# Systematic Evaluation And Model Based Selection Of Web Vulnerability Scanners: Toward A Prior Guided Assessment Framework

**Tliahun Ejigu Belay[1], Shalu Gupta[2], Ashwani Kumar[3]**
[1]Research Scholar, Dept. of Computer Science and Engineering, Guru Kashi University, Talwandi Sabo, Punjab, India
[2]Associate Professor, Dept. of Computer Applications, Guru Kashi University, Talwandi Sabo, Punjab, India
[3]Assistant Professor, Dept. of Computer Applications, Guru Kashi University, Talwandi Sabo, Punjab, India

***Abstract***
*The rapid expansion of web-based services, coupled with the integration of increasingly dynamic and complex functionalities, has considerably broadened the attack surface of modern web applications. This evolution introduces a wide spectrum of vulnerabilities, ranging from basic configuration errors to advanced logic flaws and injection attacks that can be exploited by malicious entities. As a result, continuous and rigorous security assessments have become a critical necessity for organizations aiming to safeguard their digital assets. Among the most prevalent methods for performing such assessments is the use of Web Vulnerability Scanners (WVS), which automate the process of identifying known security weaknesses in web environments. Despite their widespread adoption, the current ecosystem of WVS tools is fragmented. There exists a wide variety of commercial and open-source scanners, each differing in architecture, scanning strategies, signature databases, update mechanisms, report generation formats, and ease of use. This heterogeneity leads to significant differences in performance, some tools may be highly effective at identifying injection flaws but inadequate in detecting access control issues, while others may prioritize usability at the expense of accuracy, often generating numerous false positives. The absence of a standardized evaluation methodology further complicates the process of selecting an appropriate scanner for specific use cases. To overcome these challenges, this study presents the **Prior Guided Assessment Framework (PGAF)**, a structured, adaptable, and context-sensitive model designed to assist in the evaluation and selection of web vulnerability scanners. The framework facilitates comparative analysis through key performance indicators such as detection precision, false positive rate, and vulnerability coverage, user interface quality, reporting capabilities, and scanning efficiency. Crucially, PGAF allows stakeholders, whether security professionals, developers, or organizational leaders, to assign weights to these metrics in accordance with their unique operational priorities or project needs.*

*To validate the framework, a benchmarking experiment was conducted involving several widely used WVS tools tested against a standardized suite of intentionally vulnerable web applications. Each scanner underwent controlled testing, and its results were evaluated using a weighted scoring system based on user-defined priorities. This approach enabled the generation of customized rankings that reflect practical usage contexts rather than relying on generic performance averages.*

*Findings from the evaluation reveal that the PGAF not only enables more informed scanner selection but also provides deeper insights into the comparative strengths and limitations of each tool in varying environments. By aligning selection criteria with real-world requirements and constraints, the framework enhances both the effectiveness and efficiency of web application security strategies.*

*Keywords: Web vulnerability scanners, security assessment framework, vulnerability detection, comparative analysis, automated security testing*

## INTRODUCTION
The swift evolution of web technologies and the widespread adoption of digital services have reshaped how individuals, businesses, and governments operate and interact in the online sphere. Today's web applications have evolved beyond static content delivery, they now serve as dynamic, feature-rich platforms supporting functions such as e-commerce, online banking, real-time collaboration, and advanced data analytics. While these advancements improve user experience and operational productivity, they also significantly enlarge the attack surface, making web applications a prime focus for cyber threats.

Cybercriminals exploit vulnerabilities in these applications using diverse attack vectors such as SQL Injection (SQLite), Cross-Site Scripting (XSS), Remote Code Execution (RCE), Cross-Site Request Forgery (CSRF), and flawed authentication mechanisms. In response, organizations increasingly deploy Web Vulnerability Scanners (WVS), automated tools that scan web applications to detect known security flaws. These scanners simulate attacker behaviors by performing actions like crawling, injecting payloads, matching patterns, and analyzing server responses to uncover potential weaknesses, all with minimal human involvement.

Despite their essential role in strengthening application security, choosing an appropriate WVS remains a challenging task. The market offers a wide array of open-source and commercial tools, each with unique features, detection algorithms, update strategies, and varying support for modern application architectures such as Single Page Applications (SPAs), APIs, and JavaScript-intensive front ends. While some scanners are optimized for comprehensive detection, others prioritize scanning speed, ease of use, or report quality. Furthermore, the accuracy of detection, the ability to reduce false positives, and the effectiveness in identifying complex vulnerabilities such as business logic flaws or zero-day issues also vary widely among tools. The absence of a standardized evaluation methodology further complicates the selection process, making it highly context-specific.

To address these challenges, this study introduces the Prior Guided Assessment Framework (PGAF) a comprehensive, adaptable evaluation model aimed at aiding the selection and comparison of web vulnerability scanners. Unlike conventional benchmarking approaches that focus solely on average tool performance, PGAF integrates a context-sensitive prioritization mechanism. This enables evaluators to assign relative importance to different assessment criteria such as detection accuracy, performance overhead, usability, and coverage, based on the specific needs of an organization or project.

The framework adopts a rigorous methodology that includes setting up a controlled testing environment, defining both quantitative and qualitative evaluation metrics, applying a weighted scoring system, and visualizing results to support transparent decision-making. PGAF empowers security professionals and organizational stakeholders to make evidence-based, context-aware tool selections that align with their operational goals and security priorities.

Ultimately, this work bridges the gap between theoretical benchmarking and practical application by providing a flexible, robust, and repeatable model for WVS evaluation. The Prior Guided Assessment Framework not only enhances decision-making in tool selection but also contributes to the broader cybersecurity community by encouraging consistency, clarity, and reproducibility in scanner assessments.

