framework for distributed computing and testing that automates distributed penetration testing across cloud environments.	https://github.com/pry0cc/axiom	[146]
Beautiful Soup	Open-source	It can be used to extract data from html and xml files.	https://www.crummy.com/ software/BeautifulSoup https://git.launchpad.net/ beautifulsoup	[140]
Commix	Open-source	It is a tool that automatically detects and exploits com- mand injection vulnerabilities in web applications, in- cluding classic command injection, blind command in- jection, and other types of tools.	https://github.com/ commixproject/commix	[168]
Katalon Studio	Commercial (Free)	It is easy to install and use, has powerful built-in reporting capabilities, and supports record and playback.	https://katalon.com/	[133]
Cytestion	Open-source	It is a web-based tool for automated exploration and testing of GUIs.	https://gitlab.com/lsi- ufcg/cytestion/cytestion	[124]

 $^{^{1}}$ To measure the popularity of each WAT tool, we collected the number of times to be mentioned in the selected 314 primary studies.

Table 4. (Continued).

Tool Name	Tool Type	Overview	Resource Link	Reference
Lazyrecon	Open-source	It is an automated vulnerability scanning and data collection framework.	https://github.com/ nahamsec/lazyrecon/ blob/master/lazyrecon.sh	[164]
Jena API	Open-source	It is an open-source Java framework for building semantic web and linked data applications, supporting SPARQL queries.	https://jena.apache.org/	[51]
Binwalk	Open-source	It is a firmware analysis tool suitable for most firmware unpacking.	https://github.com/ ReFirmLabs/binwalk	[99]
jÄk	Open-source	It is a web application crawler and scanner tool that uses dynamic Java Script code analysis.	https://github.com/ ConstantinT/jAEk	[143]
Skipfish	Open-source	It is an efficient web application security vulnerability scanning tool that can quickly detect common web application security issues.	https://code.google.com/ archive/p/skipfish/	[175]
Sqlmap	Open-source	It is a tool for detecting and exploiting SQL injection vulnerabilities.	https://sqlmap.org/	[146]
Scrapy	Open-source	It is a tool for quickly crawling HTML-based static web page content.	https://github.com/scrapy/ scrapy	[153]
Metasploit	Open-source	It is an open-source framework for penetration testing and vulnerability exploitation.	https://www.metasploit. com/ https://github.com/rapid7/ metasploit-framework	[82]
OWASP AppSensor	Open-source	It is a tool that detects suspicious application layer behavior and triggers events based on predefined sensors.	https://owasp.org/www- project-appsensor/	[163]
HTMLCS	Open-source	It is a tool for checking whether HTML code complies with WCAG standards.	http://squizlabs.github.io/ HTML_CodeSniffer/	[35]
GPTFuzzer	Open-source	It is a tool for generating web application firewall (WAF) bypass payloads, testing and discovering common web attacks such as SQL injection, cross-site scripting (XSS), and remote command execution (RCE) vulnerabilities.	https://github.com/ hongliangliang/gptfuzzer	[105]
SideeX	Open-source	It is a tool that provides cross-browser Web UI automated testing, supports recording and playback functions, and optimizes the test process through an automatic waiting mechanism.	http://sideex.org	[94]
CICFlowMeter	Open-source	It is a tool that provides network flow signatures for intrusion detection and malicious behavior classification.	https://github.com/ ahlashkari/CICFlowMeter	[151]
OWASP AppSensor	Open-source	It is a tool that detects suspicious application layer behavior and triggers events based on predefined sensors.	https://owasp.org/www- project-appsensor/	[163]
WebPageTest	Commercial (Free)	It is used to analyze web page performance, providing details such as page load time, resource load time, etc.	https://www.webpagetest. org/	[7]
WebTest	Open-source	It is a tool that supports visualization of output coverage, helps testers understand covered and uncovered output parts, and optimizes test cases.	https://github.com/ git1997/VarAnalysis	[132]
GTmetrix	Commercial (Free)	It is used to measure web page loading speed and performance and provide detailed optimization suggestions.	https://gtmetrix.com/	[7]
SikuliX	Open-source	It is an image-based testing tool that can control user input methods to interact with GUI elements, suitable for both web and non-web applications.	http://sikulix.com	[110]
Hakrawler	Open-source	It is a web crawler tool that collects web links and identifies injection points.	https://github.com/ hakluke/hakrawler	[146]
Axe Studio	Commercial (Free)	It is a tool that focuses on automated detection of web accessibility issues.	https://www.deque.com/ axe/	[35]
Cornipickle	Open-source	It is a declarative testing tool that allows developers to define layout errors using logical expressions and compare states of different snapshots.	https://github.com/liflab/ cornipickle	[27]

Open-source tools like Selenium [42], JMeter [32], and ZAP [175] provide cross-platform compatibility, support multiple programming languages, and have received strong community support. They benefit from scalability, flexibility, and collaborative development through platforms like GitHub, where contributors can provide updates and improvements.

In contrast, commercial tools such as Katalon Studio [133] and GTmetrix [7] often prioritize enhanced usability and integrated reporting features, making them more suitable for testers and organizations that value refined interfaces and out-of-the-box functionality. However, limitations in the free versions may make them less suitable for certain users.

Nearly all of the reviewed tools are easily accessible through official websites or open-source platforms. This enables a smooth integration into a wide range of testing environments. Comprehensive documentation and active community support further simplify adoption and customization. The tools offer a range of diverse capabilities: From identifying security vulnerabilities with tools like Wapiti [14] and Skipfish [175], to enhancing page performance with GTmetrix [7] and Web-PageTest [7]. Their wide availability and applicability underscore their role in advancing WAT practices.

9 ANSWER TO RQ6: CHALLENGES AND FUTURE WORK

This section addresses RQ6, identifying unresolved challenges in WAT, and suggesting directions for future research. We start by outlining the key challenges still affecting the field. We then discuss potential areas for future work that could help address these issues and improve WAT practices.

9.1 Challenges

Despite notable advances over recent years, a number of significant challenges to WAT remain that prevent the full realization of its potential. Here we discuss several key challenges:

- As web applications increase in complexity, test automation faces significant challenges, particularly when scaling up to handle large test suites efficiently.
 - The recent use of machine learning and containerization to manage these large-scale testing processes has helped. However, limitations persist for dynamic content and asynchronous operations [143]. Additionally, tool fragmentation across platforms continues to affect scalability [142].
- Balancing test coverage and execution efficiency remain challenges in web testing. Virtual
 DOM coverage techniques have demonstrated the potential to enhance testing efficiency for
 dynamic web applications by optimizing coverage metrics. However, achieving a balance
 between efficiency and comprehensive coverage remains an open challenge [199]. Hyperheuristic algorithms can reduce test-case redundancy, but more research into their application
 is required [22].
- Failure diagnosis can be a bottleneck for WAT. Model-based approaches can help, but false positives and false negatives remain significant issues [8]. Furthermore, automated repair strategies may have difficulties with the dynamic content in complex web applications, indicating a need for more robust tools.
- Maintaining test suites for dynamic web applications can also present significant challenges: Frequent updates to UI and DOM structures often break test cases [156]. Although recent work has used machine learning to adapt test scripts to changes, this remains limited when dealing with significant UI overhauls.
- Web testing tools remain fragmented, leading to inefficiencies in end-to-end testing processes. Tools like Selenium and Katalon may excel in specific areas, but they also struggle to provide cohesive solutions for more comprehensive testing needs. Recent studies have emphasized the importance of integrating these tools for more seamless cross-platform testing [114]. However, ensuring that these tools can work together across different testing environments remains a challenge [142].
- The lack of standardized metrics for evaluating testing effectiveness remains a significant challenge. Although various tools provide performance and security-testing capabilities, they often lack consistent benchmarks, making comparisons difficult. There is a need for unified evaluation frameworks that provide reliable metrics for test coverage and efficiency [166].

