# UNIVERSITY OF ST.GALLEN

School of Management, Economics,
Law, Social Sciences, International Affairs and Computer Science

Type of Paper

## Title of Paper

Subtitle of Paper

Submitted by:
First name Last name
Matriculation number

Approved on Application by:

Professor:
Name and full title of professor

Date of Submission:

Day/Month/Year

# Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit.

# Contents

# List of Figures

# List of Tables

# List of Algorithms

# Listings

# 1  Introduction

Recent cybersecurity literature has highlighted the paradigm shift in offensive and defensive strategies driven by the adoption of artificial intelligence, particularly the widespread availability of large language models (LLMs) This technological evolution has fundamentally altered the threat landscape, with both attackers and defenders leveraging AI capabilities to enhance their operations. The constantly evolving landscape of cybersecurity demands more sophisticated and adaptive defense strategies, as traditional static defense mechanisms prove increasingly insufficient against advanced persistent threats and zero-day vulnerabilities **gizzarelli2024**.

The integration of Artificial Intelligence (AI) and Machine Learning (ML) represents one of the most significant emerging trends in cybersecurity, revolutionizing threat detection and response by enabling systems to analyze vast amounts of data at unprecedented speeds, identify patterns, and predict potential threats before they materialize. Moreover, the advent of generative AI, a subset of AI focused on creating new content or data, presents both opportunities and challenges in the cybersecurity landscape, opening new possibilities for both attackers and defenders **gizzarelli2024**.

Honeypot technology has evolved similarly in response to these changes. AI-enhanced honeypots now offer improved credibility against sophisticated attackers **christli2024**, enhanced threat actor analysis capabilities through advanced behavioral monitoring **Otal2024**, and the ability to detect AI-powered attack vectors including autonomous LLM-based hacking agents **reworr2024**. Recent research has demonstrated that AI-driven honeypots powered by Large Language Models can dynamically generate contextually appropriate, human-like responses in real-time, greatly improving their ability to deceive and engage attackers **christli2024**.

However, while traditional honeypots suffer from limited interactivity and predictable behavior patterns that enable easy detection, existing research on AI-enhanced honeypots has yet to adequately address practical implementation frameworks, particularly for HTTP-based attack scenarios. The LLM in the Shell (SheLLM) project has shown promising results in creating generative honeypots for command-line interfaces **sladic2023**, but comprehensive frameworks for web-based protocols remain underexplored.

This research addresses the gap between theoretical AI honeypot concepts and production-ready systems by developing and evaluating a deployable AI/LLM-enhanced honeypot specifically designed for HTTP protocol interactions. The primary research question guiding this study is: How can large language models be effectively integrated into HTTP honeypots to create more convincing, interactive decoy systems while maintaining operational feasibility for production environments?

The related work section examines various approaches to AI honeypot implementations, including model training methodologies for cybersecurity applications **gizzarelli2024**, the integration of model-context protocols (MCP) in honeypot architectures, and threat mitigation strategies. Additionally, it explores containerized deployment approaches using Docker and analyzes how open-source developments in this field diverge from current academic research efforts.

The methodology section details the tools and frameworks employed in developing the honeypot system, establishing baseline requirements for cost-effectiveness, infrastructure needs, and performance expectations. Particular attention is given to the trade-offs between AI model sophistication and operational constraints, drawing insights from recent advances in test-time scaling for language models **muennighoff2025**.

The results section presents implementation challenges and key discoveries, while addressing critical ethical and technical considerations inherent in honeypot deployment. This includes discussion of legal compliance, data privacy concerns, and responsible disclosure practices, particularly important given the sophisticated deception capabilities of AI-enhanced systems **spitzner2003**.

The ability to quickly and accurately interpret data from security systems is paramount in the ever-evolving cybersecurity landscape. To enhance the effectiveness of honeypots, this research incorporates automatic mapping of collected logs to the MITRE ATT&CK framework, a comprehensive and widely recognized knowledge base of adversary tactics and techniques **gizzarelli2024**. By making this mapping process dynamic, security teams can immediately contextualize the activities observed in honeypot logs, identifying specific attack patterns and understanding the broader strategy behind intrusion attempts.

This thesis contributes to the field by implementing a production-ready AI/LLM honeypot for HTTP protocol interactions, capable of mimicking authentic web services, intelligently analyzing incoming requests, generating contextually appropriate responses, and maintaining response consistency through intelligent caching mechanisms. The convergence of AI-driven dynamic honeypots with automatic log mapping to frameworks like MITRE ATT&CK represents a new era of evolved cybersecurity, where complex systems can grow by integrating generative AI into honeypots, creating interactions that challenge even the most experienced attackers **gizzarelli2024**.

# 2 Literature Review

This section reviews key academical papers around the subject of AI/LLM-based Honeypots allowing a better understanding of the current field and its possible future developpement.

## 2.1 Traditional Honeypot Technologies and Limitations

Honeypots are computer systems designed to capture unauthorized activity by emulating operating systems, applications or services, serving as decoy systems whose primary function is to be attacked or compromised **spitzner2003**. The foundational work by Spitzner established honeypots as critical tools for catching insider threats and understanding attacker methodologies **spitzner2003**. They are broadly categorized by their level of interaction, which determines how well they emulate a target device and how easily an attacker can identify them as a honeypot **dodson2022**.

Low-interaction honeypots typically emulate only a limited set of services, providing minimal interaction while making them easy to set up and configure for gathering basic information about known attack patterns with reduced risk **Ng2021**. The comprehensive framework analysis by Ng et al. demonstrates that these systems excel in automated attack detection but struggle with sophisticated human-operated threats **Ng2021**. Examples include Honeyd, which allows users to create virtual hosts simulating various operating systems and services to observe scanning and unauthorized access attempts, representing one of the most widely deployed low-interaction solutions **Ng2021**.

In contrast, high-interaction honeypots mimic entire systems with real operating systems and services, providing extensive interaction capabilities to capture comprehensive data and discover new types of attacks and malware **Ng2021**. Industrial Control Systems (ICS) honeypot networks, as demonstrated by Dodson et al., showcase the effectiveness of high-interaction systems in detecting targeted attacks on critical infrastructure, revealing sophisticated attack patterns that would be missed by simpler solutions **dodson2022**. However, these systems are notably more complex to deploy and manage, requiring careful planning for monitoring, logging, firewalls, and intrusion detection systems to prevent attacker escape into production environments **dodson2022**.