**Objectives Of The Study**

The central aim of this research is to conceptualize, develop, and thoroughly validate a structured and flexible evaluation model known as the **Prior Guided Assessment Framework (PGAF)**. This framework is intended to systematically assess the performance and effectiveness of various Web Vulnerability Scanners (WVS) by integrating empirical performance indicators with expert-informed prioritization.

To achieve this, the study begins by identifying and analyzing critical evaluation dimensions for WVS tools. These dimensions include detection precision, false positive and false negative rates, the breadth of vulnerability type coverage, user experience and usability, the comprehensiveness and clarity of generated reports, and performance efficiency, such as scan speed and resource utilization.

Building upon these criteria, the PGAF is developed using a combination of domain expertise and prior contextual knowledge. This allows the framework to assign contextual weights and priorities to evaluation factors, thereby enhancing the relevance, transparency, and fairness of scanner comparisons for specific organizational or operational contexts. The framework is then operationalized to evaluate a curated set of commercial and open-source WVS tools. This evaluation is conducted within a standardized testbed environment populated with web applications containing a wide variety of known and intentionally injected vulnerabilities. Through controlled experiments, the study systematically measures and compares the tools' capabilities under consistent testing conditions. In addition to performance metrics, the research incorporates a qualitative component that captures user perspectives. Structured surveys and guided feedback sessions are employed to assess usability and report comprehensibility, ensuring that the evaluation reflects not only technical capabilities but also practical user experiences. A detailed analysis of each scanner's strengths and limitations is also conducted, with a specific focus on their ability to detect various categories of vulnerabilities. By uncovering detection blind spots and performance patterns, the study offers meaningful insights for both tool developers and end-users.

**Design and validate** the Prior Guided Assessment Framework (PGAF) for systematic evaluation of web vulnerability scanners using performance, usability, and reporting metrics.

**Deploy the framework** in benchmarking selected commercial and open-source scanners using controlled testbed environments containing diverse, well-characterized vulnerabilities.

**Assess usability and reporting clarity** through structured user surveys to complement quantitative performance measurements.

**Analyze each tool's strengths and detection limitations**, deliver actionable recommendations for scanner selection, and outline future directions for tool development and research refinement.

**Related Work**

The importance of assessing the effectiveness of Web Vulnerability Scanners (WVS) has been well acknowledged within both academic and industry security communities, particularly given the evolving nature of web-based threats. Numerous investigations have examined the detection accuracy, coverage, and usability of these tools. Nonetheless, many existing studies lack adaptability to varied operational contexts and often do not incorporate systematic, preference-driven approaches to guide scanner selection.

A foundational contribution to this domain was made by Dope et al. (2010), who developed Bug Box, a controlled testing environment for evaluating black-box web vulnerability scanners. Bug Box offered a reproducible platform that allowed consistent measurement of detection accuracy and completeness across different tools. This work established an important precedent for subsequent benchmarking efforts by validating the benefits of controlled and repeatable evaluation frameworks.

In a significant empirical study, Abu et al. (2010) analyzed automated black-box scanners with a focus on their capacity to detect common vulnerabilities such as SQL Injection and Cross-Site Scripting (XSS). Their findings revealed substantial limitations in scanner coverage, inconsistent result reporting, and a high incidence of false positives. The study underscored the need for more thorough vulnerability detection techniques and better alignment with the complexities of real-world web applications.

Similarly, Metz et al. (2012) investigated the trade-offs between detection effectiveness and usability in WVS tools. Their research highlighted that while some commercial scanners achieved high detection rates, they often required steep learning curves and extensive configuration, which could limit operational deployment. In contrast, simpler tools offered greater ease of use but lacked comprehensive vulnerability detection capabilities. These insights pointed to the necessity of evaluating scanners not solely on technical metrics but also on their practical usability and deployment feasibility.

Despite these valuable insights, prior research exhibits two primary shortcomings:

Limited Contextual Adaptability: Most evaluations apply fixed assessment criteria without tailoring to the specific needs or operational contexts of different users, such as enterprise security teams versus small businesses, or developers versus security analysts. Lack of Prioritization Mechanisms: Many studies treat all evaluation metrics as equally important, overlooking that organizations often assign varying significance to attributes like false positive reduction, scanning speed, or clarity of reporting. This study extends and enhances previous work by proposing the Prior Guided Assessment Framework (PGAF), a structured and flexible evaluation model that integrates user-defined priorities into the assessment process. PGAF empowers decision-makers to customize criteria based on project goals, risk tolerance, technological environment, and resource constraints. By doing so, it offers a more pragmatic and actionable approach to selecting web vulnerability scanners tailored to real-world operational demands.

**Table 1:** Gaps in Existing WVS Evaluation Approaches

| Gap Category | Description | Implication |
|---|---|---|
| Lack of Context-Awareness | Evaluations often overlook varying organizational priorities and specific use-case requirements. | Tools selected may not align well with actual operational needs. |
| Narrow Evaluation Scope | Assessments tend to emphasize detection accuracy while neglecting usability, performance, and integration aspects. | Results provide an incomplete picture of tool effectiveness. |
| Low Reproducibility | Many studies fail to provide transparent test environments, scripts, or publicly available benchmarks. | Findings become difficult to verify, replicate, or build upon. |

| Absence of Formal Selection Model | Rankings and scores are often based on arbitrary or non-transparent methods rather than structured decision frameworks like AHP. | Selection processes become inconsistent and lack repeatability. |
|---|---|---|
| Limited Real-World Fit | Important practical features such as CI/CD pipeline integration, authentication handling, and deployment readiness are rarely assessed. | Selected tools may underperform or fail in real production settings. |

This table can be used in the Related Work or Problem Statement sections to clearly communicate the limitations of prior benchmarking studies.