Initiatives like WebEV are helping to standardize some aspects of test-behavior analysis, but further research is necessary [58].

9.2 Future Work

This section outlines some potential future WAT research directions, including enhancing scalability, the integration of Large Language Models (LLMs), enhanced failure diagnosis, and the development of standardized evaluation metrics.

- Future research should enhance the scalability of test-automation frameworks to accommodate the growing complexity of modern web applications. Machine learning models should be further optimized, reducing redundancy while maintaining coverage. The use of containerization technologies, like Docker and Kubernetes, should be expanded to better enable parallel execution and to improve resource efficiency. Integrating tools into unified frameworks will also enhance automation processes and reduce fragmentation [142].
- LLMs represent a promising approach for automating test-case generation, particularly for web forms where contextual information is key. Initial studies [190] have shown the potential of models like GPT-4 [159] and Baichuan2 [189]. However, further research is required to address security concerns and refine the models for real-world use, including for dynamic content [100].
- Failure-diagnosis tools are critical to the precision and reliability of WAT. Many current systems have difficulties with false positives and false negatives: Future research should therefore explore adaptive systems that can better interpret the complex interactions in web applications. Combining functional and cross-browser testing tools will also enhance overall efficiency [8, 142].
- Standardized metrics for WAT-tool effectiveness are also needed. The absence of universal evaluation frameworks for coverage, efficiency, and performance limits meaningful comparisons: Consistent metrics will enable more accurate assessments and foster innovation in WAT methodologies [58, 166].
- Multi-agent reinforcement learning (MARL) has also shown potential for enhancing test
 coverage and efficiency in complex web applications [54]. Initial studies have demonstrated
 MARL's effectiveness for single-agent approaches. However, further research is needed to
 refine inter-agent communication and data-sharing protocols, reducing redundancy and
 optimizing exploration across web states.
- Large vision-language models (LVLMs), such as VETL [182], can support context-aware GUI
 testing by generating relevant inputs and selecting UI elements based on visual context.
 Although LVLMs have shown potential in managing dynamic, visually complex interfaces,
 future work should explore hybrid approaches that combine LVLMs with heuristic methods
 to improve cost-efficiency and adaptability, especially in applications with frequent content
 updates.

10 CONCLUSIONS

This review paper has offered a detailed analysis of advances in WAT over the past decade. The review has examined a collection of computer science papers published between January 1, 2014, and December 31, 2023. A rigorous review-research methodology was employed: Five major academic databases were systematically searched, using targeted keywords, with each retrieved paper carefully examined to ensure its relevance to the study.

Our review was guided by eight research questions, which helped organize and synthesize the data. These questions addressed essential aspects of WAT, including test-case generation and

execution, evaluation metrics, and available tools. The paper provides a clear understanding of the key trends and challenges in WAT over the past decade.

Various test-case generation methods were reviewed, including those designed to handle the dynamic and asynchronous nature of modern web applications. Regarding test execution, the importance of testing across different environments to reflect real-world conditions was emphasized. Besides, metrics used to evaluate WAT effectiveness, such as test coverage and error-detection rates, were also explored, highlighting the testing process's efficiency and quality. Additionally, available WAT tools were categorized based on their functions and accessibility, providing a practical overview for researchers and practitioners.

This paper represents a useful resource for researchers and practitioners aiming to enhance the quality, reliability, and security of web applications.

REFERENCES

- [1] R Aditya, AJ Advaith, G Sabarinath, Rahul S Rajeev, Sreya Sunil Kurup, and T Anjali. 2023. An ensemble approach for visual testing of web applications. In *Proceedings of the 3rd International Conference on Emerging Frontiers in Electrical and Electronic Technologies (ICEFEET'23)*. 1–5.
- [2] Ole Agesen. 1995. The cartesian product algorithm: Simple and precise type inference of parametric polymorphism. In *Proceedings of the 9th European Conference on Object-Oriented Programming (ECOOP'95)*, 2–26.
- [3] Pelin Akpinar, Mehmet S Aktas, Alper Bugra Keles, Yunus Balaman, Zeynep Ozdemir Guler, and Oya Kalipsiz. 2020. Web application testing with model based testing method: Case study. In *Proceedings of the 2020 International Conference on Electrical, Communication, and Computer Engineering (ICECCE'20)*. 1–6.
- [4] Rim Akrout, Eric Alata, Mohamed Kaaniche, and Vincent Nicomette. 2014. An automated black box approach for web vulnerability identification and attack scenario generation. Journal of the Brazilian Computer Society 20 (2014), 1–16.
- [5] Saranya Alagarsamy, Chakkrit Tantithamthavorn, and Aldeida Aleti. 2024. A3test: Assertion-augmented automated test case generation. *Information and Software Technology* 176 (2024), 107565.
- [6] Amira Ali and Nagwa Badr. 2015. Performance testing as a service for web applications. In *Proceedings of the IEEE 7th International Conference on Intelligent Computing and Information Systems (ICICIS'15).* 356–361.
- [7] Malik Muhammad Ali-Shahid and Shahida Sulaiman. 2015. A case study on reliability and usability testing of a web portal. In *Proceedings of the 9th Malaysian Software Engineering Conference (MySEC'15)*. 31–36.
- [8] Hamza Alkofahi, David Umphress, and Heba Alawneh. 2022. Discovering authorization business rules toward detecting web applications logic flaws. In Proceedings of the 2022 International Arab Conference on Information Technology (ACIT'22). 1–7.
- [9] Shayma Ahmed Altayaran and Wael Elmedany. 2021. Integrating web application security penetration testing into the software development life cycle: A systematic literature review. In Proceedings of the 2021 International Conference on Data Analytics for Business and Industry (ICDABI'21). 671–676.
- [10] Richard Amankwah, Jinfu Chen, Patrick Kwaku Kudjo, Beatrice Korkor Agyemang, and Alfred Adutwum Amponsah. 2020. An automated framework for evaluating open-source web scanner vulnerability severity. Service Oriented Computing and Applications 14 (2020), 297–307.
- [11] Saule Amanzholova, Darya Akhmetova, and Azhar Sagymbekova. 2021. Development of a web-resources testing system for compliance with GDPR regulation. In *Proceedings of the 7th International Conference on Engineering and MIS (ICEMIS'21)*). 1–6.
- [12] Paul Ammann and Jeff Offutt. 2017. Introduction to software testing. Cambridge University Press.
- [13] Mohammadhossein Amouei, Mohsen Rezvani, and Mansoor Fateh. 2021. RAT: Reinforcement-learning-driven and adaptive testing for vulnerability discovery in web application firewalls. *IEEE Transactions on Dependable and Secure Computing* 19, 5 (2021), 3371–3386.
- [14] Karthik Anagandula and Pavol Zavarsky. 2020. An analysis of effectiveness of black-box web application scanners in detection of stored SQL injection and stored XSS vulnerabilities. In *Proceedings of the 3rd International Conference on Data Intelligence and Security (ICDIS'20)*. 40–48.
- [15] Saswat Anand, Edmund K Burke, Tsong Yueh Chen, John Clark, Myra B Cohen, Wolfgang Grieskamp, Mark Harman, Mary Jean Harrold, Phil McMinn, Antonia Bertolino, et al. 2013. An orchestrated survey of methodologies for automated software test case generation. *Journal of systems and software* 86, 8 (2013), 1978–2001.
- [16] Dennis Appelt, Cu D Nguyen, Annibale Panichella, and Lionel C Briand. 2018. A machine-learning-driven evolutionary approach for testing web application firewalls. *IEEE Transactions on Reliability* 67, 3 (2018), 733–757.
- [17] Jinat Ara, Cecilia Sik-Lanyi, and Arpad Kelemen. 2024. Accessibility engineering in web evaluation process: A systematic literature review. *Universal Access in the Information Society* 23, 2 (2024), 653–686.

