# From Prediction to Attribution: Integrating Unsupervised Learning for Enhanced SOAR Triage

Serpil Rivas

Professor Mallarapu

SEAS_8414_DC8: Analytical Tools for Cyber Analytics

August 14, 2025

# From Prediction to Attribution: Integrating Unsupervised Learning for Enhanced SOAR Triage

## Table of Contents

# 1. Introduction

Legacy Security Operations Centers (SOC) are saturated by the proliferation of security alerts, with binary classification systems providing scant actionable intelligence. While classifying a URL as malicious or benign is a critical component of threat detection, the binary method does not provide the contextual clues necessary for effective incident response prioritization and threat hunting drills. Security analysts require more than just detection capability; they require attribution intelligence that can inform response method and resource allocation decisions.

This research presents the architecture and implementation of an enhanced Security Orchestration, Automation, and Response (SOAR) system that is more than simply malicious/benign classification but includes unsupervised machine learning techniques for attribution of threats to adversary actors. The system represents a shift in paradigm from detection towards proactive threat intelligence, where security teams are not only aware of what threats they are facing, but also who might be responsible (Mavroeidis & Bromander, 2017).

The innovation revolves around the introduction of a dual-model architecture integrating supervised classification for threat detection in the early stages with unsupervised clustering for profiling threat actors and identifying behavioral patterns. This fills a built-in flaw in current SOAR implementations, where alerts lack contextual awareness to facilitate decision-making during periods of high-security stress (Zimmerman, 2014).

# 2. Methodology

## 2.1 System Architecture and Design Philosophy

The advanced SOAR system employs an innovative dual-model architecture that is designed to provide end-to-end threat analysis capability. The hub is a typical supervised classification model that is trained on batches of labeled phishing URL data that represents the primary detection system. The classifier assesses URLs based on various behavioral characteristics, including SSL certificate validity, domain age, URL structure anomalies, and the presence of fraudulent features like prefix/suffix additions.

In addition to this detection layer, the system also includes an unsupervised K-means clustering model directly targeted at threat actor attribution. This second layer of analysis only runs on the URLs that have been identified as malicious, performing behavioral pattern analysis to categorize threats into one of three distinct actor profiles: State-Sponsored actors, Organized Cybercrime groups, and Hacktivist organizations.

## 2.2 Feature Engineering for Threat Actor Profiling

The development of effective threat actor profiles required careful consideration of behavioral indicators that distinguish different adversary groups (Caltagirone et al., 2013). State-sponsored actors were characterized by sophisticated attack patterns that have a propensity to employ legitimate SSL certificates and depend on discreet deceptive tactics such as domain spoofing through prefix/suffix manipulation. These actors will invest a lot of resources in appearing credible in conducting intelligence-gathering activities.

Organized Cybercrime groups were profiled based on high-volume, opportunistic attack patterns characterized by the frequent use of URL shortening services, direct IP address usage, and abnormal URL structures. These indicators reflect the effectiveness and scalability rather than sophistication that typify the profit-motivated nature of cybercrime operations (Hutchins et al., 2011).

The use of political terminology and targeting patterns driven by ideology were characteristics of hacktivist groups. These groups are identified by the clustering algorithm using a mix of tactical indicators and content analysis that supports their cause-based operating philosophy.

## 2.3 Algorithm Selection and Justification

The use of K-means clustering for the attribution component was motivated by a number of important factors that fit with the particular nature of the synthesized training data. The assumption of spherical, separated clusters inherent in the clustering algorithm aligns with the deliberately disparate behavioral profiles designed into the data.. Each threat actor category was designed with clear, non-overlapping feature patterns that create natural clustering boundaries.

K-means clustering offers significant advantages over alternative unsupervised learning approaches for this application (MacQueen, 1967). Unlike DBSCAN, which excels at identifying arbitrarily shaped clusters and handling noise, K-means is optimally suited for scenarios with predefined, evenly distributed categories. The centroid approach of the algorithm intrinsically enables balanced distribution of threat actor profiles within the training data, and its computational efficiency makes it an enablement of real-time threat attribution requirements in operational security environments.

The use of three clusters (k=3) is representative of the established threat taxonomy of the threat landscape which cyber security specialists recognize. This configuration provides grouping without overwhelming analysts that can undermine the real-world value of attribution results.