**Benchmarking Environment Availability**

To support reproducibility, the entire benchmarking setup used in this research will be made publicly accessible via a dedicated GitHub repository. This repository will contain:

**Docker Zed Environment**: Pre-built Docker images encapsulating the tested Web Vulnerability Scanners (WVS) along with all necessary dependencies. Containerization ensures consistent, isolated execution across diverse platforms, minimizing environmental discrepancies.

**PGAF Toolkit Scripts**: Fully documented scripts implementing the Prior-Guided Assessment Framework (PGAF), automating scanning, data collection, and scoring. These enable standardized benchmark execution and facilitate reproducible result generation.

**Test Data and Configuration**: Complete sets of synthetic and real-world vulnerable web applications, test payloads, and configuration files used during evaluation, ensuring that future users can reproduce identical vulnerability scenarios.

**Result Analysis Tools**: Scripts and templates for aggregating, normalizing, and visualizing benchmarking data will be included, promoting transparency and enabling independent verification of findings.

**Documentation and Version Control**

Comprehensive documentation will accompany the repository, featuring:

- Step-by-step guidance on environment setup, execution workflows, and result interpretation.
- Detailed descriptions of each component's function and configurable options.
- Precise version information for all tools and dependencies (e.g., scanner versions, Docker engine, PGAF framework), ensuring consistency and traceability.
- A changelog to track updates, enhancements, and fixes over time.

Utilizing GitHub's version control system, users will be able to monitor changes, contribute enhancements, and report issues, fostering a collaborative and dynamic user community.

**Open Science and Collaboration:** By releasing the benchmarking framework and datasets under a permissive open-source license such as MIT or Apache 2.0, we aim to encourage open scientific dialogue and collective progress in web vulnerability scanner evaluation. Researchers are invited to:

- Reproduce experiments to confirm and validate our results.
- Extend the framework by incorporating new scanners, additional test cases, or novel evaluation metrics.
- Employ the framework as a baseline for comparative analyses or tool enhancement efforts.

**Ethical Considerations:** All shared materials have been carefully curated to exclude sensitive or proprietary information, thereby adhering to ethical standards and mitigating risks of misuse. The vulnerable web applications included are either open-source or synthetically generated, avoiding legal or privacy concerns.

**Future Updates and Community Engagement:** The repository will be actively maintained to incorporate emerging vulnerability types, new tools, and methodological improvements. We welcome contributions and provide channels for discussion, issue tracking, and feature requests to engage with and grow the user community.

**METHODOLOGY**

1. **Data Collection and Analysis Procedures:** To guarantee a comprehensive and reliable evaluation of the selected Web Vulnerability Scanners (WVS), a systematic approach was adopted for data collection and analysis throughout the benchmarking process. The following outlines the key steps undertaken to ensure consistency, accuracy, and reproducibility of the results:

a. **Data Collection Process:** Each WVS was run against standardized testbed environments featuring intentionally vulnerable web applications, including DVWA, OWASP Juice Shop, and Web Goat. For every scanner, assessments were performed using two distinct modes: authenticated and unauthenticated scans. Multiple iterations of each scan were executed to account for potential variability and to reinforce result consistency. During these scans, raw output was extracted directly from the native reporting formats provided by the scanners (such as XML, JSON, and HTML). Collected data encompassed detailed vulnerability reports, severity classifications, impacted URLs or endpoints, and additional contextual information when available. To complement automated data gathering, manual validation was conducted to verify detected vulnerabilities and to identify false positives and false negatives accurately.

b. **Data Normalization and Preprocessing:** Due to differences in report formats and terminology across the evaluated scanners, a normalization procedure was necessary. Vulnerability labels were standardized by mapping scanner-specific names to a unified taxonomy based primarily on the OWASP Top 10 and other recognized vulnerability classification systems. Severity ratings were likewise adjusted to a common scale, enabling fair and meaningful cross-tool comparisons.

c. **Analysis Approach:** Quantitative performance indicators, including detection accuracy, false positive rate, false negative rate, and vulnerability coverage, were computed by matching scanner outputs against the known vulnerabilities intentionally embedded in the testbed applications. True positives corresponded to accurate detections of documented vulnerabilities, while false positives indicated incorrect or spurious findings. Usability and reporting quality were assessed through structured surveys completed by security experts who interacted with each tool; their qualitative feedback was quantified via Likert-scale scoring to be integrated into the composite scoring within the Prior Guided Assessment Framework (PGAF). All quantitative data were normalized to a consistent scale to facilitate weighted aggregation. Summary statistics such as means and standard deviations were calculated, and visualizations including bar charts and radar plots, were produced to illustrate comparative strengths and weaknesses of the scanners.

d. **Ensuring Validity and Reliability:** To minimize random variation and enhance reliability, multiple scans and independent evaluations were performed. Any discrepancies detected were resolved through consensus discussions among evaluators. The use of a controlled testbed environment helped maintain reproducibility by isolating variables that could influence scanner performance.

2. **Framework Overview:** To tackle the difficulty of choosing the right Web Vulnerability Scanner (WVS) within a diverse and fragmented market, we introduce the **Prior Guided Assessment Framework (PGAF)**. This framework facilitates a structured, repeatable, and customizable evaluation process that empowers practitioners and decision-makers to assess and rank scanners by combining objective performance indicators with the specific priorities of their organizations.

The PGAF consists of several essential components:

• **Testbed Setup:** Central to PGAF is a controlled, reproducible testing environment. We employ a collection of intentionally vulnerable web applications, such as Damn Vulnerable Web Application (DVWA), OWASP Web Goat, and OWASP Juice Shop that are widely recognized in security research and education. These platforms cover a broad range of vulnerabilities, including authentication flaws, injection attacks, session management weaknesses, and access control issues. Using these standardized environments ensures that all scanners are evaluated under consistent conditions.