- [18] Andrea Arcuri, Man Zhang, Amid Golmohammadi, Asma Belhadi, Juan P Galeotti, Bogdan Marculescu, and Susruthan Seran. 2023. EMB: A curated corpus of web/enterprise applications and library support for software testing research. In Proceedings of the 16th IEEE Conference on Software Testing, Verification and Validation (ICST'23). 433–442.
- [19] Anuja Arora and Madhavi Sinha. 2013. Dynamic content testing of web application using user session based state testing. In *Proceedings of the 4th International Conference-Confluence The Next Generation Information Technology Summit (Confluence'13)*. 22–28.
- [20] Nor Fatimah Awang, Ahmad Dahari Jarno, Syahaneim Marzuki, Nor Azliana Akmal Jamaludin, Khairani Abd Majid, and Taniza Tajuddin. 2019. Method for generating test data for detecting SQL injection vulnerability in web application. In Proceedings of the 7th International Conference on Cyber and IT Service Management (CITSM'19), Vol. 7. 1–5.
- [21] Murat Aydos, Çiğdem Aldan, Evren Coşkun, and Alperen Soydan. 2022. Security testing of web applications: A systematic mapping of the literature. Journal of King Saud University-Computer and Information Sciences 34, 9 (2022), 6775–6792.
- [22] Juliana Marino Balera and Valdivino Alexandre de Santiago Júnior. 2022. Multiperspective web testing supported by a generation hyper-heuristic. In *Proceedings of the 22nd International Conference on Computational Science and Its Applications (ICCSA'22)*. 447–462.
- [23] Sebastian Balsam and Deepti Mishra. 2024. Web application testing—Challenges and opportunities. *Journal of Systems and Software* (2024), 112186.
- [24] Giuseppe Beltrano, Claudia Greco, Michele Ianni, and Giancarlo Fortino. 2023. Deep learning-based detection of CSRF vulnerabilities in web applications. In Proceedings of the 2023 IEEE International Conference on Dependable, Autonomic and Secure Computing, International Conference on Pervasive Intelligence and Computing, International Conference on Cloud and Big Data Computing, International Conference on Cyber Science and Technology Congress (DASC/PiCom/CBDCom/CyberSciTech'23). 0916–0920.
- [25] Wafa Ben Jaballah and Nizar Kheir. 2016. A grey-box approach for detecting malicious user interactions in web applications. In Proceedings of the 8th ACM CCS International Workshop on Managing Insider Security Threats (MIST'16). 1–12.
- [26] T.J. Berners-Lee. 1992. The world-wide web. Computer networks and ISDN systems 25, 4 (1992), 454-459.
- [27] Oussama Beroual, Francis Guérin, and Sylvain Hallé. 2020. Detecting responsive web design bugs with declarative specifications. In *Proceedings of the 20th International Conference on Web Engineering (ICWE'20)*. 3–18.
- [28] Thilini Bhagya, Jens Dietrich, and Hans Guesgen. 2019. Generating mock skeletons for lightweight web-service testing. In *Proceedings of the 26th Asia-Pacific Software Engineering Conference (APSEC'19)*. 181–188.
- [29] SM Bindu Bhargavi and V Suma. 2022. A survey of the software test methods and identification of critical success factors for automation. SN Computer Science 3, 6 (2022), 449.
- [30] Josip Bozic, Bernhard Garn, Ioannis Kapsalis, Dimitris Simos, Severin Winkler, and Franz Wotawa. 2015. Attack pattern-based combinatorial testing with constraints for web security testing. In Proceedings of the 2015 IEEE International Conference on Software Quality, Reliability and Security (QRS'15). 207–212.
- [31] Josip Bozic, Bernhard Garn, Dimitris E Simos, and Franz Wotawa. 2015. Evaluation of the IPO-family algorithms for test case generation in web security testing. In *Proceedings of the IEEE 8th International Conference on Software Testing, Verification and Validation Workshops (ICSTW'15)*. 1–10.
- [32] Jianhua Cai and Qingchun Hu. 2014. Analysis for cloud testing of web application. In *Proceedings of the 2nd International Conference on Systems and Informatics (ICSAI'14)*. 293–297.
- [33] Stefano Calzavara, Mauro Conti, Riccardo Focardi, Alvise Rabitti, and Gabriele Tolomei. 2020. Machine learning for web vulnerability detection: The case of cross-site request forgery. IEEE Security and Privacy 18, 3 (2020), 8–16.
- [34] Stefano Calzavara, Hugo Jonker, Benjamin Krumnow, and Alvise Rabitti. 2021. Measuring web session security at scale. *Computers & Security* 111 (2021), 102472.
- [35] Milton Campoverde-Molina, Sergio Luján-Mora, and Llorenç Valverde. 2021. Process model for continuous testing of web accessibility. IEEE Access 9 (2021), 139576–139593.
- [36] Nazanin Bayati Chaleshtari, Fabrizio Pastore, Arda Goknil, and Lionel C Briand. 2023. Metamorphic testing for web system security. *IEEE Transactions on Software Engineering* 49, 6 (2023), 3430–3471.
- [37] Xiaoning Chang, Zheheng Liang, Yifei Zhang, Lei Cui, Zhenyue Long, Guoquan Wu, Yu Gao, Wei Chen, Jun Wei, and Tao Huang. 2023. A reinforcement learning approach to generating test cases for web applications. In Proceedings of the 2023 IEEE/ACM International Conference on Automation of Software Test (AST'23). 13–23.
- [38] Tsong Yueh Chen, Fei-Ching Kuo, Huai Liu, Pak-Lok Poon, Dave Towey, TH Tse, and Zhi Quan Zhou. 2018. Metamorphic testing: A review of challenges and opportunities. Comput. Surveys 51, 1 (2018), 4–1.
- [39] Wei Chen, Hanyang Cao, and Xavier Blanc. 2021. An improving approach for DOM-based web test suite repair. In *Proceedings of the 21st International Conference on Web Engineering (ICWE'21)*. 372–387.
- [40] Anna Corazza, Sergio Di Martino, Adriano Peron, and Luigi Libero Lucio Starace. 2021. Web application testing: Using tree kernels to detect near-duplicate states in automated model inference. In *Proceedings of the 15th ACM/IEEE*