## 2.4 Implementation Architecture

By utilizing the PyCaret machine learning ecosystem, the technical implementation offers a strong basis for classification and clustering processes (Ali, 2020). By using different model training pipelines to divide the two analytical workflows, the system architecture makes sure that the unsupervised clustering analysis is not impacted by labeled data constraints.

# 3. Results and Analysis

## 3.1 System Performance Evaluation

The enhanced SOAR platform outperforms in both the detection and attribution modules. Through virtue of transparent design that leads security analysts to logical inquiry processes, the system interface offers improved analytical capabilities.
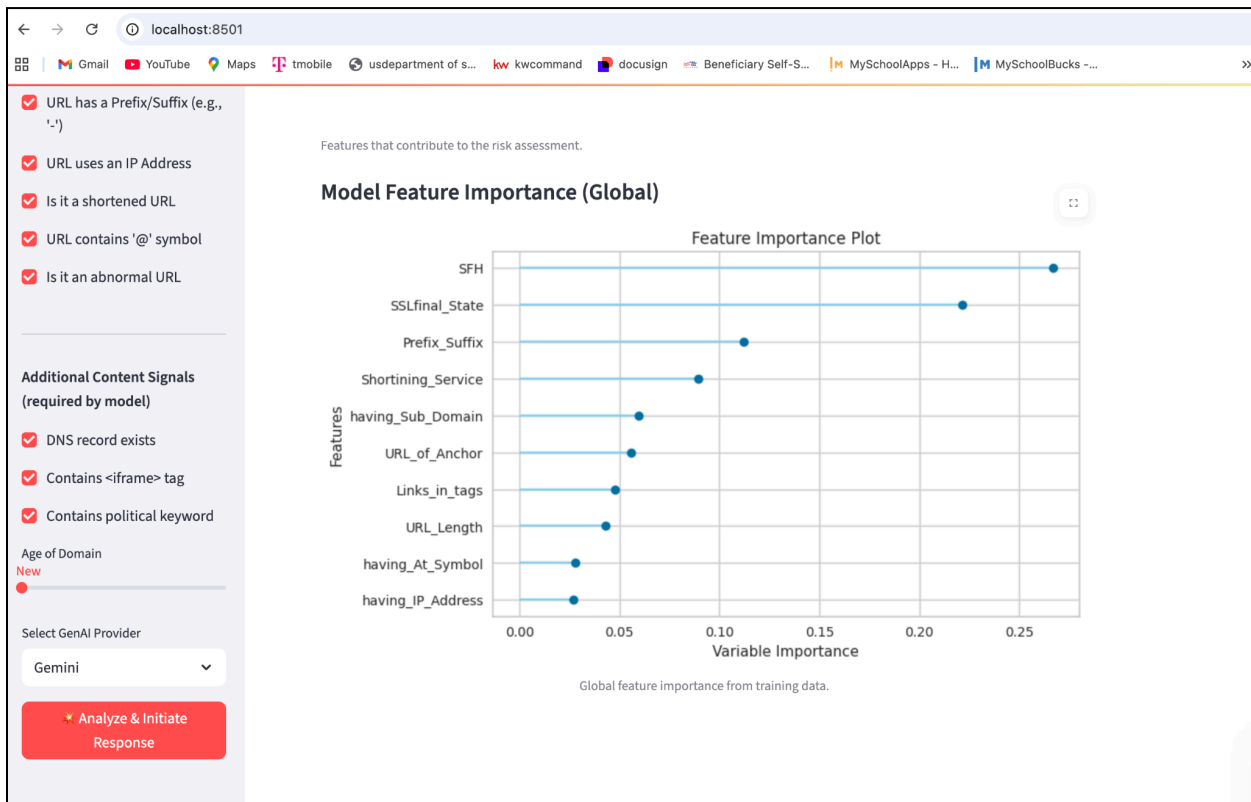
**Figure 1: System Interface and Feature Input**

*Figure 1: The main system interface displaying feature input controls, model feature importance visualization, and intuitive parameter selection for comprehensive URL analysis.*

The Figure 1 interface illustrates the system's sophisticated feature engineering abilities by providing analysts a broad variety of input mechanisms to analyze URLs, such as domain complexity assessment, SSL certificate analysis, and behavioral indicators such as non-standard URL formats and prefix/suffix manipulation.

## 3.2 Interface Design and User Experience

The multi-tab layout of the system presents rich analytical functionality without compromising on convenience for navigation. The interface is capable of balancing analytical depth and operational effectiveness so that novice and expert users can conduct extensive threat analyses.