Adaptive honeypot configuration strategies have emerged as a response to traditional limitations, with Fraunholz et al. proposing dynamic deployment and maintenance approaches that can adjust to evolving threat landscapes **frauenholz2017**. Their research demonstrates that adaptive configurations can significantly improve detection rates while reducing maintenance overhead compared to static deployments **frauenholz2017**.

Performance analysis of traditional honeypot systems has demonstrated both their effectiveness and limitations in real-world deployments. Fuzi's comprehensive analysis of T-Pot honeypots revealed significant capabilities in network intrusion detection, providing valuable insights into attack patterns and threat behaviors while highlighting the challenges of processing large volumes of honeypot data **fuzi2024**. The study showed that traditional honeypots excel at capturing automated attacks but struggle with sophisticated, human-operated campaigns that can adapt their behavior based on system responses **fuzi2024**.

The taxonomy for dynamic honeypot measures of effectiveness established by Pittman et al. provides a structured framework for evaluating honeypot performance, identifying key limitations in traditional approaches **pittman2020taxonomy**. Their research reveals that static honeypots lack the ability to adapt to evolving attack patterns and may become obsolete as attackers develop new techniques for honeypot detection **pittman2020taxonomy**. Limited interaction capabilities in low-interaction systems consistently fail to capture sophisticated multi-stage attacks that require deeper system engagement and sustained attacker interaction.

Detection by experienced attackers remains a significant concern, as traditional honeypots often exhibit predictable behavior patterns that can reveal their deceptive nature to skilled adversaries. The scalability challenges of high-interaction honeypots make them resource-intensive to deploy and maintain across large networks, while manual configuration and maintenance requirements create operational overhead that limits their practical deployment in many organizations **pittman2020taxonomy**. These fundamental limitations have driven the cybersecurity community toward developing more sophisticated, AI-enhanced honeypot systems that can address many of these traditional shortcomings through adaptive behavior and intelligent response generation.

## 2.2 AI and Machine Learning in Cybersecurity

The integration of Artificial Intelligence (AI) and Machine Learning (ML) has revolutionized cybersecurity by enabling systems to analyze vast amounts of data at unprecedented speeds, identify patterns, and predict potential threats before they materialize **gizzarelli2024**. AI-driven tools can automatically detect anomalies, reduce false positives, and prioritize the most critical threats, thereby enhancing the efficiency and effectiveness of cybersecurity efforts.

Machine learning applications in threat detection and analysis have demonstrated significant improvements over traditional signature-based approaches. Supervised learning algorithms trained on labeled datasets can identify known attack patterns with high accuracy, while unsupervised learning techniques excel at detecting previously unknown threats by identifying deviations from normal behavior patterns. Deep learning models, particularly neural networks, have shown remarkable success in analyzing complex data structures such as network traffic, system logs, and malware binaries.

Behavioral analysis using ML techniques has emerged as a particularly powerful approach for identifying sophisticated attacks that evade traditional detection methods. These systems analyze user behavior, network patterns, and system activities to establish baseline behaviors and detect anomalous activities that may indicate compromise. Machine learning models can adapt to new attack patterns by continuously

learning from observed behaviors, making them more resilient against evolving threats.

Adversarial AI in cybersecurity contexts presents both opportunities and challenges. While AI systems can be used to enhance defensive capabilities, they can also be exploited by attackers to develop more sophisticated attack methods. AI-powered attack generation techniques can create adaptive malware, generate convincing phishing content, and automate reconnaissance activities. This arms race between AI-enhanced attacks and defenses necessitates continuous advancement in defensive AI technologies.

The application of AI in cybersecurity extends to automated incident response, where machine learning algorithms can analyze security events, correlate threat indicators, and recommend or automatically execute response actions. Natural Language Processing (NLP) techniques enable the automated analysis of threat intelligence reports, extracting relevant indicators of compromise (IoCs) and mapping them to established threat frameworks such as MITRE ATT&CK.

## 2.3 AI-Enhanced Honeypot Systems

The convergence of artificial intelligence and honeypot technology represents a significant advancement in cybersecurity defense strategies. Early AI honeypot implementations focused primarily on automated response generation and basic behavioral adaptation. However, recent developments have leveraged advanced machine learning techniques, particularly large language models (LLMs), to create more sophisticated and convincing deception environments.

Machine learning for honeypot data analysis has transformed how organizations process and interpret the vast amounts of data collected by honeypot systems. Traditional honeypots generate extensive logs that require manual analysis to extract meaningful insights. AI-enhanced systems can automatically analyze this data, identify attack patterns, classify threats, and generate actionable intelligence for security teams.

Dynamic honeypot adaptation using AI enables systems to modify their behavior in real-time based on attacker actions and emerging threat patterns. These adaptive systems can change their apparent vulnerabilities, modify response patterns, and even alter their simulated operating environment to maintain attacker engagement and gather more comprehensive intelligence **Otal2024**.

Comparison of different AI approaches reveals distinct advantages and limitations. Rule-based systems offer predictable behavior and easier debugging but lack the flexibility to handle novel attack scenarios. Machine learning approaches provide better adaptability and can learn from new attack patterns but may require extensive training data and can be more difficult to interpret. Deep learning models, particularly transformer-based architectures, offer the most sophisticated response generation capabilities but require significant computational resources.

Recent research has demonstrated the effectiveness of LLM-powered honeypots in creating convincing interactive environments. The SheLLM project represents a pioneering effort in this field, leveraging large language models to simulate realistic command-line interfaces that can engage attackers in extended interactions while maintaining believable system behaviors **sladic2023**. These systems can generate contextually appropriate responses to complex commands, maintain session state across multiple inter-

actions, and adapt their behavior based on attacker actions.