- International Symposium on Empirical Software Engineering and Measurement (ESEM'21). 1-6.
- [41] Yanpeng Cui, Junjie Cui, and Jianwei Hu. 2020. A survey on XSS attack detection and prevention in web applications. In *Proceedings of the 12th International Conference on Machine Learning and Computing (ICMLC'20)*. 443–449.
- [42] Xiao Dawei, Jiang Liqiu, Xu Xinpeng, and Wang Yuhang. 2016. Web application automatic testing solution. In Proceedings of the 3rd International Conference on Information Science and Control Engineering (ICISCE'16). 1183–1187.
- [43] Romulo de Almeida Neves, Willian Massami Watanabe, and Rafael Oliveira. 2022. Morpheus web testing: A tool for generating test cases for widget based web applications. *Journal of Web Engineering* 21, 2 (2022), 119–144.
- [44] Flávio Rezende de Jesus, Leandro Guarino de Vasconcelos, and Laércio A Baldochi. 2015. Leveraging task-based data to support functional testing of web applications. In Proceedings of the 30th Annual ACM Symposium on Applied Computing (SAC'15). 783–790.
- [45] Jessica Lasch de Moura, Andrea Schwertner Charao, João Carlos Damasceno Lima, and Benhur de Oliveira Stein. 2017. Test case generation from BPMN models for automated testing of Web-based BPM applications. In Proceedings of the 17th International Conference on Computational Science and Its Applications (ICCSA'17). 1–7.
- [46] Vidroha Debroy, Lance Brimble, Matthew Yost, and Archana Erry. 2018. Automating web application testing from the ground up: Experiences and lessons learned in an industrial setting. In *Proceedings of the IEEE 11th International Conference on Software Testing, Verification and Validation (ICST'18)*. 354–362.
- [47] Dilek Yilmazer Demirel and Mehmet Tahir Sandikkaya. 2023. ACUM: An approach to combining unsupervised methods for detecting malicious web sessions. In Proceedings of the 8th International Conference on Computer Science and Engineering (UBMK'23). 288–293.
- [48] R Sri Devi and M Mohan Kumar. 2020. Testing for security weakness of web applications using ethical hacking. In 2020 4th International Conference on Trends in Electronics and Informatics (ICOEI)(48184). IEEE, 354–361.
- [49] Arilo C Dias Neto, Rajesh Subramanyan, Marlon Vieira, and Guilherme H Travassos. 2007. A survey on model-based testing approaches: A systematic review. In Proceedings of the 1st ACM international workshop on Empirical assessment of software engineering languages and technologies: held in conjunction with the 22nd IEEE/ACM International Conference on Automated Software Engineering (ASE'07). 31–36.
- [50] Serdar Doğan, Aysu Betin-Can, and Vahid Garousi. 2014. Web application testing: A systematic literature review. *Journal of Systems and Software* 91 (2014), 174–201.
- [51] K Naveen Durai, R Subha, and Anandakumar Haldorai. 2021. A novel method to detect and prevent SQLIA using ontology to cloud web security. Wireless Personal Communications 117, 4 (2021), 2995–3014.
- [52] Sebastian Elbaum, Hui Nee Chin, Matthew B Dwyer, and Matthew Jorde. 2008. Carving and replaying differential unit test cases from system test cases. *IEEE Transactions on Software Engineering* 35, 1 (2008), 29–45.
- [53] Gerald D Everett and Raymond McLeod Jr. 2007. Software testing: Testing across the entire software development life cycle. John Wiley and Sons.
- [54] Yujia Fan, Sinan Wang, Zebang Fei, Yao Qin, Huaxuan Li, and Yepang Liu. 2024. Can cooperative multi-agent reinforcement learning boost automatic web testing? An exploratory study. In Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering (ASE'24). 14–26.
- [55] Yujia Fan, Siyi Wang, Sinan Wang, Yepang Liu, Guoyao Wen, and Qi Rong. 2023. A Comprehensive Evaluation of Q-Learning Based Automatic Web GUI Testing. In Proceedings of the 10th International Conference on Dependable Systems and Their Applications (DSA'23). 12–23.
- [56] Michael Felderer, Matthias Büchler, Martin Johns, Achim D Brucker, Ruth Breu, and Alexander Pretschner. 2016. Security testing: A survey. In Advances in Computers. Vol. 101. 1–51.
- [57] Ana Fidalgo, Ibéria Medeiros, Paulo Antunes, and Nuno Neves. 2020. Towards a deep learning model for vulnerability detection on web application variants. In Proceedings of the 13th IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW'20). 465–476.
- [58] Mridha Md Nafis Fuad and Kazi Sakib. 2023. WebEV: A dataset on the behavior of testers for web application end to end testing. In Proceedings of the IEEE/ACM 31st International Conference on Program Comprehension (ICPC'23), 79–83.
- [59] Tommaso Fulcini and Luca Ardito. 2022. Gamified exploratory GUI testing of web applications: A preliminary evaluation. In Proceedings of the 2022 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW'22). 215–222.
- [60] Disha Garg, Abhishek Singhal, and Abhay Bansal. 2015. A framework for testing web applications using action word based testing. In Proceedings of the 1st International Conference on Next Generation Computing Technologies (NGCT'15). 593–598.
- [61] Bernhard Garn, Ioannis Kapsalis, Dimitris E Simos, and Severin Winkler. 2014. On the applicability of combinatorial testing to web application security testing: A case study. In *Proceedings of the 2014 Workshop on Joining AcadeMiA and Industry Contributions to Test Automation and Model-Based Testing (JAMAICA'14)*. 16–21.
- [62] Vahid Garousi, Alper Buğra Keleş, Yunus Balaman, Zeynep Özdemir Güler, and Andrea Arcuri. 2021. Model-based testing in practice: An experience report from the web applications domain. Journal of Systems and Software 180

- (2021), 111032.
- [63] Vahid Garousi, Ali Mesbah, Aysu Betin-Can, and Shabnam Mirshokraie. 2013. A systematic mapping study of web application testing. *Information and Software Technology* 55, 8 (2013), 1374–1396.
- [64] Jesse James Garrett et al. 2005. Ajax: A new approach to web applications. Technical Report (2005).
- [65] Bronjon Gogoi, Tasiruddin Ahmed, and Ratnaboli Ghorai Dinda. 2022. PHP web shell detection through static analysis of AST using LSTM based deep learning. In Proceedings of the 1st International Conference on Artificial Intelligence Trends and Pattern Recognition (ICAITPR'22). 1–6.
- [66] Amid Golmohammadi, Man Zhang, and Andrea Arcuri. 2023. Testing restful APIs: A survey. ACM Transactions on Software Engineering and Methodology 33, 1 (2023), 1–41.
- [67] Maciej Grzenda, Stanisław Kaźmierczak, Marcin Luckner, Grzegorz Borowik, and Jacek Mańdziuk. 2023. Evaluation of machine learning methods for impostor detection in web applications. Expert Systems with Applications 231 (2023), 120736.
- [68] Zhibin Guan, Jiajie Wang, Xiaomeng Wang, Wei Xin, Jing Cui, and Xiangping Jing. 2021. A comparative study of RNN-based methods for web malicious code detection. In Proceedings of the IEEE 6th International Conference on Computer and Communication Systems (ICCCS'21). 769–773.
- [69] Marco Guarnieri, Petar Tsankov, Tristan Buchs, Mohammad Torabi Dashti, and David Basin. 2017. Test execution checkpointing for web applications. In Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'17). 203–214.
- [70] Mukesh Kumar Gupta, M.C. Govil, and Girdhari Singh. 2014. Static analysis approaches to detect SQL injection and cross site scripting vulnerabilities in web applications: A survey. In Proceedings of the 1st International Conference on Recent Advances and Innovations in Engineering (ICRAIE'14). 1–5.
- [71] Mukesh Kumar Gupta, Mahesh Chandra Govil, and Girdhari Singh. 2015. Text-mining based predictive model to detect XSS vulnerable files in web applications. In *Proceedings of the 2015 Annual IEEE India Conference (INDICON'15)*. 1–6.
- [72] Nishant Gupta, Vibhash Yadav, and Mayank Singh. 2018. Automated regression test case generation for web application: A survey. *Comput. Surveys* 51, 4 (2018), 1–25.
- [73] Rashmi Gupta and Neha Bajpai. 2014. A keyword-driven tool for testing Web applications (KeyDriver). *IEEE Potentials* 33, 5 (2014), 35–42.
- [74] Elahe Habibi and Seyed-Hassan Mirian-Hosseinabadi. 2015. Event-driven web application testing based on model-based mutation testing. *Information and Software Technology* 67 (2015), 159–179.
- [75] Axel Halin, Alexandre Nuttinck, Mathieu Acher, Xavier Devroey, Gilles Perrouin, and Benoit Baudry. 2019. Test them all, is it worth it? Assessing configuration sampling on the JHipster web development stack. *Empirical Software Engineering* 24 (2019), 674–717.
- [76] Samer Hanna and Amro Al-Said Ahmad. 2022. Web applications testing techniques: A systematic mapping study. International Journal of Web Engineering and Technology 17, 4 (2022), 372–412.
- [77] Samer Hanna and Hayat Jaber. 2019. An approach for web applications test data generation based on analyzing client side user input fields. In *Proceedings of the 2nd International Conference on New Trends in Computing Sciences (ICTCS'19)*. 1–6.
- [78] Samer Hanna and Malcolm Munro. 2018. Test case generation for semantic-based user input validation of web applications. *International Journal of Web Engineering and Technology* 13, 3 (2018), 225–254.
- [79] Mark Harman, Yue Jia, and Yuanyuan Zhang. 2015. Achievements, open problems and challenges for search based software testing. In *Proceedings of the IEEE 8th International Conference on Software Testing, Verification and Validation (ICST'15)*. 1–12.
- [80] Rubing Huang, Weifeng Sun, Yinyin Xu, Haibo Chen, Dave Towey, and Xin Xia. 2021. A survey on adaptive random testing. *IEEE Transactions on Software Engineering* 47, 10 (2021), 2052–2083.
- [81] CP Indumathi and A Sabeena Begum. 2016. Web database testing using ER diagram and state transition model. In Proceedings of the 2016 International Conference on Communication and Signal Processing (ICCSP'16). 1937–1942.
- [82] Indranil Jana and Alina Oprea. 2019. AppMine: Behavioral analytics for web application vulnerability detection. In *Proceedings of the 2019 ACM SIGSAC Conference on Cloud Computing Security Workshop (CCSW'19)*. 69–80.
- [83] Beakcheol Jang, Myeonghwi Kim, Gaspard Harerimana, and Jong Wook Kim. 2019. Q-learning algorithms: A comprehensive classification and applications. IEEE access 7 (2019), 133653–133667.
- [84] Mehdi Jazayeri. 2007. Some trends in web application development. In *Proceedings of the Future of Software Engineering* (FOSE'07). 199–213.
- [85] Yue Jia and Mark Harman. 2010. An analysis and survey of the development of mutation testing. *IEEE Transactions on Software Engineering* 37, 5 (2010), 649–678.
- [86] Ali Karimoddini, Mubbashar Altaf Khan, Solomon Gebreyohannes, Mike Heiges, Ethan Trewhitt, and Abdollah Homaifar. 2022. Automatic test and evaluation of autonomous systems. *IEEE Access* 10 (2022), 72227–72238.