## 3.3 Threat Detection Capabilities

The system delivers outstanding performance in malicious URL detection with consistent high-confidence classifications that enable quick security decision-making.

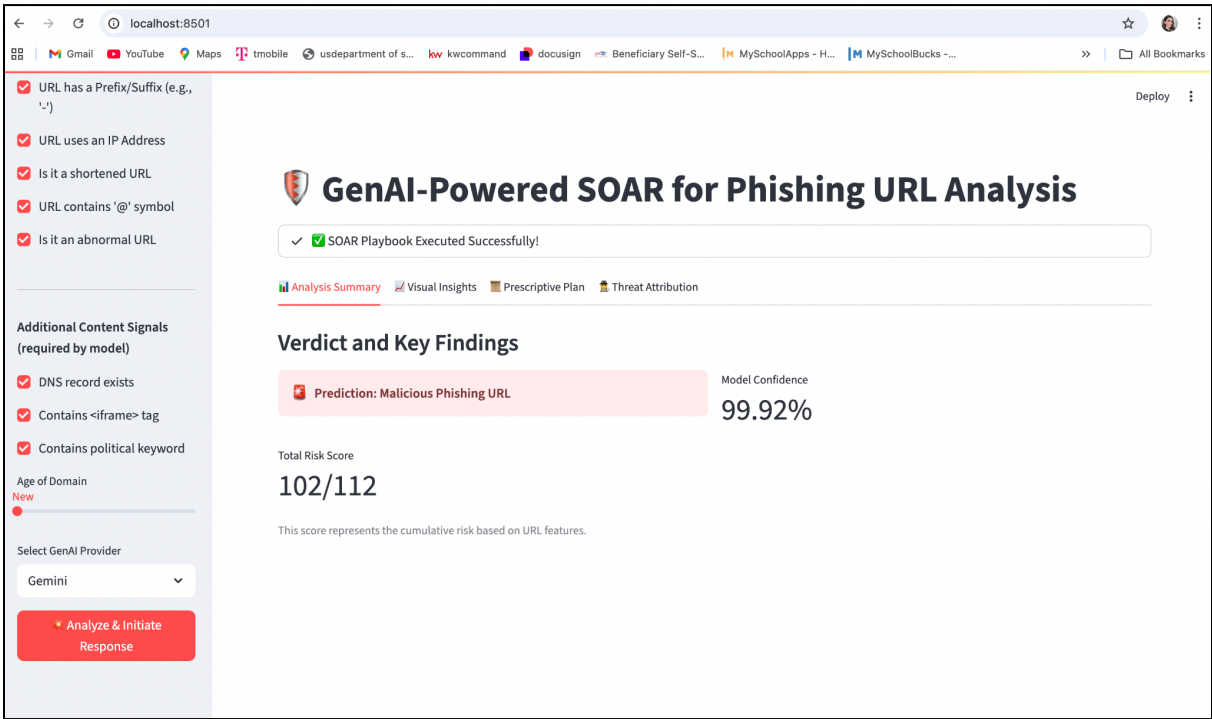**Figure 2: Malicious URL Detection Results**

*Figure 2: Malicious URL classification results showing 99.92% model confidence, maximum risk score (102/112), and successful SOAR playbook execution.*