Advanced LLM honeypot implementations have shown significant improvements in attacker engagement and intelligence gathering capabilities **Otal2024**; **reworr2024**. These systems demonstrate the potential for large language models to create sophisticated interactive deception environments that can maintain prolonged attacker engagement while collecting comprehensive behavioral data.

The integration of generative AI technologies has enabled honeypots to simulate various operating systems, applications, and services with unprecedented realism. These systems can generate dynamic content, respond to novel queries, and maintain consistent personalities across extended engagements. The ability to process natural language inputs and generate human-like responses makes AI-enhanced honeypots particularly effective against social engineering attacks and human-operated threats **christli2024**.

## 2.4   Model Training for Cybersecurity Applications

Training data requirements and datasets for security models present unique challenges in the cybersecurity domain. Unlike many other machine learning applications, cybersecurity models require datasets that accurately represent real-world attack patterns while maintaining ethical and legal boundaries. The sensitive nature of cybersecurity data often limits the availability of comprehensive training datasets, necessitating innovative approaches to data collection and sharing.

Domain-specific model fine-tuning approaches have emerged as essential techniques for adapting general-purpose AI models to cybersecurity applications. Pre-trained language models can be fine-tuned on cybersecurity-specific datasets to improve their understanding of technical terminology, attack methodologies, and appropriate response patterns. This approach leverages the broad knowledge encoded in large-scale pre-trained models while specializing them for cybersecurity tasks.

Transfer learning in cybersecurity contexts enables the application of knowledge gained from one security domain to another. Models trained on network intrusion data can be adapted for malware analysis, and techniques developed for one type of honeypot can be transferred to different deployment scenarios. This approach is particularly valuable given the limited availability of labeled cybersecurity datasets.

Model training for deception and social engineering scenarios requires careful consideration of ethical implications and potential misuse. Training data must include realistic but appropriately sanitized examples of social engineering attacks, phishing attempts, and other human-targeted threats. The challenge lies in creating models that can effectively simulate these interactions for defensive purposes without enabling malicious applications.

Evaluation metrics for security-focused AI models must account for the unique requirements of cybersecurity applications. Traditional machine learning metrics such as accuracy and precision remain important, but security-specific metrics such as false positive rates, detection time, and attacker engagement duration become critical measures of effectiveness. The development of standardized evaluation frameworks for AI-enhanced cybersecurity systems remains an active area of research.

Challenges in creating realistic training datasets include the need to balance realism with privacy and

security concerns. Synthetic data generation techniques offer promising solutions by creating artificial datasets that maintain the statistical properties of real cybersecurity data without exposing sensitive information. Advanced generative models can create realistic network traffic patterns, system logs, and attack scenarios for training purposes.

Few-shot and zero-shot learning applications are particularly relevant in cybersecurity, where new attack patterns emerge continuously. These techniques enable AI models to adapt to new threats with minimal training data, leveraging their existing knowledge to understand and respond to novel attack scenarios. This capability is crucial for maintaining effectiveness against rapidly evolving cyber threats **muennighoff2025**.

Model training innovations such as simple test-time scaling have shown promise in improving the performance of AI models in complex reasoning tasks, which has direct applications in cybersecurity scenarios where models must analyze and respond to sophisticated attack patterns in real-time **muennighoff2025**.

## 2.5 HTTP Protocol Security and Web-Based Attacks

Common HTTP-based attack vectors continue to represent a significant portion of cybersecurity threats, making web-focused honeypots essential components of modern defense strategies. Cross-site scripting (XSS), SQL injection, cross-site request forgery (CSRF), and various injection attacks exploit vulnerabilities in web applications and services. Understanding these attack patterns is crucial for developing effective honeypot systems that can attract and analyze web-based threats.

Web application security challenges stem from the complexity of modern web technologies and the diverse attack surface they present. The interaction between client-side and server-side components, the use of multiple programming languages and frameworks, and the integration of third-party services create numerous potential vulnerabilities. Honeypots designed to simulate these complex environments must accurately replicate these vulnerabilities while maintaining security boundaries.

HTTP honeypot implementations require careful consideration of protocol compliance and realistic behavior simulation. These systems must handle various HTTP methods, status codes, headers, and content types while maintaining believable responses to both legitimate and malicious requests. The challenge lies in creating systems that appear vulnerable to automated scanners while providing meaningful interaction for human attackers.

API security and monitoring represent increasingly important aspects of web-based cybersecurity. Modern applications rely heavily on REST APIs, GraphQL endpoints, and other web services that present unique attack vectors. Honeypot systems must simulate these interfaces convincingly, handling authentication mechanisms, parameter validation, and response formatting while collecting intelligence about API-focused attacks.

The integration of AI technologies into web-based honeypots enables more sophisticated simulation of web application behaviors. Machine learning models can generate dynamic content, simulate database interactions, and respond to complex query patterns. Natural language processing capabilities allow these

systems to handle text-based inputs in forms, comments, and user-generated content areas, providing more realistic interaction surfaces for attackers.

## 2.6  Containerization and Deployment Frameworks

Docker in cybersecurity applications has revolutionized the deployment and management of security tools, including honeypot systems. Containerization provides isolation, scalability, and reproducibility benefits that are particularly valuable for honeypot deployments. Docker containers can encapsulate entire honeypot environments, making them easily deployable across different infrastructure platforms while maintaining consistent behavior.

Container security considerations are crucial when deploying honeypot systems, as the goal of attracting attackers creates unique security challenges. Proper container isolation, resource limitations, and network segmentation become critical for preventing attackers from escaping the honeypot environment and accessing production systems. Security best practices for containerized honeypots include minimal base images, read-only file systems where possible, and comprehensive monitoring of container activities.

Scalable honeypot deployments benefit significantly from containerization technologies. Container orchestration platforms such as Kubernetes enable the automated deployment, scaling, and management of honeypot systems across large infrastructure environments. These platforms can dynamically create and destroy honeypot instances based on demand, distribute them across multiple nodes, and provide centralized logging and monitoring capabilities.

Cloud-based honeypot architectures leverage the scalability and flexibility of cloud computing platforms to create distributed honeypot networks. These deployments can span multiple geographic regions, simulate various infrastructure configurations, and scale automatically based on attack patterns and resource requirements. Cloud-native technologies such as serverless computing and managed services can reduce operational overhead while improving scalability.