- [87] Rijwan Khan and Mohd Amjad. 2016. Smoke testing of web application based on ALM tool. In Proceedings of the 2016 International Conference on Computing, Communication and Automation (ICCCA'16). 862–866.
- [88] Sandhya Kiran, Akshyansu Mohapatra, and Rajashekara Swamy. 2015. Experiences in performance testing of web applications with unified authentication platform using Jmeter. In Proceedings of the 2015 International Symposium on Technology Management and Emerging Technologies (ISTMET'15). 74–78.
- [89] Vinodkumar H Kiranagi and Gopal Krishna Shyam. 2017. Feature driven hybrid test automation framework (FDHTAF) for web based or cloud based application testing. In *Proceedings of the 2017 International Conference On Smart Technologies For Smart Nation (SmartTechCon'17)*. 1555–1559.
- [90] Hiroyuki Kirinuki, Shinsuke Matsumoto, Yoshiki Higo, and Shinji Kusumoto. 2021. NLP-assisted web element identification toward script-free testing. In Proceedings of the 37th IEEE International Conference on Software Maintenance and Evolution (ICSME'21). 639–643.
- [91] Ambreen Kousar, Saif Ur Rehman Khan, Shahid Hussain, M Abdul Basit Ur Rahim, Wen-Li Wang, and Naseem Ibrahim. 2023. A systematic review on pattern-based GUI testing of android and web apps: State-of-the-art, taxonomy, challenges and future directions. In *Proceedings of the 25th International Multitopic Conference (INMIC'23)*. 1–7.
- [92] Deepak Kumar, Zane Ma, Zakir Durumeric, Ariana Mirian, Joshua Mason, J Alex Halderman, and Michael Bailey. 2017. Security challenges in an increasingly tangled web. In Proceedings of the 26th International Conference on World Wide Web (WWW'17). 677–684.
- [93] Pratap Kumar and Ravi K Sheth. 2016. A review on 0-day vulnerability testing in web application. In Proceedings of the 2nd International Conference on Information and Communication Technology for Competitive Strategies (ICTCS'16).
- [94] Shin-Jie Lee, Yu-Xian Chen, Shang-Pin Ma, and Wen-Tin Lee. 2018. Test command auto-wait mechanisms for record and playback-style web application testing. In *Proceedings of the IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC'18)*, Vol. 2. 75–80.
- [95] Maurizio Leotta, Matteo Biagiola, Filippo Ricca, Mariano Ceccato, and Paolo Tonella. 2020. A family of experiments to assess the impact of page object pattern in web test suite development. In *Proceedings of the IEEE 13th International Conference on Software Testing, Validation and Verification (ICST'20)*. 263–273.
- [96] Maurizio Leotta, Davide Paparella, and Filippo Ricca. 2024. Mutta: A novel tool for E2E web mutation testing. *Software Quality Journal* 32, 1 (2024), 5–26.
- [97] Maurizio Leotta, Andrea Stocco, Filippo Ricca, and Paolo Tonella. 2015. Meta-heuristic generation of robust XPath locators for web testing. In *Proceedings of the IEEE/ACM 8th International Workshop on Search-Based Software Testing (SBST'15)*. 36–39.
- [98] Guodong Li, Esben Andreasen, and Indradeep Ghosh. 2014. SymJS: Automatic symbolic testing of JavaScript web applications. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE'14)*. 449–459.
- [99] Lukai Li, Ruijie Cai, Yongguang Zhang, and Xiaokang Yin. 2023. Facilitating web vulnerability detection on embedded devices with root path pruning. In Proceedings of the 3rd International Conference on Computer Science, Electronic Information Engineering and Intelligent Control Technology (CEI'23). 41–46.
- [100] Tao Li, Chenhui Cui, Lei Ma, Dave Towey, Yujie Xie, and Rubing Huang. 2024. Leveraging large language models for automated web-form-test generation: An empirical study. arXiv 2405.09965 (2024).
- [101] Tao Li, Rubing Huang, Chenhui Cui, Dave Towey, and Lei Ma. 2024. More details for a survey on web application testing. https://github.com/abelli1024/wat-survey. Accessed: 2024-10-15.
- [102] Xiaowei Li, Xujie Si, and Yuan Xue. 2014. Automated black-box detection of access control vulnerabilities in web applications. In Proceedings of the 4th ACM Conference on Data and Application Security and Privacy (SP'14). 49–60.
- [103] Xinxin Li and Hongwei Zeng. 2014. Modeling web application for cross-browser compatibility testing. In Proceedings of the 15th IEEE/ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD'14). 1–5.
- [104] Yuan-Fang Li, Paramjit K Das, and David L Dowe. 2014. Two decades of Web application testing-A survey of recent advances. *Information Systems* 43 (2014), 20–54.
- [105] Hongliang Liang, Xiangyu Li, Da Xiao, Jie Liu, Yanjie Zhou, Aibo Wang, and Jin Li. 2023. Generative pre-trained transformer-based reinforcement learning for testing web application firewalls. IEEE Transactions on Dependable and Secure Computing 21, 1 (2023), 309–324.
- [106] Jun-Wei Lin, Farn Wang, and Paul Chu. 2017. Using semantic similarity in crawling-based web application testing. In Proceedings of the 10th IEEE International Conference on Software Testing, Verification and Validation (ICST'17). 138–148.
- [107] Cuauhtémoc López-Martín. 2022. Machine learning techniques for software testing effort prediction. *Software Quality Journal* 30, 1 (2022), 65–100.