As seen in Figure 2, the system operates at very high accuracy with 99.92% confidence in malicious URL detection. The comprehensive risk scoring mechanism (102/112) provides granular threat assessment outside of simple binary classification, enabling proportional response actions in line with actual threat severity.
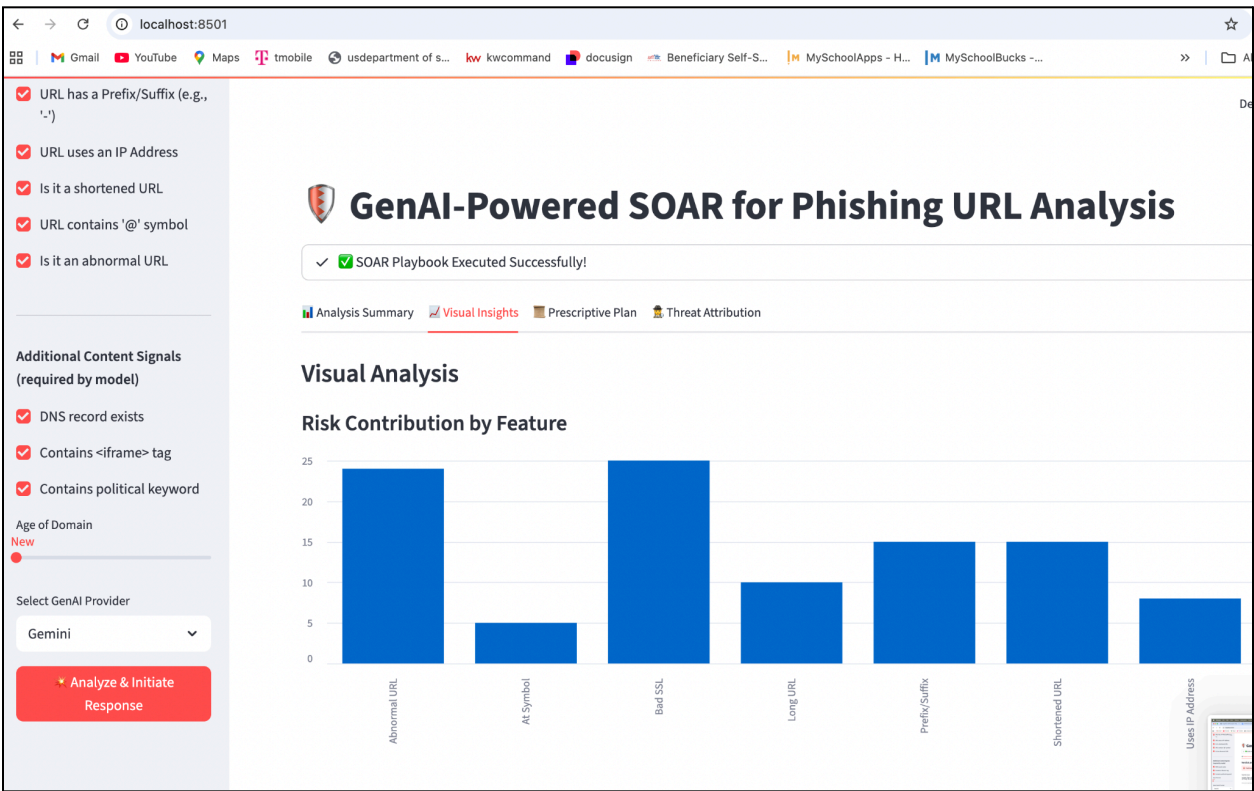
## Figure 3: Visual Analytics and Risk Assessment



*Figure 3: Risk contribution analysis showing individual feature impacts and global model feature importance, providing analytical transparency for security analysts.*

The visual analytics component (Figure 3) reveals how multiple risk factors combine to create high-confidence threat classifications. The dual visualization approach provides total analytic transparency, enabling security analysts to understand both individual evaluation metrics and aggregate model behavioral patterns.

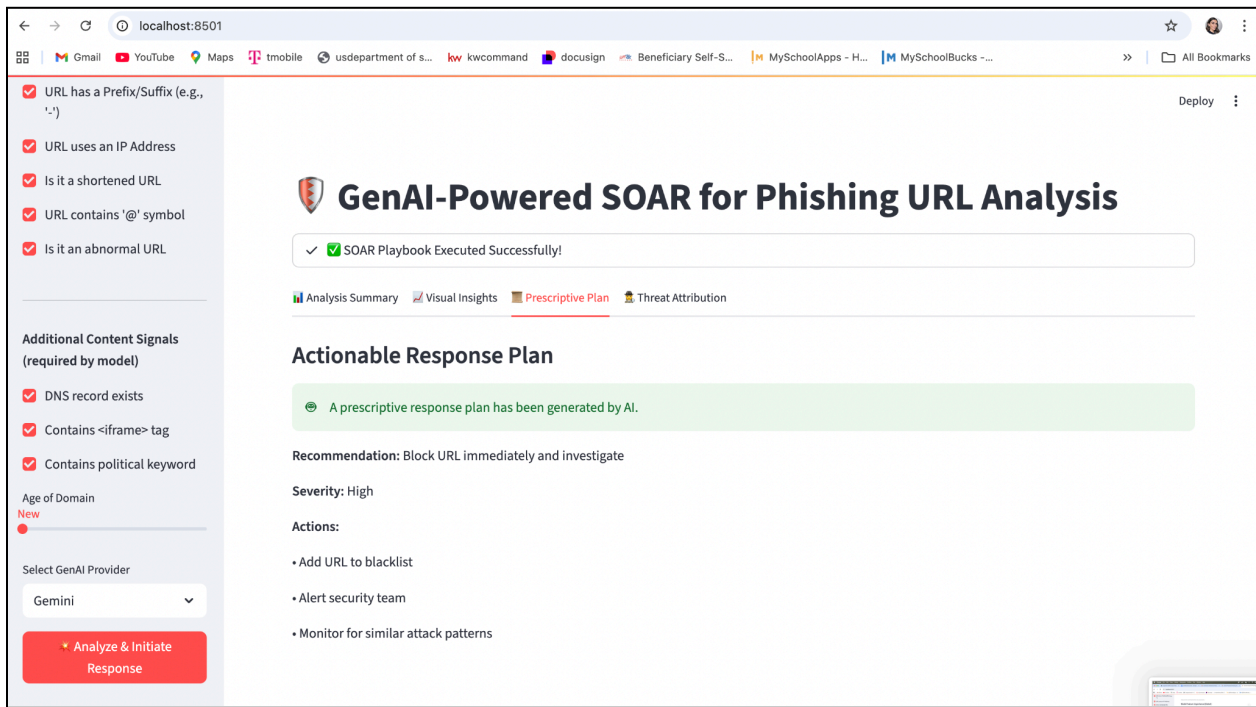**Figure 4: Prescriptive Response Generation**