The integration of AI technologies with containerized honeypot deployments enables sophisticated orchestration and management capabilities. Machine learning algorithms can analyze attack patterns and automatically adjust honeypot configurations, spin up additional instances in response to increased attack activity, and optimize resource allocation based on threat intelligence. Container platforms provide the infrastructure foundation for these AI-driven capabilities.

## 2.7  Threat Intelligence and Analysis

Automated threat intelligence gathering has become essential for maintaining situational awareness in the rapidly evolving cybersecurity landscape. AI-enhanced systems can continuously monitor threat feeds, analyze security reports, and extract relevant indicators of compromise (IoCs) from various sources. Natural language processing techniques enable the automated analysis of unstructured threat intelligence reports, converting human-readable descriptions into machine-processable data structures.

Threat actor profiling and attribution represent complex analytical challenges that benefit significantly

from machine learning approaches. AI systems can analyze attack patterns, tool usage, infrastructure preferences, and tactical approaches to identify similarities between different campaigns and potentially attribute them to specific threat actors or groups. This analysis relies on comprehensive datasets of historical attack data and sophisticated pattern recognition algorithms.

IoC (Indicators of Compromise) extraction from unstructured sources requires advanced natural language processing capabilities. AI systems must identify IP addresses, domain names, file hashes, registry keys, and other technical indicators within free-text reports while understanding their context and relevance. Machine learning models trained on cybersecurity datasets can achieve high accuracy in extracting and classifying these indicators.

Real-time threat analysis frameworks integrate multiple data sources, including honeypot logs, network telemetry, endpoint data, and external threat intelligence, to provide comprehensive threat assessment capabilities. AI algorithms can correlate events across these diverse data sources, identify complex attack patterns, and prioritize threats based on their potential impact and likelihood of success.

The application of honeypot data to threat intelligence analysis provides unique insights into attacker behavior and tactics. Unlike other security data sources that may only capture successful attacks or detection events, honeypots provide detailed logs of attacker interactions, tool usage, and tactical approaches. This data is particularly valuable for understanding the human element of cyber attacks and developing behavioral models of threat actors.

## 2.8   MITRE ATT&CK Framework

The MITRE ATT&CK (Adversarial Tactics, Techniques, and Common Knowledge) framework represents one of the most comprehensive and widely adopted approaches to understanding and categorizing cyber adversary behavior. Developed by the MITRE Corporation, this globally accessible knowledge base provides a structured approach to documenting adversary tactics and techniques based on real-world observations **gizzarelli2024**.

The framework organizes adversary behavior into a matrix structure that maps tactics (the "why" of attacks) to techniques (the "how" of attacks). This organization provides security professionals with a common language for discussing and analyzing cyber threats, facilitating improved communication and collaboration across the cybersecurity community. The framework encompasses multiple technology domains, including Enterprise, Mobile, and Industrial Control Systems (ICS), each representing different operational environments where adversaries operate.

The MITRE ATT&CK framework describes adversarial TTPs (Tactics, Techniques, and Procedures) using a structure of four increasingly granular levels: tactics, techniques, sub-techniques, and procedures. Tactics represent the high-level goals or objectives of adversaries, such as gaining credential access or establishing persistence. Techniques describe the specific methods used to achieve tactical objectives, while sub-techniques provide additional granularity for complex techniques that can be implemented in multiple ways.

Automatic mapping of cyber threat intelligence to the MITRE ATT&CK framework addresses one of the most significant challenges in applying this framework at scale. Traditional manual mapping processes are time-consuming and error-prone, particularly when dealing with large volumes of unstructured threat intelligence data. AI-powered mapping systems leverage natural language processing and machine learning techniques to automatically analyze threat reports and map relevant behaviors to appropriate ATT&CK techniques.

The integration of honeypot data with MITRE ATT&CK mapping provides unprecedented insights into attacker behavior and tactics. Systems like SYNAPSE-to-MITRE demonstrate the potential for real-time mapping of honeypot interactions to the ATT&CK framework, enabling security teams to understand attack patterns in the context of established threat models. This integration facilitates rapid threat assessment, incident response planning, and defensive strategy development.

Machine learning models for ATT&CK mapping face unique challenges related to the complexity and ambiguity of natural language descriptions of cyber threats. The SYNAPSE-to-MITRE extension addresses these challenges by training multilayer perceptron classifiers on curated datasets of cyber threat intelligence reports. The system's ability to provide multiple potential mappings for each analyzed sentence acknowledges the inherent ambiguity in threat descriptions while providing security analysts with comprehensive classification options.

## 2.9   Ethical and Legal Considerations

Legal frameworks for honeypot deployment vary significantly across jurisdictions and present complex challenges for organizations implementing deception technologies. The intentional creation of vulnerable systems designed to attract attackers raises questions about legal liability, data protection compliance, and the appropriate scope of defensive activities. Organizations must carefully consider applicable laws regarding computer fraud, unauthorized access, and data privacy when deploying honeypot systems.

Ethical guidelines for deception technologies require balancing the legitimate security benefits of honeypots against potential risks and unintended consequences. The use of deception in cybersecurity contexts raises questions about proportionality, necessity, and the potential for misuse. Professional organizations and industry groups have developed guidelines for the ethical deployment of honeypots, emphasizing the importance of proper authorization, scope limitation, and responsible data handling.

Privacy concerns and data protection represent significant considerations for honeypot deployments, particularly in jurisdictions with strict data protection regulations such as the European Union's General Data Protection Regulation (GDPR). Honeypots may collect personal information from attackers, including IP addresses, user agents, and potentially personally identifiable information entered during attacks. Organizations must implement appropriate privacy safeguards and ensure compliance with applicable data protection requirements.

Responsible disclosure practices become particularly important when honeypots discover new vulnerabilities or attack techniques. Organizations must balance the need to protect their own systems and share threat intelligence with the cybersecurity community against the potential risks of disclosing sensitive in-

formation. Established frameworks for responsible disclosure provide guidance for handling discoveries made through honeypot deployments.