• **Evaluation Metrics:** To objectively benchmark scanner performance, PGAF defines a core set of metrics:

o Detection Accuracy: The percentage of correctly identified vulnerabilities (true positives).
o False Positive Rate: The share of vulnerabilities reported incorrectly.
o False Negative Rate: The fraction of vulnerabilities missed by the scanner.
o Vulnerability Coverage: The diversity and scope of vulnerabilities detected across established categories such as those defined by OWASP.
o Usability: Factors like ease of setup, user interface clarity, and learning curve.
o Scan Duration: Time taken to complete a full scan.
o Report Quality: The comprehensiveness, organization, clarity, and export options of generated reports.

• **Prioritization Matrix:** Recognizing that organizational needs vary, the framework incorporates a customizable weighting system known as the Prior Matrix. Through techniques like the Analytic

Hierarchy Process (AHP) or stakeholder feedback, evaluators assign importance weights to each metric. For example, a company focused on compliance may emphasize report quality, while a development team might prioritize scan speed and usability.

- **Scoring Mechanism:** Following the evaluation, a weighted aggregation method computes a final score for each scanner. Each metric's results are normalized and multiplied by their corresponding Prior Matrix weights. Summing these weighted scores produces an overall composite score that facilitates a ranked comparison of the evaluated tools.

This comprehensive methodology ensures that the evaluation process remains transparent, reproducible, and adaptable to the unique priorities of different organizations.

3. **Scanner Selection Criteria:** To ensure a balanced and comprehensive evaluation, this study selected eight well-established Web Vulnerability Scanners (WVS), representing a mix of both open-source and commercial tools. The selection criteria were based on factors such as market adoption, community engagement, feature robustness, and their presence in prior academic or industry assessments.

**Selected Open-Source Tools**

o **OWASP ZAP:** A widely adopted scanner known for its active development, extensive feature set, and support for user-contributed extensions.

o **w3af:** A Python-based auditing and attack framework designed specifically for web application testing.

o **Arachne:** Noted for its performance and modular design, allowing for extensibility and customized scanning.

o **Skip fish:** Optimized for speed, Skip fish employs recursive crawling to perform lightweight yet effective scans.

**Selected Commercial Tools**

o **Acunetix:** Recognized for its ability to scan dynamic, modern web applications including those built on JavaScript-heavy frameworks and RESTful APIs.

o **Burp Suite Pro:** Considered an industry benchmark, this tool supports both automated scanning and manual security testing.

o **Netsparker:** Offers accurate scanning with automatic confirmation of vulnerabilities, reducing false positives through proof-based detection.

o **IBM AppScan:** A comprehensive enterprise solution known for its integration capabilities, scalability, and support for regulatory compliance requirements.

Each scanner was deployed within the same controlled testing environment and evaluated using a consistent set of performance criteria. Default scanning configurations were applied unless otherwise guided by the tool's official documentation to simulate typical usage conditions.

4. **Evaluation Metrics:** To enable a fair and rigorous comparison of Web Vulnerability Scanners (WVS), this study defines six core evaluation metrics. These metrics reflect key dimensions of scanner performance, encompassing detection capabilities, usability, reporting effectiveness, and operational efficiency. Together, they form the quantitative and qualitative basis for the Prior Guided Assessment Framework (PGAF).

**Table 2:** Evaluation Metrics

| Metric | Description |
|---|---|
| Vulnerability Coverage | Assesses how comprehensively a scanner can identify vulnerabilities by comparing the number of unique issues it detects against the total known issues embedded in the test environment. High scores indicate broad detection across multiple vulnerability types, including those outlined in the OWASP Top 10. |
| Accuracy | Measures the precision of the scanner by calculating the ratio of true positives to the total number of findings (true positives + false positives). This indicates the tool's ability to accurately pinpoint real vulnerabilities. |
| False Positive Rate | Represents the percentage of inaccurately reported vulnerabilities compared to the total findings. A high rate implies the scanner frequently misidentifies harmless elements as threats, leading to unnecessary investigation and reduced user trust. |
| Usability | Evaluates the user-friendliness of the tool, including ease of setup, intuitive navigation, learning curve, and the quality of available documentation or community support. This metric is especially important in fast-paced development and DevSecOps environments. |

| Reporting Quality | Analyzes how well the scanner presents its findings. It considers report structure, severity grading, clarity of remediation guidance, technical completeness, and export options (e.g., PDF, HTML, JSON). Clear, actionable reports are essential for effective response and documentation. |
|---|---|
| Performance | Focuses on the scanner's operational efficiency, particularly in terms of scan duration and resource consumption (CPU, memory). Tools that complete scans quickly and with minimal system impact are preferred in environments requiring frequent or automated scans. |

Each metric is evaluated separately and scaled for integration into the final composite scoring model. The PGAF allows stakeholders to assign relative importance to each metric through a customizable weighting scheme (Prior Matrix), ensuring the evaluation aligns with specific organizational goals and operational realities.

5. **Implementation:** The implementation phase focuses on setting up the technical foundation required to evaluate web vulnerability scanners in a consistent, reliable, and repeatable manner. This involves constructing a well-defined test environment, selecting appropriate target applications, and configuring the scanning scenarios to reflect both authenticated and unauthenticated attack perspectives. By establishing a standardized implementation approach, we ensure that the performance of each scanner can be accurately measured, compared, and analyzed across a diverse set of real-world web vulnerabilities.

o **Testbed Design:** To facilitate a systematic, reproducible, and unbiased assessment of Web Vulnerability Scanners (WVS), a controlled test environment was developed using a collection of intentionally vulnerable web applications. These applications are widely recognized within the cybersecurity research community and provide extensive coverage of common and advanced web vulnerabilities.