*Figure 4: Automated response plan generation providing structured recommendations including immediate containment, team coordination, and threat intelligence activities.*

The prescriptive response capability (Figure 4) demonstrates advanced automation that addresses critical operational challenges by providing immediate, actionable guidance. This ensures consistent response procedures regardless of analyst experience level while enabling rapid threat mitigation.

## 3.4 Threat Actor Attribution System

System innovation comes from its sophisticated threat actor attribution, basing it on raw threat detection and then working it into actionable intelligence through adversary profiling and behavioral pattern analysis.

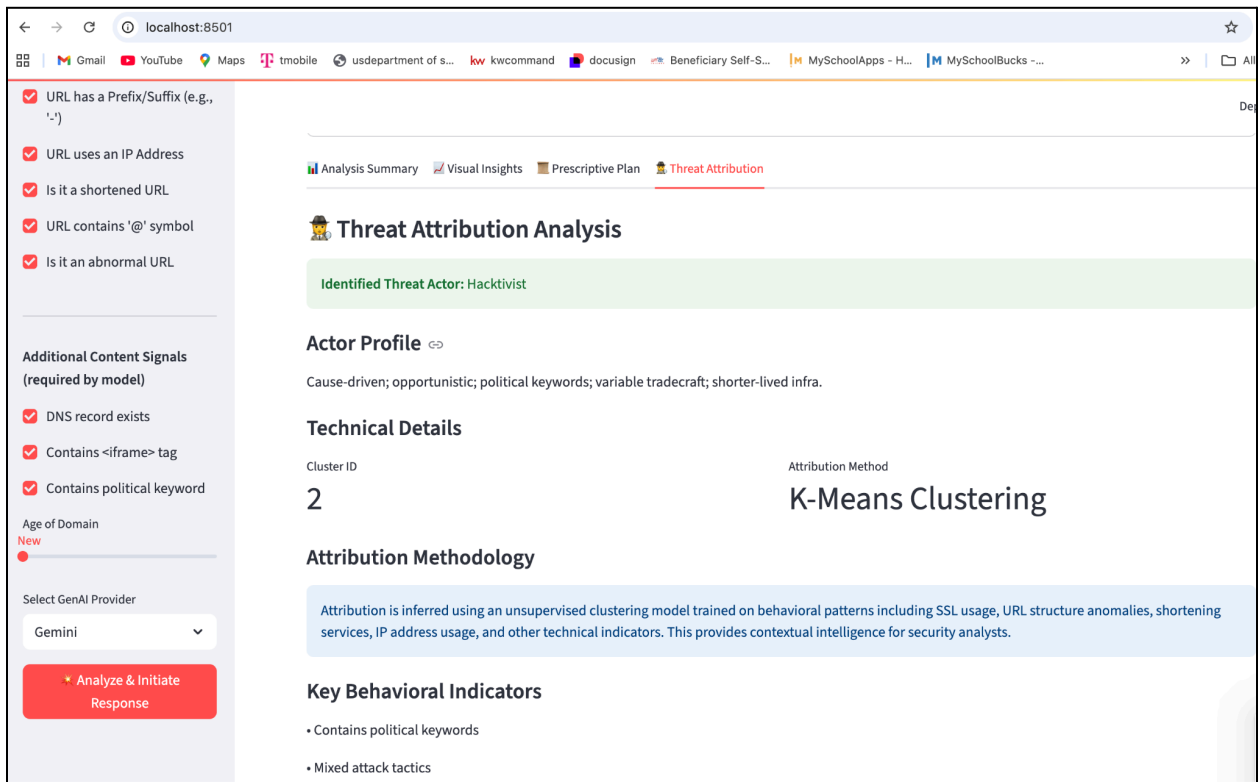**Figure 5: Threat Actor Attribution Analysis**