- [108] Ouarda Lounis, Salah Eddine Bouhouita Guermeche, and Lalia Saoudi. 2014. A new algorithm for detecting SQL injection attack in Web application. In Proceedings of the 2014 Science and Information Conference (SAI'14). 589–594.
- [109] Romaric Ludinard, Eric Totel, Frédéric Tronel, Vincent Nicomette, Mohamed Kaâniche, Eric Alata, Rim Akrout, and Yann Bachy. 2018. An invariant-based approach for detecting attacks against data in web applications. In Application Development and Design: Concepts, Methodologies, Tools, and Applications. 1073–1094.
- [110] Federico Macchi, Pierpaolo Rosin, Juan Marcos Mervi, and Luca Turchet. 2021. Image-based approaches for automating GUI testing of interactive web-based applications. In *Proceedings of the 28th Conference of Open Innovations Association (FRUCT'21)*. 278–285.
- [111] Mostafa Mahdieh, Seyed-Hassan Mirian-Hosseinabadi, and Mohsen Mahdieh. 2022. Test case prioritization using test case diversification and fault-proneness estimations. *Automated Software Engineering* 29, 2 (2022), 50.
- [112] Abdalla Wasef Marashdih, Zarul Fitri Zaaba, Khaled Suwais, and Nur Azimah Mohd. 2019. Web application security: An investigation on static analysis with other algorithms to detect cross site scripting. *Procedia Computer Science* 161 (2019), 1173–1181.
- [113] Leonardo Mariani, Mauro Pezzè, Oliviero Riganelli, and Mauro Santoro. 2014. Link: Exploiting the web of data to generate test inputs. In Proceedings of the 23rd International Symposium on Software Testing and Analysis (ISSTA'14). 373–384.
- [114] Fabian Marquardt and Lennart Buhl. 2021. Déjà vu? Client-side fingerprinting and version detection of web application software. In *Proceedings of the IEEE 46th Conference on Local Computer Networks (LCN'21)*. 81–89.
- [115] Beatriz Martins and Carlos Duarte. 2023. A large-scale web accessibility analysis considering technology adoption. *Universal Access in the Information Society* (2023), 1–16.
- [116] Achmad Fahrurrozi Maskur and Yudistira Dwi Wardhana Asnar. 2019. Static code analysis tools with the taint analysis method for detecting web application vulnerability. In *Proceedings of the 2019 International Conference on Data and Software Engineering (ICoDSE'19)*. 1–6.
- [117] Guilherme Ricken Mattiello and André Takeshi Endo. 2022. Model-based testing leveraged for automated web tests. *Software Quality Journal* 30, 3 (2022), 621–649.
- [118] Ibéria Medeiros, Nuno Neves, and Miguel Correia. 2016. DEKANT: A static analysis tool that learns to detect web application vulnerabilities. In *Proceedings of the 25th International Symposium on Software Testing and Analysis* (ISSTA'16). 1–11.
- [119] Ali Mesbah, Arie Van Deursen, and Stefan Lenselink. 2012. Crawling Ajax-based web applications through dynamic analysis of user interface state changes. *ACM Transactions on the Web* 6, 1 (2012), 1–30.
- [120] Amin Milani Fard, Mehdi Mirzaaghaei, and Ali Mesbah. 2014. Leveraging existing tests in automated test generation for web applications. In *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering (ASE'14)*. 67–78.
- [121] Shabnam Mirshokraie, Ali Mesbah, and Karthik Pattabiraman. 2014. Guided mutation testing for JavaScript web applications. *IEEE Transactions on Software Engineering* 41, 5 (2014), 429–444.
- [122] Fawaz Mahiuob Mohammed Mokbal, Wang Dan, Azhar Imran, Lin Jiuchuan, Faheem Akhtar, and Wang Xiaoxi. 2019. MLPXSS: An integrated XSS-based attack detection scheme in web applications using multilayer perceptron technique. IEEE Access 7 (2019), 100567–100580.
- [123] Humaid Mollah and Petra van den Bos. 2023. From user stories to end-to-end web testing. In *Proceedings of the 16th IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW'23).* 140–148.
- [124] Thiago Santos de Moura, Everton LG Alves, Hugo Feitosa de Figueirêdo, and Cláudio de Souza Baptista. 2023. Cytestion: Automated GUI testing for web applications. In Proceedings of the XXXVII Brazilian Symposium on Software Engineering (SBSI'23). 388–397.
- [125] Joydeep Mukherjee, Mea Wang, and Diwakar Krishnamurthy. 2014. Performance testing web applications on the cloud. In *Proceedings of the IEEE 7th International Conference on Software Testing, Verification and Validation Workshops* (ICSTW'14). 363–369.
- [126] Muhammad Takdir Muslihi and Daniyal Alghazzawi. 2020. Detecting SQL injection on web application using deep learning techniques: A systematic literature review. In *Proceedings of the 3rd International Conference on Vocational Education and Electrical Engineering (ICVEE'20)*. 1–6.
- [127] Miguel Nabuco and Ana CR Paiva. 2014. Model-based test case generation for web applications. In *Proceedings of the* 14th International Conference on Computational Science and Its Applications (ICCSA'14). 248–262.
- [128] Leckraj Nagowah and Kreshnah Kora-Ramiah. 2017. Automated complete test case coverage for web based applications. In Proceedings of the 2017 International Conference on Infocom Technologies and Unmanned Systems (Trends and Future Directions) (ICTUS'17). 383–390.
- [129] Sangeeta Nagpure and Sonal Kurkure. 2017. Vulnerability assessment and penetration testing of web application. In Proceedings of the 2017 International Conference on Computing, Communication, Control and Automation (ICCUBEA'17). 1–6.