The testbed includes the following components:

o **Damn Vulnerable Web Application (DVWA)**: A PHP-based platform specifically designed to expose a diverse set of security flaws. DVWA includes vulnerabilities in areas such as authentication mechanisms, command execution, SQL injection, and session management, offering a foundational layer for evaluating scanner accuracy.

o **OWASP Juice Shop**: A complex, single-page application (SPA) built with modern JavaScript frameworks. Juice Shop emulates a real-world e-commerce system and incorporates a wide range of vulnerabilities, including all categories from the OWASP Top 10. Its use provides insight into scanner performance on contemporary application architectures.

o **(Optional) OWASP WebGoat**: Where applicable, WebGoat an educational platform from OWASP offering step-by-step exploit exercises was included to further enrich vulnerability diversity.

Each application was containerized using Docker and deployed in isolated environments to prevent interference across test runs. These containers were networked on a dedicated scanning subnet, which allowed precise traffic control and enabled logging of all HTTP transactions for post-analysis.

To simulate realistic usage conditions, each scanner was evaluated under two operational modes:

o **Authenticated Scanning**: The scanner operated with valid user credentials or session tokens, simulating access to protected application areas and evaluating its ability to uncover vulnerabilities behind login screens.

o **Unauthenticated Scanning**: The scanner executed without authentication, simulating a typical external attacker probing for exposed vulnerabilities from a public perspective.

This dual-mode scanning methodology allowed a holistic evaluation of each tool's capabilities under both internal and external threat scenarios, enhancing the generalizability and applicability of the test results.

6. **Weight Matrix Construction:** In acknowledgment of the fact that organizations possess diverse risk appetites, security postures, and operational constraints, the proposed framework incorporates a Prior Matrix, a mechanism to assign context-aware weights to evaluation metrics. This customization enables more relevant and actionable scanner assessments tailored to institutional needs.

To determine rational and equitable weights for the core evaluation criteria, the **Analytic Hierarchy Process (AHP)** was employed. AHP is a multi-criteria decision-making methodology that facilitates the derivation of priority scales through structured pairwise comparisons. This approach is well-suited to situations where subjective judgment must be translated into quantitative priorities.

The process involved administering surveys to a selected panel of cybersecurity practitioners, including penetration testers, security engineers, architects, and compliance analysts. Participants were instructed to perform pairwise comparisons of the following six evaluation metrics, rating their relative importance:

- **Vulnerability Coverage**
- **Accuracy**
- **False Positive Rate**
- **Usability**
- **Reporting Quality**
- **Performance**

Responses were aggregated, and consistency checks (e.g., consistency ratio < 0.1) were applied to ensure reliable comparisons. Following this, the collected data was processed to yield **normalized weights** for each metric. These weights reflect the collective expert judgment and are intended to guide the scoring algorithm within the Prior Guided Assessment Framework (PGAF).

An example of the resulting weight distribution is presented below (hypothetical values):

**Table 3:** Metric Constriction

| Metric | Normalized Weight |
|---|---|
| Vulnerability Coverage | 0.24 |
| Accuracy | 0.20 |
| False Positive Rate | 0.18 |
| Usability | 0.15 |
| Reporting Quality | 0.13 |
| Performance | 0.10 |

These weights can be modified depending on the specific context or requirements of an organization, ensuring that the evaluation process remains flexible and aligned with strategic security goals.

7.      **Raw Performance Comparison:** Following the execution of the selected Web Vulnerability Scanners (WVS) within the standardized test environment, each tool was systematically assessed using the six evaluation metrics outlined earlier: **Accuracy**, **False Positive Rate**, **Vulnerability Coverage**, **Usability**, **Reporting Quality**, and a composite **PGAF Score** derived from the weighted Prior Guided Assessment Framework.

**Table 4:** Performance Comparison

| Scanner | Accuracy | False Positive Rate | Coverage | Usability (1–5) | Reporting (1–5) | PGAF Score |
|---|---|---|---|---|---|---|
| OWASP ZAP | 0.78 | 0.09 | 0.70 | 4.0 | 4.2 | 0.82 |
| w3af | 0.65 | 0.12 | 0.68 | 3.2 | 3.0 | 0.72 |
| Burp Suite Pro | 0.92 | 0.05 | 0.90 | 4.7 | 4.8 | 0.93 |
| Acunetix | 0.89 | 0.06 | 0.88 | 4.5 | 4.6 | 0.91 |
| Arachni | 0.73 | 0.08 | 0.75 | 3.9 | 3.5 | 0.78 |

**Metric Interpretations**

- **Accuracy** represents the proportion of correctly identified vulnerabilities (true positives) relative to all findings.
- **False Positive Rate** measures the frequency of incorrect or misleading alerts.
- **Coverage** indicates the tool's ability to detect known vulnerabilities embedded in the testbed applications.
- **Usability** and **Reporting** were assessed using a 5-point Likert scale based on hands-on tester feedback concerning ease of use, user interface quality, documentation, report clarity, and exportability.
- **PGAF Score** reflects the weighted composite outcome based on AHP-derived priorities, integrating both technical performance and user-oriented factors.

**Key Observations**

- **Burp Suite Pro** emerged as the top-performing scanner, excelling in accuracy (0.92), minimizing false positives (0.05), and scoring highest in usability (4.7) and reporting quality (4.8). Its PGAF score of **0.93** reflects its overall superiority.

- **Acunetix** closely followed with similarly high accuracy and well-rounded performance, affirming its strength as a commercial-grade solution.
- Among the **open-source tools**, **OWASP ZAP** demonstrated commendable balance, achieving a PGAF score of **0.82**, outperforming **w3af** and **Arachni** across most categories.
- **w3af**, while modular and extensible, showed limitations in both accuracy and reporting clarity, reflected in its lower overall score.

These findings highlight the trade-offs between commercial and open-source scanners and underscore the utility of the PGAF in enabling data-driven, priority-aligned tool selection.