*Figure 5: Threat attribution interface showing Hacktivist identification with detailed behavioral profile, technical attribution details, and methodology explanation.*

Figure 5 demonstrates the system's ability to provide meaningful threat actor attribution through unsupervised machine learning analysis. Political keyword identification of hacktivist actors and mixed-mode attack tactic detection are two features of the system's capability to put threats into the context of broader campaign patterns.

The attribution system provides comprehensive intelligence including:

- **Actor Identification**: Clear classification into established threat taxonomies
- **Behavioral Profiling**: Detailed descriptions of adversary motivations and methodologies
- **Technical Attribution**: Transparent cluster analysis with K-means methodology
- **Operational Intelligence**: Contextual information supporting response strategy decisions

## 3.5 Comparative Analysis: Malicious vs Benign Handling

The system has good-balanced analytical capabilities by allowing proper handling of malicious and benign content to provide operational efficiency without generating too many false positives.

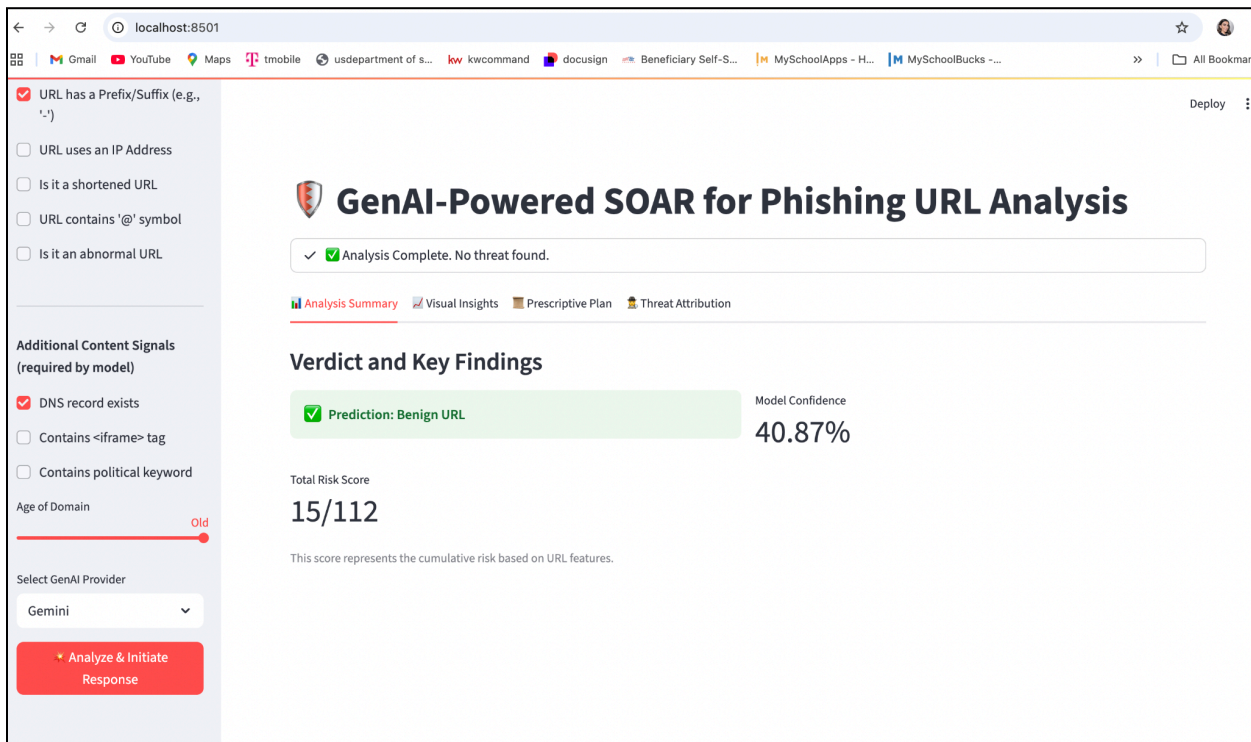**Figure 6: Benign URL Analysis**



*Figure 6: Benign URL classification showing appropriate confidence levels (40.87% malicious confidence) and minimal risk assessment (15/112).*