The deployment of AI-enhanced honeypots introduces additional ethical considerations related to the use of artificial intelligence in security contexts. The potential for AI systems to engage in sophisticated deception raises questions about the appropriate limits of automated defensive activities. Organizations must consider the implications of using AI technologies that can learn from and adapt to attacker behavior, potentially leading to increasingly sophisticated deception capabilities.

Recent research on LLM agent honeypots has highlighted the importance of monitoring AI hacking agents in the wild, providing insights into how malicious actors might leverage AI technologies for offensive purposes **reworr2024**. This research underscores the need for ethical guidelines that address the dual-use nature of AI technologies in cybersecurity contexts.

International cooperation and information sharing through honeypot networks require careful consideration of legal and regulatory frameworks across multiple jurisdictions. The global nature of cyber threats necessitates international collaboration, but differences in legal systems, data protection requirements, and law enforcement approaches can complicate information sharing arrangements. Organizations participating in honeypot networks must navigate these complex legal landscapes while maintaining effective threat intelligence sharing capabilities.

# 3 Galah: The Foundation LLM-Enhanced Honeypot System

The Galah honeypot system represents a paradigm shift from traditional static deception technologies toward dynamic, AI-powered interactive environments. This chapter provides a comprehensive analysis of the Galah architecture, its core components, and the foundational capabilities that enable sophisticated attacker engagement through large language model integration.

## 3.1 System Architecture and Design Principles

Galah's architecture is built upon several key design principles that differentiate it from traditional honeypot implementations. The system employs a service-oriented architecture that separates concerns between HTTP request processing, LLM integration, caching mechanisms, and event logging. This modular approach enables flexible deployment configurations while maintaining consistent behavior across different operational environments.

The core service layer encapsulates all components required for response generation, implementing a provider pattern that abstracts LLM-specific implementations through a unified interface. This abstraction enables seamless switching between different LLM providers, including OpenAI's GPT models, Google's Gemini, Anthropic's Claude, and local deployments through Ollama, without requiring changes to the core honeypot logic.

The system's event-driven architecture enables real-time processing of HTTP requests while maintaining detailed audit trails of all interactions. Each request undergoes systematic analysis to extract relevant features, generate appropriate prompts for the LLM, and construct responses that maintain consistency with the simulated environment. The architecture supports both synchronous and asynchronous processing patterns, enabling optimal resource utilization under varying load conditions.

Caching strategies form a critical component of the system architecture, addressing both cost optimization and performance requirements. The port-specific caching mechanism ensures that responses generated for particular services remain consistent while preventing cross-contamination between different simulated applications. The SQLite-based caching layer provides persistent storage with minimal operational overhead, enabling cache persistence across system restarts and deployments.

## 3.2 LLM Integration and Provider Abstraction

The integration of large language models into the Galah system required careful consideration of provider-specific capabilities, API limitations, and cost optimization strategies. The system implements a comprehensive provider abstraction layer that standardizes interactions across different LLM services while preserving access to provider-specific features and configurations.

The provider abstraction supports major commercial LLM services including OpenAI's GPT-4 and GPT-3.5-turbo models, Google's Gemini Pro and Flash variants, Anthropic's Claude models, Cohere's command models, and Google Cloud's Vertex AI platform. Additionally, the system supports local deployments through Ollama integration, enabling organizations to maintain complete control over model deployment and data privacy.

Temperature control and sampling parameters are standardized across providers to ensure consistent response generation characteristics. The system implements adaptive temperature scaling based on request complexity and context requirements, enabling more creative responses for complex scenarios while maintaining deterministic behavior for standard interactions.

Model context management represents a significant technical challenge in LLM integration. The system implements sophisticated prompt engineering techniques that maximize the effective use of available context windows while ensuring that critical system instructions and behavioral guidelines remain intact throughout extended interactions. Dynamic context pruning algorithms selectively remove less relevant historical context when approaching model limits while preserving essential state information.

Error handling and failover mechanisms ensure system resilience in the face of LLM service disruptions. The system implements exponential backoff strategies for rate limiting scenarios, automatic provider switching for service outages, and graceful degradation to cached responses when LLM services are unavailable. These mechanisms ensure continuous honeypot operation even during periods of LLM service instability.

## 3.3 HTTP Request Processing and Analysis

Galah's HTTP request processing pipeline implements sophisticated analysis capabilities that extract meaningful features from incoming requests while maintaining the performance characteristics required for real-time response generation. The processing pipeline begins with comprehensive request parsing that captures HTTP methods, headers, query parameters, request bodies, and metadata relevant to behavioral analysis.

The system implements advanced payload analysis techniques that assess request complexity, identify potential attack vectors, and extract semantic meaning from request content. Natural language processing capabilities enable the system to understand human-readable content in form submissions, comments, and other user-generated inputs, facilitating more contextually appropriate responses.

Header analysis algorithms examine request headers for indicators of automated tools, browser finger-

prints, and potential evasion techniques. The system maintains extensive databases of known tool signatures, enabling identification of common security testing tools, web crawlers, and automated scanners. This analysis informs response generation strategies, allowing the system to adapt its behavior based on detected attacker characteristics.

Request correlation and session management enable the system to maintain consistent behavior across multiple interactions with the same attacker. The sessionization algorithm creates stable session identifiers based on connection characteristics while respecting privacy considerations. Session state management preserves context across requests, enabling the simulation of stateful web applications and maintaining narrative consistency in multi-step interactions.

The system implements comprehensive logging of all request processing activities, capturing detailed metrics about processing times, feature extraction results, and decision points throughout the analysis pipeline. This telemetry data enables continuous optimization of processing algorithms and provides valuable insights into system performance characteristics under different load conditions.

## 3.4   Response Generation and Consistency Management

Response generation in Galah represents the culmination of sophisticated analysis and AI-powered content creation. The system constructs detailed prompts that provide the LLM with comprehensive context about the simulated environment, the specific request being processed, and behavioral guidelines for maintaining realistic interactions.

The prompt engineering framework implements templating systems that standardize prompt structure while enabling customization for different deployment scenarios. Templates include system-level instructions that define the honeypot's personality and behavioral characteristics, context-specific information about the simulated application or service, and request-specific details that guide response generation for the particular interaction.