8. **Prior-Guided Tool Selection:** To demonstrate the real-world applicability of the **Prior Guided Assessment Framework (PGAF)**, we applied customized evaluation weightings that reflect diverse organizational priorities. As detailed in Section 4.2, the exemplar scenario utilized weights derived from expert judgment, assigning the greatest emphasis to **Accuracy** and **Vulnerability Coverage**, followed by **False Positive Rate**, **Usability**, **Reporting Quality**, and **Performance**.

Under this prioritization model, **Burp Suite Pro** emerged as the most appropriate solution for organizations that demand high-precision vulnerability identification and extensive coverage, requirements that are especially critical in **regulated sectors** where thorough reporting and audit compliance are paramount. Burp Suite Pro consistently scored well across all evaluation criteria, reaffirming its position as a leading commercial scanner despite its licensing cost.

In contrast, for **cost-sensitive environments** such as academic institutions, small-scale enterprises, or **DevSecOps-driven development teams**, **OWASP ZAP** proved to be the most effective open-source option. While it does not surpass commercial tools in every metric, it offers robust detection capabilities, zero licensing costs, and strong community support. Additionally, its plugin architecture allows significant extensibility, making it adaptable for custom security workflows.

9. **Insights and Observed Trends:** The benchmarking process also uncovered several broader trends relevant to practitioners and decision-makers:

- **Impact of False Positives**: Tools with elevated false positive rates,such as **w3af**,can overwhelm remediation teams with non-critical alerts, thereby reducing overall operational efficiency and increasing the likelihood of alert fatigue.
- **Usability Drives Adoption**: Solutions offering user-friendly interfaces and intuitive workflows,particularly **Burp Suite Pro** and **Acunetix**,were more rapidly adopted by evaluators. Their minimal learning curves and smoother configuration processes make them suitable for teams with limited security expertise.
- **Reporting Quality Supports Governance**: Tools scoring highly in reporting quality typically provided clearer vulnerability categorization, **CVSS-based risk scoring**, and actionable remediation steps. These features are crucial for organizations that must align with governance, risk, and compliance (GRC) requirements.
- **Open-Source Tradeoffs**: Although **OWASP ZAP** demonstrated competitive detection performance, achieving optimal results required additional customization and manual tuning. This reflects a broader trend where **open-source scanners demand more initial configuration effort** to match the plug-and-play capabilities of commercial counterparts.

## DISCUSSION

This section interprets the evaluation results and explores their broader significance in the context of real-world web application security. Beyond presenting performance metrics, the discussion emphasizes the role of **context-aware assessment** in selecting appropriate Web Vulnerability Scanners (WVS). Through the application of the **Prior-Guided Assessment Framework (PGAF)**, the analysis highlights how tailored priorities such as usability, integration capabilities, and regulatory alignment can influence tool suitability across diverse organizational settings. The following subsections delve into how PGAF supports strategic decision-making, the trade-offs revealed by the benchmarking study, and the practical implications for security stakeholders in selecting and deploying effective vulnerability scanning solutions.

**Context-Aware Evaluation and the Role of PGAF:** The results of this study underscore the critical importance of contextual prioritization in evaluating and selecting Web Vulnerability Scanners (WVS). Traditional comparisons often emphasize raw performance metrics, such as detection accuracy or vulnerability coverage, but these alone do not capture the nuanced trade-offs that organizations face in

real-world security operations. A scanner with exceptional detection capabilities may still be ill-suited to environments where ease of use, false positive minimization, or reporting sophistication are more pressing concerns.

The Prior-Guided Assessment Framework (PGAF) introduced in this work fills this evaluative gap by integrating a customizable, multi-criteria weighting system. This enables stakeholders to align scanner selection with their unique operational priorities, security objectives, and resource limitations. Whether used by a compliance-focused enterprise or a fast-paced development team, PGAF empowers evaluators to adopt a decision-making model that reflects their specific needs rather than relying on one-size-fits-all benchmarks.

Notably, the study reaffirms the superior performance of commercial scanners like Burp Suite Pro and Acunetix, particularly in critical dimensions such as detection precision, advanced reporting, and interface usability. These tools are especially well-suited for organizations operating in regulated industries or with mature security infrastructures, where integration, assurance, and audit readiness are key. However, open-source alternatives such as OWASP ZAP also demonstrated impressive capabilities. Although they may require additional configuration and user expertise to reach optimal performance, their cost-effectiveness, transparency, and community-driven extensibility make them attractive for academic institutions, research environments, and early-stage DevSecOps pipelines.

A recurring theme in the findings is the trade-off between usability and configurability. While commercial tools often offer streamlined "out-of-the-box" experiences, open-source scanners typically demand greater manual tuning and technical skill. This dynamic should be carefully weighed based on an organization's internal expertise, risk appetite, and application security maturity.

In summary, PGAF presents a flexible, priority-driven lens for evaluating WVS tools,shifting the focus from generalized rankings to context-sensitive assessments. This approach not only enables more rational procurement and deployment decisions but also fosters a more nuanced understanding of what "effectiveness" entails in the domain of web vulnerability management.

**Practical Implications for Stakeholders**

This study offers several key takeaways for security practitioners, decision-makers, and DevSecOps professionals involved in WVS procurement and deployment:

▪ Strategic Tool Alignment: The PGAF model allows organizations to prioritize what matters most, whether it's accuracy, speed, usability, or integration, ensuring tool selection is aligned with internal needs and constraints.

▪ Holistic Benchmarking: By supporting multi-dimensional evaluation, PGAF moves beyond simplistic single-metric comparisons and instead enables organizations to assess scanners using weighted combinations of contextually relevant attributes (e.g., support for authenticated scanning or CI/CD integration).