- [130] Andy Neumann, Nuno Laranjeiro, and Jorge Bernardino. 2018. An analysis of public REST web service APIs. IEEE Transactions on Services Computing 14, 4 (2018), 957–970.
- [131] Hanh-Phuc Nguyen, Thanh-Nhan Luong, and Ninh-Thuan Truong. 2022. Generating test paths to detect XSS vulnerabilities of web applications. In Proceedings of the 9th NAFOSTED Conference on Information and Computer Science (NICS'22). 287–293.
- [132] Hung Viet Nguyen, Hung Dang Phan, Christian Kästner, and Tien N Nguyen. 2019. Exploring output-based coverage for testing PHP web applications. *Automated Software Engineering* 26 (2019), 59–85.
- [133] Vu Nguyen, Thanh To, and Gia-Han Diep. 2021. Generating and selecting resilient and maintainable locators for web automated testing. *Software Testing, Verification and Reliability* 31, 3 (2021), e1760.
- [134] Changhai Nie and Hareton Leung. 2011. A survey of combinatorial testing. Comput. Surveys 43, 2 (2011), 1-29.
- [135] K Nirmal, B Janet, and Rajagopal Kumar. 2021. Analyzing and eliminating phishing threats in IoT, network and other Web applications using iterative intersection. *Peer-to-Peer Networking and Applications* 14 (2021), 2327–2339.
- [136] Paulo Nunes, Ibéria Medeiros, José Fonseca, Nuno Neves, Miguel Correia, and Marco Vieira. 2019. An empirical study on combining diverse static analysis tools for web security vulnerabilities based on development scenarios. *Computing* 101 (2019), 161–185.
- [137] Jeff Offutt and Sunitha Thummala. 2019. Testing concurrent user behavior of synchronous web applications with Petri nets. *Software and Systems Modeling* 18 (2019), 913–936.
- [138] Agnija Onukrane, Heinrihs Kristians Skrodelis, Galina Merkurjeva, and Andrejs Romanovs. 2023. Navigating web application security: A survey of vulnerabilities and detection solutions. In Proceedings of the IEEE 64th International Scientific Conference on Information Technology and Management Science of Riga Technical University (ITMS'23). 1–6.
- [139] Vikas Panthi and Durga Prasad Mohapatra. 2017. An approach for dynamic web application testing using MBT. *International Journal of System Assurance Engineering and Management* 8 (2017), 1704–1716.
- [140] Nicey Paul and Robin Tommy. 2018. An approach of automated testing on web based platform using machine learning and Selenium. In Proceedings of the 2018 International Conference on Inventive Research in Computing Applications (ICIRCA'18). 851–856.
- [141] Silvio Pavanetto and Marco Brambilla. 2020. Generation of realistic navigation paths for web site testing using recurrent neural networks and generative adversarial neural networks. In *Proceedings of the 20 the International Conference on Web Engineering (ICWE'20)*. 244–258.
- [142] Elis Pelivani and Betim Cico. 2021. A comparative study of automation testing tools for web applications. In *Proceedings* of the 10th Mediterranean Conference on Embedded Computing (MECO'21). 1–6.
- [143] Giancarlo Pellegrino, Constantin Tschürtz, Eric Bodden, and Christian Rossow. 2015. jäk: Using dynamic analysis to crawl and test modern web applications. In *Proceedings of the 18th International Symposium on Research in Attacks, Intrusions, and Defenses (RAID'15)*. 295–316.
- [144] I Putu Agus Eka Pratama and Alvin Maulana Rhusuli. 2022. Penetration testing on web application using insecure direct object references (IDOR) method. In *Proceedings of the 2022 International Conference on ICT for Smart Society* (ICISS'22). 01–07.
- [145] Irfan Prazina, Šeila Bećirović, Emir Cogo, and Vensada Okanović. 2023. Methods for automatic web page layout testing and analysis: A review. *IEEE Access* 11 (2023), 13948–13964.
- [146] Pratama Aji Prisadi, Setiyo Cahyono, Ryan Muhammad Azizulfiqar, and Alfido Osdie. 2023. Implementation of distributed attack penetration testing automation using dynamic infrastructure framework Axiom on web-based systems. In Proceeings of the 2023 International Conference on Informatics, Multimedia, Cyber and Informations System (ICIMCIS'23). 183–188.
- [147] Mayang Anglingsari Putri, Hilman Nuril Hadi, and Fatwa Ramdani. 2017. Performance testing analysis on web application: Study case student admission web system. In *Proceedings of the 2017 International Conference on Sustainable Information Engineering and Technology (SIET'17)*. 1–5.
- [148] Xiao-Fang Qi, Zi-Yuan Wang, Jun-Qiang Mao, and Peng Wang. 2017. Automated testing of web applications using combinatorial strategies. *Journal of Computer Science and Technology* 32, 1 (2017), 199–210.
- [149] Sajjad Rafique, Mamoona Humayun, Bushra Hamid, Ansar Abbas, Muhammad Akhtar, and Kamil Iqbal. 2015. Web application security vulnerabilities detection approaches: A systematic mapping study. In Proceedings of the IEEE/ACIS 16th International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD'15). 1–6.
- [150] Karishma Rahman and Clemente Izurieta. 2022. A mapping study of security vulnerability detection approaches for web applications. In Proceedings of the 48th Euromicro Conference on Software Engineering and Advanced Applications (SEAA'22). 491–494.
- [151] Branislav Rajić, Žarko Stanisavljević, and Pavle Vuletić. 2023. Early web application attack detection using network traffic analysis. *International Journal of Information Security* 22, 1 (2023), 77–91.

- [152] Paresh Rathod, Viljami Julkunen, Tero Kaisti, and Janne Nissilä. 2015. Automatic acceptance testing of the web application security with ITU-T X. 805 framework. In Proceedings of the 2nd International Conference on Computer Science, Computer Engineering, and Social Media (CSCESM'15). 103–108.
- [153] Marc Rennhard, Malte Kushnir, Olivier Favre, Damiano Esposito, and Valentin Zahnd. 2022. Automating the detection of access control vulnerabilities in web applications. *SN Computer Science* 3, 5 (2022), 376.
- [154] Jesús-Ángel Román-Gallego, María-Luisa Pérez-Delgado, Marcos Luengo Viñuela, and María-Concepción Vega-Hernández. 2023. Artificial intelligence web application firewall for advanced detection of web injection attacks. Expert Systems (2023), e13505.
- [155] Fernando Román Muñoz, Iván Israel Sabido Cortes, and Luis Javier García Villalba. 2018. Enlargement of vulnerable web applications for testing. *The Journal of Supercomputing* 74 (2018), 6598–6617.
- [156] Yeonhee Ryou and Sukyoung Ryu. 2018. Automatic detection of visibility faults by layout changes in HTML5 web pages. In Proceedings of the IEEE 11th International Conference on Software Testing, Verification and Validation (ICST'18). 182–192.
- [157] Divya Saharan, Yogesh Kumar, and Rahul Rishi. 2018. Analytical study and implementation of web performance testing tools. In *Proceedings of the 2018 International Conference on Recent Innovations in Electrical, Electronics and Communication Engineering (ICRIEECE'18).* 2370–2377.
- [158] Sreedevi Sampath and Sara Sprenkle. 2016. Advances in web application testing, 2010–2014. In Advances in Computers. Vol. 101. 155–191.
- [159] Katharine Sanderson. 2023. GPT-4 is here: What scientists think. Nature 615, 7954 (2023), 773-773.
- [160] Sergio Segura, Gordon Fraser, Ana B Sanchez, and Antonio Ruiz-Cortés. 2016. A survey on metamorphic testing. IEEE Transactions on Software Engineering 42, 9 (2016), 805–824.
- [161] Saad Shafiq, Atif Mashkoor, Christoph Mayr-Dorn, and Alexander Egyed. 2021. A literature review of using machine learning in software development life cycle stages. *IEEE Access* 9 (2021), 140896–140920.
- [162] Muzammil Shahbaz, Phil McMinn, and Mark Stevenson. 2015. Automatic generation of valid and invalid test data for string validation routines using web searches and regular expressions. Science of Computer Programming 97 (2015), 405–425.
- [163] Pojan Shahrivar, Stuart Millar, and Ezzeldin Shereen. 2023. Detecting web application DAST attacks with machine learning. In *Proceedings of the 2023 IEEE Conference on Dependable and Secure Computing (DSC'23)*. 1–8.
- [164] Elwin Shaji and Narayanan Subramanian. 2021. Assessing non-intrusive vulnerability scanning methodologies for detecting web application vulnerabilities on large scale. In Proceedings of the 2021 International Conference on System, Computation, Automation and Networking (ICSCAN'21). 1–5.
- [165] Navneet Singh, Vishtasp Meherhomji, and BR Chandavarkar. 2020. Automated versus manual approach of web application penetration testing. In *Proceedings of the 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT'20)*. 1–6.
- [166] Ferda Özdemir Sönmez and Banu Günel Kiliç. 2021. Holistic web application security visualization for multi-project and multi-phase dynamic application security test results. *IEEE Access* 9 (2021), 25858–25884.
- [167] Dimitri Michel Stallenberg and Annibale Panichella. 2019. JCOMIX: A search-based tool to detect XML injection vulnerabilities in web applications. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE'19)*. 1090–1094.
- [168] Anastasios Stasinopoulos, Christoforos Ntantogian, and Christos Xenakis. 2019. Commix: Automating evaluation and exploitation of command injection vulnerabilities in Web applications. *International Journal of Information Security* 18 (2019), 49–72.
- [169] Andrea Stocco, Maurizio Leotta, Filippo Ricca, and Paolo Tonella. 2016. Automatic page object generation with APOGEN. In *Proceedings of the 16th International Conference on Web Engineering (ICWE'16)*. 533–537.
- [170] Xin Su and Xiaohui Li. 2021. Elastic performance test method of web server in cloud computing environment. *Journal of Web Engineering* 20, 5 (2021), 1641–1658.
- [171] Chungha Sung, Markus Kusano, Nishant Sinha, and Chao Wang. 2016. Static DOM event dependency analysis for testing web applications. In Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE'16). 447–459.
- [172] Nezih Sunman, Yiğit Soydan, and Hasan Sözer. 2022. Automated web application testing driven by pre-recorded test cases. *Journal of Systems and Software* 193 (2022), 111441.
- [173] Sunitha Thummala and Jeff Offutt. 2016. Using Petri nets to test concurrent behavior of web applications. In Proceedings of the IEEE 9th International Conference on Software Testing, Verification and Validation Workshops (ICSTW'16). 189–198.
- [174] Vatsya Tiwari, Sachin Upadhyay, Jayati Krishna Goswami, and Sanjiv Agrawal. 2023. Analytical evaluation of web performance testing tools: Apache JMeter and SoapUI. In Proceedings of the IEEE 12th International Conference on Communication Systems and Network Technologies (CSNT'23). 519–523.