JSON-structured response generation ensures consistent formatting and enables reliable parsing of LLM outputs. The system enforces strict JSON schema validation to prevent malformed responses and implements fallback mechanisms for cases where LLM outputs do not conform to expected formats. This structured approach enables reliable extraction of HTTP headers, status codes, and response content while maintaining flexibility in content generation.

Consistency management algorithms ensure that responses maintain logical coherence across extended interactions. The system tracks state information about simulated environments, user accounts, application data, and system configurations, ensuring that subsequent responses remain consistent with previously established context. This state management enables the simulation of complex, stateful web applications that maintain realistic behavior patterns throughout extended attacker engagement.

The system implements sophisticated content filtering and safety mechanisms that prevent the generation of inappropriate content while maintaining the deceptive effectiveness of the honeypot. These mechanisms include automated content analysis, keyword filtering, and safety classification algorithms that

ensure generated responses remain within acceptable bounds for organizational deployment.

## 3.5 Caching Architecture and Optimization

Galah's caching architecture addresses critical performance and cost optimization requirements while maintaining response quality and consistency. The system implements a multi-layered caching strategy that operates at request, session, and system levels to maximize cache hit rates while minimizing LLM API costs.

The primary caching layer utilizes SQLite for persistent storage of generated responses, implementing sophisticated key generation algorithms that create stable cache keys based on request characteristics while accounting for relevant context variables. The cache key generation process considers HTTP methods, normalized URLs, relevant headers, and request content while filtering out ephemeral data that would prevent effective cache utilization.

Port-specific caching ensures that responses generated for different services remain isolated, preventing inappropriate response reuse across different simulated applications. This isolation is critical for maintaining the believability of the honeypot environment, as responses generated for a database administration interface should not be reused for a content management system on a different port.

Cache invalidation strategies balance freshness requirements with cost optimization objectives. The system implements time-based expiration with configurable timeout periods, usage-based invalidation for frequently accessed entries, and manual invalidation capabilities for administrative control. Advanced cache warming algorithms can pre-generate responses for common request patterns, improving response times for initial attacker interactions.

The system provides comprehensive cache analytics that track hit rates, response time improvements, cost savings, and cache efficiency metrics. These analytics enable administrators to optimize cache configurations for their specific deployment environments and attack patterns. Cache performance monitoring alerts administrators to cache misses that may indicate new attack patterns or system configuration changes requiring attention.

## 3.6 Event Logging and Audit Capabilities

Comprehensive event logging forms a foundational capability of the Galah system, enabling detailed analysis of attacker behavior while providing audit trails for security analysis and compliance requirements. The logging system captures complete interaction records that include request details, processing decisions, response generation processes, and timing information.

The structured logging format utilizes JSON encoding to ensure machine-readable logs that facilitate automated analysis and integration with security information and event management (SIEM) systems. Each log entry includes standardized fields for temporal information, request characteristics, response details, processing metadata, and system performance metrics.

The system implements privacy-preserving logging techniques that capture necessary information for security analysis while respecting privacy considerations. IP address anonymization, request sanitization, and selective field redaction ensure that logs contain valuable security intelligence without exposing unnecessary personal information.

Log rotation and retention policies ensure that log files remain manageable while preserving historical data for trend analysis and forensic investigation. The system supports configurable retention periods, automatic compression of archived logs, and secure deletion of expired log data. Integration with external logging systems enables centralized log management for large-scale deployments.

Real-time log streaming capabilities enable integration with security orchestration and automated response (SOAR) systems, facilitating immediate response to critical security events. The streaming interface supports multiple output formats and delivery mechanisms, including syslog, JSON over HTTP, and message queue integration.

## 3.7   Integration with Security Infrastructure

Galah's design prioritizes seamless integration with existing security infrastructure components, enabling organizations to incorporate the honeypot system into their broader security architecture without requiring significant changes to established processes and tools.

The system implements standardized interfaces for security information sharing, including support for structured threat intelligence formats such as STIX (Structured Threat Information eXpression) and indicators of compromise (IoC) extraction. These capabilities enable automatic sharing of threat intelligence derived from honeypot interactions with threat intelligence platforms and security vendors.

API endpoints provide programmatic access to honeypot data and configuration, enabling integration with security orchestration platforms and custom security tools. The RESTful API design follows OpenAPI specifications and includes comprehensive authentication and authorization mechanisms to ensure secure access to sensitive honeypot data.

The system supports deployment in containerized environments through comprehensive Docker support, enabling integration with container orchestration platforms such as Kubernetes. Container-based deployment facilitates scalable honeypot networks and simplifies integration with cloud-native security architectures.

Network integration capabilities enable deployment in various network architectures, including DMZ networks, internal network segments, and cloud environments. The system includes support for network address translation (NAT), proxy deployments, and multi-interface configurations to accommodate different network topology requirements.

## 3.8 Performance Characteristics and Scalability

Performance analysis of the Galah system reveals characteristics that enable effective deployment in production environments while maintaining acceptable response times and resource utilization. The system's performance profile balances the computational requirements of LLM integration with the responsiveness requirements of convincing honeypot operation.

Response time analysis demonstrates that cached responses achieve sub-millisecond response times, while uncached responses requiring LLM generation typically complete within 2-10 seconds depending on model selection and provider performance. These response times remain within acceptable bounds for most web application scenarios while providing sufficient interaction time to maintain attacker engagement.

Memory utilization patterns show that the system maintains stable memory consumption under normal operating conditions, with periodic spikes during LLM response generation and caching operations. The SQLite caching layer demonstrates excellent memory efficiency, with cache databases typically consuming less than 100MB for systems handling thousands of unique requests.

CPU utilization remains minimal during cache hit scenarios, with the majority of computational load occurring during LLM API interactions and JSON processing. The system's event-driven architecture enables efficient handling of concurrent requests while maintaining resource isolation between different processing activities.

Scalability testing reveals that single-instance deployments can effectively handle hundreds of concurrent attackers while maintaining response quality and consistency. Horizontal scaling through container orchestration enables support for thousands of concurrent sessions across distributed deployments, with linear scaling characteristics for most operational scenarios.