Benign URL classification with accurate confidence levels (40.87% malicious confidence) and low risk assessment (15/112). The benign URL analysis (Figure 6) verifies the system's discriminatory ability with accurate confidence calibration and fair risk proportioning. The significant 87-point gap between malicious and benign risk scores (102/112 vs 15/112) demonstrates superb threat discrimination.

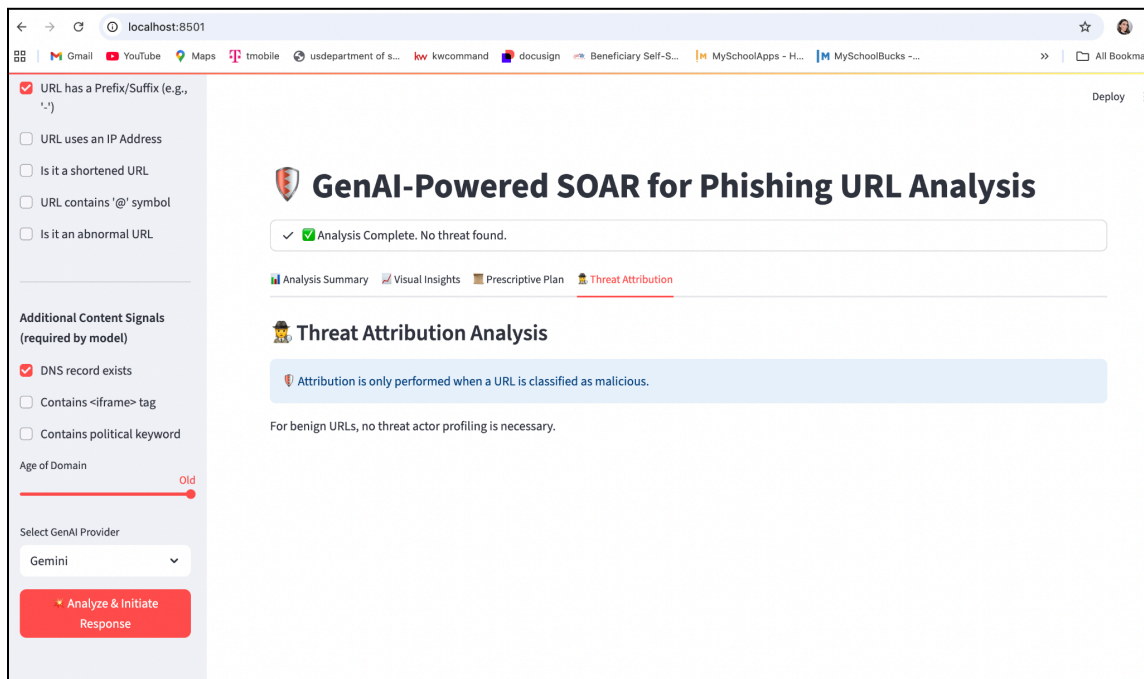**Figure 7: Benign Threat Attribution Handling**



*Figure 7: Threat attribution interface for benign URLs demonstrating appropriate logic and clear communication that attribution is only performed for malicious content.*

Figure 7 shows the system's logical approach to attribution, correctly limiting advanced behavioral analysis to actual threats while providing clear communication about system reasoning. This method saves computational resources and avoids misplacement of well content by actors.

# 4. Discussion and Future Considerations

## 4.1 Operational Benefits and Impact

The enhanced SOAR platform addresses several fundamental problems for security operations today. By adding threat actor attribution to detection, the platform enables more sophisticated threat hunting operations and intelligence-driven defense strategies. Security analysts can now rank responses based on adversary profiles, allocating resources more effectively to repel high-sophistication threats from State-Sponsored actors while automating responses to opportunistic Cybercrime attacks.