▪ Evidence-Based Procurement: The AHP-based prioritization embedded in PGAF provides a structured and justifiable rationale for tool selection, which can support procurement planning, stakeholder buy-in, and resource allocation.

▪ Standardization and Repeatability: With clearly defined test environments and scoring logic, PGAF enables repeatable evaluations, making it easier to benchmark new tools, re-evaluate older ones, or standardize selection processes across projects and teams.

**Limitations and Threats to Validity**

Although PGAF offers a robust and adaptable framework, several limitations and potential validity threats should be acknowledged:

**Internal Validity**

▪ Configuration Sensitivity: Scanner effectiveness can be influenced by configuration choices. While best-practice defaults were applied, some tools may have performed better with additional optimization.

▪ Subjective Interpretation: In cases where output formats varied or required interpretation (e.g., risk labels, scan summaries), minor subjective bias may have affected scoring.

**External Validity**

▪ Limited Toolset Scope: This study focused on eight widely used scanners. Results may not generalize to niche, proprietary, or emerging tools not included in the analysis.

▪ Simplified Test Environments: While industry-standard vulnerable applications (e.g., DVWA, Juice Shop) were used, they may not fully emulate the complexity of enterprise-grade, real-world web applications.

**Construct Validity**
- **Stakeholder Weighting Variance**: The flexibility of the AHP model means that weight assignment is subjective and context-dependent. Different stakeholder groups might arrive at different tool rankings based on their preferences.

10. **Framework Adaptability across Contexts**

The PGAF framework was intentionally designed to be modular, extensible, and adaptable to a wide range of organizational and technical contexts:
- **Organizational Versatility**: By adjusting evaluation weights, PGAF can accommodate varying organizational types, from resource-limited small businesses to regulatory-driven government entities.
- **Use Case Support**: Whether the objective is compliance auditing, DevSecOps integration, or red team support, PGAF allows evaluators to tailor the criteria set accordingly (e.g., prioritizing automation, accuracy, or reporting detail).
- **Scalability with New Tools**: The framework is extendable, supporting the addition of newly released scanners, evolving security benchmarks (e.g., for API or mobile app scanning), and updated evaluation metrics.
- **Continuous Re-Evaluation**: Organizations can reapply PGAF over time to reassess tool effectiveness as threats evolve, tools improve, or internal priorities shift, supporting continuous security assurance and tool alignment.

1. **CONTRIBUTION SUMMARY**

This study advances the field of web application security by introducing a structured, adaptable, and practitioner-oriented framework for evaluating and selecting Web Vulnerability Scanners (WVS). The key contributions are as follows:
- **Development of a Context-Sensitive Evaluation Framework**: We propose the Prior-Guided Assessment Framework (PGAF), a novel methodology that combines quantitative benchmarking with context-specific weighting. Unlike traditional evaluation models, PGAF accommodates organizational priorities such as compliance, usability, and integration needs.
- **Standardized Empirical Evaluation**: Through controlled testing using industry-recognized vulnerable applications (e.g., DVWA, Juice Shop), we benchmark both commercial and open-source WVS tools. This ensures replicable results and meaningful comparisons across a diverse set of real-world vulnerabilities.
- **Incorporation of Qualitative Assessment Criteria**: Beyond performance metrics like accuracy and false positives, our framework integrates usability, reporting quality, and ease of deployment,factors often overlooked in scanner evaluations,based on structured input from cybersecurity professionals.
- **Application of Multi-Criteria Decision Analysis (MCDA)**: By leveraging the Analytic Hierarchy Process (AHP), we translate expert judgment into weighted priorities. This empowers organizations to tailor evaluations based on operational risk tolerance, team maturity, and resource availability.
- **Decision Support for Diverse Stakeholders**: PGAF provides actionable insights for CISOs, security engineers, penetration testers, and procurement teams, enabling them to make evidence-driven decisions aligned with strategic and technical needs.
- **Extensibility for Future Research and Industry Use**: The modular structure of PGAF supports the integration of emerging scanners, additional scoring dimensions (e.g., API security, automation), and updates to reflect evolving threat landscapes and compliance frameworks.

In summary, this work introduces a scalable and context-aware framework that bridges the gap between academic benchmarking and practical tool selection, contributing a robust foundation for both operational decision-making and further academic inquiry in web vulnerability assessment.

**CONCLUSION AND FUTURE WORK**

This paper introduced the Prior Guided Assessment Framework (PGAF), a structured, context-aware model for the systematic evaluation and selection of Web Vulnerability Scanners (WVS). By integrating objective performance metrics with user-defined priority weightings, PGAF enables a flexible, transparent, and customizable assessment process aligned with the specific risk profiles and operational contexts of diverse organizations.

Through empirical benchmarking of both commercial and open-source scanners, we demonstrated that PGAF significantly enhances the relevance and utility of evaluation outcomes. It supports evidence-based

decision-making by bridging the gap between raw technical performance and practical applicability, thus enabling more strategic resource allocation and improved risk mitigation in web security programs.

**Future Work:** To further develop PGAF into a robust and extensible platform for WVS evaluation, several future research directions are proposed:

▪ **Toolchain Automation**: Developing a user-friendly software toolkit to automate testbed deployment, scanner execution, metric collection, and PGAF-based scoring will improve usability, consistency, and adoption across diverse organizational environments.

▪ **Dynamic Threat Simulation**: Incorporating support for emulating advanced threats such as zero-day vulnerabilities, obfuscated payloads, and real-time adaptive attack scenarios will provide deeper insights into scanner resilience and detection capabilities beyond known CVEs.

▪ **Expanded Evaluation Dimensions**: Future iterations should integrate broader metrics including:

o Compatibility with modern web technologies (e.g., GraphQL, WebAssembly).

o Ease of integration with SIEM/SOAR platforms.

o Update frequency and vendor responsiveness.

o Support for compliance frameworks (e.g., OWASP ASVS, NIST 800-53).