## 3.9 Security Considerations and Hardening

The deployment of honeypot systems requires careful attention to security considerations that prevent attackers from escaping the deception environment and accessing production systems. Galah implements comprehensive security hardening measures that address both technical and operational security requirements.

Container isolation mechanisms ensure that honeypot processes remain contained within their designated execution environments. The system utilizes Docker security features including resource limitations, read-only filesystems where appropriate, and network isolation to prevent lateral movement from honeypot containers to host systems.

Input validation and sanitization algorithms process all incoming data to prevent injection attacks against the honeypot system itself. While the system is designed to simulate vulnerable applications, the underlying honeypot infrastructure implements robust security controls to prevent actual compromise of the honeypot system.

Network segmentation requirements ensure that honeypot deployments remain isolated from production networks while maintaining necessary connectivity for management and data collection. The system supports deployment in dedicated network segments with carefully controlled ingress and egress rules.

Data protection mechanisms ensure that sensitive information captured during honeypot interactions remains secure throughout collection, processing, and storage. Encryption at rest and in transit protects captured data, while access controls ensure that only authorized personnel can access honeypot intelligence.

The system implements comprehensive security monitoring that tracks both honeypot interactions and the security posture of the honeypot system itself. Security telemetry includes authentication events, system resource utilization, network connection patterns, and indicators of potential compromise of the honeypot infrastructure.

# 4 AI/LLM Extensions to Galah: Advanced Intelligence and Adaptive Response Systems

This chapter presents the comprehensive AI/LLM extensions developed for the Galah honeypot system, transforming it from a dynamic response generator into a sophisticated threat intelligence and adaptive defense platform. The extensions introduce four major capability enhancements: enhanced research data collection, MITRE ATT&CK framework integration, behavioral analysis and attacker profiling, and context-aware response generation with machine learning adaptation.

## 4.1 Enhanced Research Data Collection Framework

The enhanced research data collection framework represents a fundamental advancement in honeypot intelligence gathering capabilities, transforming raw HTTP interactions into structured, analyzable data suitable for academic research and operational threat intelligence. This framework implements sophisticated attack vector detection algorithms, payload analysis techniques, and comprehensive session tracking mechanisms that provide unprecedented insights into attacker behavior and tactics.

### 4.1.1 Advanced Attack Vector Detection

The attack vector detection system implements a comprehensive taxonomy of web-based threats, utilizing pattern recognition algorithms and heuristic analysis to identify attack patterns in real-time. The detection engine recognizes multiple attack categories including SQL injection variants, cross-site scripting (XSS) attacks, directory traversal attempts, command injection techniques, XML external entity (XXE) attacks, server-side request forgery (SSRF), and deserialization vulnerabilities.

SQL injection detection utilizes sophisticated pattern matching algorithms that recognize both classic and advanced injection techniques. The system identifies union-based injections, boolean-based blind injections, time-based blind injections, and error-based injection attempts through comprehensive regex patterns and semantic analysis. The detection algorithms account for encoding variations, comment injection evasion techniques, and case manipulation attempts commonly used to bypass security controls.

Cross-site scripting detection encompasses reflected, stored, and DOM-based XSS variants through content analysis that examines request parameters, headers, and body content for malicious script injection attempts. The system recognizes both traditional script tag injections and advanced techniques utilizing event handlers, JavaScript protocol handlers, and CSS-based injection methods.

Command injection detection algorithms analyze request content for operating system command patterns, identifying attempts to execute system commands through web application vulnerabilities. The system recognizes command chaining techniques, encoded command injection, and platform-specific command variants across Windows, Linux, and Unix systems.

## 4.1.2 Sophisticated Payload Analysis

The payload analysis subsystem implements advanced mathematical and linguistic analysis techniques to quantify attack sophistication and extract meaningful characteristics from malicious requests. The system calculates Shannon entropy for payload randomness assessment, implementing sliding window analysis to identify encrypted or encoded content segments within larger payloads.

Complexity scoring algorithms evaluate payload sophistication through multiple dimensions including syntactic complexity, semantic diversity, encoding techniques utilized, and evasion mechanisms employed. The scoring system considers factors such as multi-stage payload construction, polymorphic encoding techniques, and anti-analysis measures embedded within attack payloads.

Pattern recognition capabilities identify common attack frameworks, automated tool signatures, and custom exploitation techniques through comprehensive signature databases and machine learning classification algorithms. The system maintains extensive databases of tool-specific patterns, enabling identification of popular security testing frameworks, automated vulnerability scanners, and custom exploitation tools.

Encoding detection algorithms identify and analyze various encoding schemes employed in attack payloads, including URL encoding, HTML entity encoding, Base64 encoding, hexadecimal encoding, and Unicode escape sequences. The system implements recursive decoding capabilities that can process multiple encoding layers while maintaining awareness of encoding-based evasion techniques.

## 4.1.3 Comprehensive Session Tracking and Timeline Analysis

Session management capabilities enable longitudinal analysis of attacker behavior through sophisticated session identification and state tracking mechanisms. The system creates stable session identifiers that persist across multiple requests while respecting privacy considerations and avoiding personally identifiable information exposure.

Timeline analysis algorithms construct detailed attack progression narratives that document the evolution of attack techniques throughout extended engagement periods. The system tracks reconnaissance activities, vulnerability probing attempts, exploitation phases, and post-exploitation activities to create comprehensive attack timelines suitable for forensic analysis.

Behavioral consistency tracking identifies patterns in attacker behavior including request timing patterns, tool switching behavior, skill progression indicators, and tactical adaptations. The system maintains historical context across sessions to identify returning attackers and correlate activities across multiple engagement periods.

## 4.2   MITRE ATT&CK Framework Integration

The integration of the MITRE ATT&CK framework into the Galah system provides automatic classification of observed attacker behaviors according to established adversarial tactics and techniques. This integration enables real-time mapping of honeypot interactions to the globally recognized framework, facilitating immediate contextualization of attack activities and enabling integration with existing threat intelligence and incident response processes.