- [175] Pariwish Touseef, Khubaib Amjad Alam, Abid Jamil, Hamza Tauseef, Sahar Ajmal, Rimsha Asif, Bisma Rehman, and Sumaira Mustafa. 2019. Analysis of automated web application security vulnerabilities testing. In *Proceedings of the 3rd International Conference on Future Networks and Distributed Systems (ICFNDS'19)*. 1–8.
- [176] Ganesh Vaidyanathan and Steven Mautone. 2009. Security in dynamic web content management systems applications. Commun. ACM 52, 12 (2009), 121–125.
- [177] Nisal Madhushan Vithanage and Neera Jeyamohan. 2016. WebGuardia-An integrated penetration testing system to detect web application vulnerabilities. In *Proceedings of the 2016 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET'16)*. 221–227.
- [178] Fatima Waheed, Farooque Azam, Muhammad Waseem Anwar, and Yawar Rasheed. 2020. Model driven approach for automatic script generation in stress testing of web applications. In Proceedings of the 6th International Conference on Computer and Technology Applications (ICCTA'20). 46-50.
- [179] Thomas A Walsh, Phil McMinn, and Gregory M Kapfhammer. 2015. Automatic detection of potential layout faults following changes to responsive web pages (N). In *Proceedings of the 30th IEEE/ACM International Conference on Automated Software Engineering (ASE'15)*. 709–714.
- [180] Shengye Wan, Yue Li, and Kun Sun. 2019. PathMarker: Protecting web contents against inside crawlers. *Cybersecurity* 2, 1 (2019), 9.
- [181] Qian Wang, Jinan Sun, Chen Wang, Shikun Zhang, Sisi Xuanyuan, and Bin Zheng. 2020. Access control vulnerabilities detection for web application components. In Proceedings of the IEEE 6th International Conference on Big Data Security on Cloud (BigDataSecurity'20), IEEE International Conference on High Performance and Smart Computing (HPSC'20), and IEEE International Conference on Intelligent Data and Security (IDS'20). 24–28.
- [182] Siyi Wang, Sinan Wang, Yujia Fan, Xiaolei Li, and Yepang Liu. 2024. Leveraging large vision language model for better automatic web GUI testing. *arXiv* 2410.12157 (2024).
- [183] Shu-Yan Wang, Jia-Ze Sun, and Ju Zhang. 2016. The method of generating web link security testing scenario based on UML diagram. In Proceedings of the 15th IEEE International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom'16). 1831–1838.
- [184] Weiwei Wang, Xiaohong Guo, Zheng Li, and Ruilian Zhao. 2019. Test case generation based on client-server of web applications by memetic algorithm. In *Proceedings of the IEEE 30th International Symposium on Software Reliability Engineering (ISSRE'19)*. 206–216.
- [185] Wenhua Wang, Sreedevi Sampath, Yu Lei, Raghu Kacker, Richard Kuhn, and James Lawrence. 2016. Using combinatorial testing to build navigation graphs for dynamic web applications. *Software Testing, Verification and Reliability* 26, 4 (2016), 318–346.
- [186] Weiwei Wang, Shumei Wu, Zheng Li, and Ruilian Zhao. 2023. Parallel evolutionary test case generation for web applications. *Information and Software Technology* 155 (2023), 107113.
- [187] Jane Webster and Richard T Watson. 2002. Analyzing the past to prepare for the future: Writing a literature review. MIS quarterly (2002), xiii–xxiii.
- [188] Te-En Wei, Hahn-Ming Lee, Albert B Jeng, Hemank Lamba, and Christos Faloutsos. 2019. WebHound: a data-driven intrusion detection from real-world web access logs. Soft Computing 23 (2019), 11947–11965.
- [189] Jianfei Xiao, Yancan Chen, Yimin Ou, Hanyi Yu, Kai Shu, and Yiyong Xiao. 2024. Baichuan2-Sum: Instruction Finetune Baichuan2-7B Model for Dialogue Summarization. In 2024 International Joint Conference on Neural Networks (IJCNN). IEEE, 1–8.
- [190] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. A survey on large language model (LLM) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing* (2024), 100211.
- [191] Nazish Yousaf, Farooque Azam, Wasi Haider Butt, Muhammad Waseem Anwar, and Muhammad Rashid. 2019. Automated model-based test case generation for web user interfaces (WUI) from interaction flow modeling language (IFML) models. IEEE Access 7 (2019), 67331–67354.
- [192] Bing Yu, Lei Ma, and Cheng Zhang. 2015. Incremental web application testing using page object. In *Proceedings of the* 3rd IEEE Workshop on Hot Topics in Web Systems and Technologies (HotWeb'15). 1–6.
- [193] Bing Zhang, Jingyue Li, Jiadong Ren, and Guoyan Huang. 2021. Efficiency and effectiveness of web application vulnerability detection approaches: A review. *Comput. Surveys* 54, 9 (2021), 1–35.
- [194] Ming Zhang, Shuaibing Lu, and Boyi Xu. 2017. An anomaly detection method based on multi-models to detect web attacks. In Proceedings of the 10th International Symposium on Computational Intelligence and Design (ISCID'17), Vol. 2. 404–409.
- [195] Meng-ni Zhang, Can Wang, Jia-jun Bu, Zhi Yu, Yu Zhou, and Chun Chen. 2015. A sampling method based on URL clustering for fast web accessibility evaluation. Frontiers of Information Technology and Electronic Engineering 16, 6 (2015), 449–456.
- [196] Simin Zhang, Bo Li, Jianxin Li, Mingming Zhang, and Yang Chen. 2015. A novel anomaly detection approach for mitigating web-based attacks against clouds. In *Proceedings of the IEEE 2nd International Conference on Cyber Security*

- and Cloud Computing (CSCloud'15). 289-294.
- [197] Jingling Zhao and Rulin Gong. 2015. A new framework of security vulnerabilities detection in PHP web application. In Proceedings of the 9th International Conference on Innovative Mobile and Internet Services in Ubiquitous Computing (IMIS'15). 271–276.
- [198] Junzan Zhou, Bo Zhou, and Shanping Li. 2014. LTF: A model-based load testing framework for web applications. In *Proceedings of the 14th International Conference on Quality Software (QSIC'14)*. 154–163.
- [199] Yunxiao Zou, Zhenyu Chen, Yunhui Zheng, Xiangyu Zhang, and Zebao Gao. 2014. Virtual DOM coverage for effective testing of dynamic web applications. In *Proceedings of the 23rd International Symposium on Software Testing and Analysis (ISSTA'14).* 60–70.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009