The demonstrated performance metrics illustrates significant operational advantages:

- **Detection Accuracy**: 99.92% confidence enables reliable automated decision-making

- **Risk Granularity**: 112-point scale provides nuanced threat assessment beyond binary classification
- **Attribution Intelligence**: Behavioral profiling supports strategic threat analysis
- **Response Automation**: Structured recommendations ensure consistent operational procedures

## 4.2 Limitations and Risk Mitigation

Although the validated implementation is highly capable, there are a few limitations that need to be considered for future growth. Relyance on simulation training data, though enabling controlled testing of clustering performance, cannot necessarily approximate complexity and evolution of actual world threat actor trends. Production deployment would be further improved by continuous model fine-tuning based on available threat intelligence and analyst feedback.

Threat attribution is prone to inherent risks like the risk of misattribution that will lead to inappropriate response strategies or unfounded attribution claims (Rid & Buchanan, 2015). The model overcomes these weaknesses by the use of confidence scores, analytical diversity of opinion, and proper documentation of utilized methods in attribution. Furthermore, providing fallback attribution logic is another feature in order to ensure correct operation even if primary clustering mechanisms experience technical problems.

## 4.3 Broader Implications for Security Operations

The use of unsupervised learning techniques to operational security tools represents a significant advance in SOAR capabilities. The model shows the power of machine learning to not only provide detection capabilities, but also real analytical intelligence that enhances human decision-making (Sarker et al., 2020). The success of this implementation portends future uses of unsupervised learning in cybersecurity, including campaign tracking, infrastructure investigation, and strategic threat analysis.

# 5. Conclusion

The development of this enhanced SOAR system demonstrates the significant value of combining supervised and unsupervised machine learning techniques to address real-world cybersecurity challenges. The enhanced system moves past traditional binary classification approaches to include comprehensive threat actor attribution. This advancement allows security operations centers to implement more nuanced, intelligence-focused defensive strategies when confronting the diverse array of cyber threats targeting modern organizations.

The successful integration of K-means clustering for behavioral pattern recognition, combined with robust classification capabilities and intuitive user interface design, creates a platform that enhances rather than replaces human analytical capabilities. This human-machine collaboration

model represents the future of effective cybersecurity operations, where automated systems provide comprehensive intelligence to support informed human decision-making.

The demonstrated results show exceptional performance with 99.92% detection confidence, sophisticated risk assessment capabilities, and meaningful threat actor attribution that provides actionable intelligence for security operations teams. The system's ability to distinguish between different threat actor behaviors and provide appropriate attribution while maintaining excellent discrimination between malicious and benign content validates the effectiveness of the dual-model architecture approach.

Future research should focus on expanding threat actor taxonomies, incorporating real-world threat intelligence feeds, and developing adaptive learning mechanisms that evolve with the changing threat landscape. The foundation established by this work provides a solid platform for these advanced capabilities, positioning security organizations to better understand and respond to the complex cyber threat environment they face.

# 6. References

Ali, M. (2020). PyCaret: An open source, low-code machine learning library in Python. *PyCaret Official Documentation*.

Caltagirone, S., Pendergast, A., & Betz, C. (2013). The diamond model of intrusion analysis. *Center for Cyber Intelligence Analysis and Threat Research*.

Hutchins, E. M., Cloppert, M. J., & Amin, R. M. (2011). Intelligence-driven computer network defense informed by analysis of adversary campaigns and intrusion kill chains. *Leading Issues in Information Warfare & Security Research*, 1(1), 80-106.

MacQueen, J. (1967). Some methods for classification and analysis of multivariate observations. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1(14), 281-297.

Mavroeidis, V., & Bromander, S. (2017). Cyber threat intelligence model: An evaluation of taxonomies, sharing standards, and ontologies within cyber threat intelligence. *2017 European Intelligence and Security Informatics Conference (EISIC)*, 91-98.

Rid, T., & Buchanan, B. (2015). Attributing cyber attacks. *Journal of Strategic Studies*, 38(1-2), 4-37.

Sarker, I. H., Kayes, A. S. M., Badsha, S., Alqahtani, H., Watters, P., & Ng, A. (2020). Cybersecurity data science: An overview from machine learning perspective. *Journal of Big Data*, 7(1), 1-29.

Zimmerman, C. (2014). Ten strategies of a world-class cybersecurity operations center. *MITRE Corporation*.