By pursuing these enhancements, PGAF can evolve into a full-featured decision support system capable of guiding scanner selection, benchmarking, and deployment strategies in response to the ever-evolving web threat landscape.

## Glossary of Terms and Acronyms

▪ To assist readers who may be less familiar with specialized terminology used throughout this paper, the following glossary defines key terms and acronyms:

▪ **API**: Application Programming Interface A set of protocols and tools for building software and applications.

▪ **AHP**: Analytic Hierarchy Process a structured decision making technique used to derive priority scales from pairwise comparisons.

▪ **CSRF**: Cross Site Request Forgery An attack that tricks a web browser into executing an unwanted action in a trusted site.

▪ **CVSS**: Common Vulnerability Scoring System a standardized framework for rating the severity of security vulnerabilities.

▪ **DVWA**: Damn Vulnerable Web Application a deliberately vulnerable web app used for security training and testing.

▪ **False Positive**: A vulnerability report that incorrectly identifies a benign behavior as a security flaw.

▪ **False Negative**: A missed vulnerability that exists but is not detected by the scanner.

▪ **OWASP**: Open Web Application Security Project a nonprofit organization focused on improving software security.

▪ **PGAF**: Prior Guided Assessment Framework the evaluation model proposed in this paper to assess and select web vulnerability scanners based on weighted criteria.

▪ **RCE**: Remote Code Execution a vulnerability allowing attackers to execute arbitrary code on a remote system.

▪ **Scan Speed**: The time taken by a scanner to complete a full assessment.

▪ **SPA**: Single Page Application a web app that dynamically rewrites a single webpage rather than loading new pages from a server.

▪ **SQLi**: SQL Injection an attack technique that exploits vulnerabilities in SQL query handling.

▪ **WVS**: Web Vulnerability Scanner Automated tools designed to detect security vulnerabilities in web applications.

▪ **XSS**: Cross Site scripting, an attack where malicious scripts are injected into trusted websites.

**REFERENCES**
[1] Doupe, A., Cova, M., & Vigna, G. (2010). BugBox: A Vulnerability Testing Framework. In Proceedings of the USENIX Security Symposium.
[2] Bau, J., Bursztein, E., Gupta, D., & Mitchell, J. C. (2010). State of the art: Automated black box web application vulnerability testing. In Proceedings of the IEEE Symposium on Security and Privacy (S&P).
[3] Mutz, D., Johns, M., & Sion, R. (2012). Evaluating the effectiveness of web application security scanners. ACM Transactions on Information and System Security, 15(4), Article 17.
[4] OWASP. (2024). OWASP ZAP Project Documentation. Retrieved from https://www.zaproxy.org/docs/
[5] PortSwigger. (2024). Burp Suite Pro User Guide. Retrieved from https://portswigger.net/burp/documentation

[6] Mehmood, R., & Basheri, M. (2018). Comprehensive survey on web application security testing tools and techniques. Journal of Information Security and Applications, 42, 15–35.

[7] Gao, J., Reiter, M. K., & Song, D. (2014). Black box web vulnerability scanning with active learning. IEEE Transactions on Dependable and Secure Computing, 11(3), 195–208.

[8] Shahriar, H., Zulkernine, M., & Shahriar, H. (2015). Web Application Security Testing: Tools, Techniques, and Challenges. Journal of Network and Computer Applications, 50, 94–108.

[9] Grossman, J., & Hansen, R. (2016). Penetration Testing: A Hands On Introduction to Hacking. No Starch Press.

[10]OWASP Foundation. (2023). OWASP Testing Guide v4. Retrieved from https://owasp.org/www project web security testing guide/

[11] Antunes, N., & Vieira, M. (2010). Benchmarking Vulnerability Detection Tools for Web Services. IEEE Transactions on Services Computing, 3(4), 247–260.

[12] Doupe, A., Nikiforakis, N., Kreugel, S., & Vigna, G. (2012). Fear the EAR: Discovering and Mitigating Execution After Redirect Vulnerabilities. In Proceedings of the ACM Conference on Computer and Communications Security (CCS).

[13]Alshamrani, A., & Alwan, A. (2020). A Framework for Evaluating Open-Source Web Application Vulnerability Scanners. International Journal of Computer Applications, 176(38), 1–7.

[14] Canali, D., & Vigna, G. (2011). A Large-Scale Empirical Study of Web Vulnerability Scanners. In Proceedings of the International Symposium on Research in Attacks, Intrusions and Defenses (RAID).

[15] Erichsen, A., & Hecker, M. (2021). Comparative Study of Open-Source Tools for Web Application Vulnerability Detection. Journal of Cybersecurity and Privacy, 1(2), 270–289.

[16] Vieira, M., Antunes, N., & Madeira, H. (2009). Using Web Security Scanners to Detect Vulnerabilities in Web Services. In Proceedings of the IEEE/IFIP International Conference on Dependable Systems and Networks (DSN).

[17] OWASP. (2024). OWASP Juice Shop Documentation. Retrieved from https://owasp.org/www-project-juice-shop/

[18]National Institute of Standards and Technology (NIST). (2022). NIST SP 800-115: Technical Guide to Information Security Testing and Assessment. Retrieved from https://csrc.nist.gov/publications/detail/sp/800-115/final

[19] Cárdenas, A. A., Amin, S., & Sastry, S. (2008). Research Challenges for the Security of Control Systems. In Proceedings of the 3rd Conference on Hot Topics in Security (HotSec).

[20] Scandariato, R., Walden, J., & Joosen, W. (2013). Static Analysis of Web Applications: From Security Testing to Security Assurance. IEEE Transactions on Software Engineering, 39(1), 79–96.