### 4.2.1   Automated Technique Classification

The MITRE classification engine implements comprehensive mapping algorithms that analyze HTTP requests and automatically identify relevant ATT&CK techniques based on observed behaviors. The system maintains an extensive database of technique-to-HTTP-pattern mappings that cover the complete range of web-based attack vectors within the ATT&CK framework scope.

Initial Access techniques (TA0001) are identified through analysis of exploitation attempts targeting public-facing applications, including technique T1190 (Exploit Public-Facing Application) classification based on detected vulnerability exploitation attempts. The system recognizes various initial access vectors including web application exploitation, service-specific attacks, and protocol-level vulnerabilities.

Execution techniques (TA0002) are classified through analysis of command injection attempts, server-side code execution attacks, and script injection behaviors. The system identifies technique T1059 (Command and Scripting Interpreter) variants including PowerShell execution, command prompt usage, and Unix shell command execution attempts.

Persistence techniques (TA0003) are detected through analysis of web shell deployment attempts, configuration file modifications, and account manipulation activities. The system recognizes technique T1505 (Server Software Component) activities including web shell installation and malicious plugin deployment.

Defense Evasion techniques (TA0005) are identified through analysis of obfuscation attempts, encoding techniques, and anti-detection measures. The system classifies technique T1027 (Obfuscated Files or Information) activities including payload encoding, string obfuscation, and anti-analysis techniques.

### 4.2.2   Confidence Scoring and Evidence Collection

The classification system implements sophisticated confidence scoring algorithms that assess the reliability of technique identifications based on the strength of observed evidence and the specificity of detected patterns. Confidence scores range from 0.0 to 1.0, with higher scores indicating stronger evidence for specific technique classifications.

Evidence collection mechanisms capture detailed information supporting each technique classification, including specific request patterns that triggered the classification, relevant payload segments, header

information, and contextual factors that contributed to the identification. This evidence provides auditable justification for automated classifications and enables manual review of system decisions.

The system implements multi-factor confidence assessment that considers pattern specificity, false positive likelihood, contextual relevance, and historical accuracy metrics. Pattern specificity measures evaluate how uniquely a detected pattern corresponds to a specific technique, with more specific patterns receiving higher confidence scores.

### 4.2.3   Attack Campaign Tracking and Threat Actor Identification

Campaign tracking capabilities enable identification of coordinated attack activities through correlation of techniques, timing patterns, infrastructure indicators, and behavioral characteristics. The system groups related activities into campaign clusters that represent sustained attack efforts by individual actors or groups.

Threat actor profiling algorithms analyze attack patterns to identify potential threat actor characteristics including skill levels, tool preferences, target preferences, and operational patterns. The system compares observed behaviors against known threat actor profiles to suggest potential attribution hypotheses.

The system maintains comprehensive attack timeline reconstructions that document the progression of techniques throughout campaign lifecycles. Timeline analysis enables identification of attack pattern evolution, tool switching behaviors, and tactical adaptations that provide insights into threat actor capabilities and intentions.

## 4.3   Behavioral Analysis and Attacker Profiling System

The behavioral analysis and attacker profiling system represents a significant advancement in honeypot intelligence capabilities, implementing sophisticated machine learning algorithms and statistical analysis techniques to create detailed psychological and technical profiles of attackers. This system enables differentiation between various attacker types, skill level assessment, and prediction of future attack behaviors.

### 4.3.1   Sophisticated Attacker Classification

The attacker classification system utilizes machine learning algorithms trained on comprehensive datasets of attacker behaviors to categorize attackers into distinct types based on observable characteristics. The classification taxonomy includes professional penetration testers, advanced persistent threat (APT) actors, script kiddies, automated scanners, opportunistic attackers, and insider threats.

Professional penetration tester identification relies on behavioral patterns indicating systematic methodology, comprehensive tool usage, and measured approach to vulnerability discovery. The system recognizes patterns associated with professional security testing including methodical reconnaissance, targeted vulnerability assessment, and controlled exploitation attempts.

APT actor identification focuses on indicators of sophisticated, persistent attack campaigns including advanced evasion techniques, custom tool usage, multi-stage attack progression, and operational security

measures. The system analyzes patterns of behavior that indicate nation-state or organized criminal group involvement.

Script kiddie classification identifies attackers utilizing basic attack tools without deep understanding of underlying vulnerabilities. The system recognizes patterns associated with automated tool usage, basic payload manipulation, and unsophisticated attack sequences.

Automated scanner detection differentiates between human-operated attacks and automated vulnerability scanning activities through analysis of request timing patterns, tool signatures, and behavioral consistency indicators.

### 4.3.2   Advanced Behavioral Metrics and Analysis

The behavioral analysis engine implements comprehensive metrics that quantify various aspects of attacker behavior including request rate patterns, attack diversity measures, timing consistency analysis, error handling behaviors, and persistence indicators.

Request rate analysis utilizes statistical techniques to identify patterns in attack timing that indicate automation levels, human operator characteristics, and tactical approaches. The system calculates rolling averages, variance measures, and pattern correlation metrics to assess request timing characteristics.

Attack diversity metrics implement Shannon entropy calculations to measure the variety of attack techniques employed by individual attackers. Higher diversity scores indicate more sophisticated attackers with broader knowledge of vulnerability categories, while lower scores suggest focused or automated attack approaches.

Timing consistency analysis evaluates the regularity of request intervals to distinguish between human-operated attacks and automated scanning activities. The system calculates coefficient of variation measures and implements statistical tests to identify automation indicators.

Persistence scoring algorithms assess attacker dedication and long-term engagement characteristics through analysis of session duration, return visit patterns, and sustained attack activities. Persistence scores provide insights into attacker motivation and threat assessment priorities.

### 4.3.3   Tool Detection and Signature Analysis

The tool detection system maintains comprehensive databases of security tool signatures, enabling identification of specific tools employed by attackers including commercial security scanners, open-source vulnerability assessment tools, custom exploitation frameworks, and manual testing techniques.

User-Agent analysis algorithms identify security tools through comprehensive pattern matching against known tool signatures. The system recognizes popular tools including SQLMap, Nikto, Burp Suite, OWASP ZAP, Nessus, OpenVAS, and custom tool variants.

Request pattern analysis identifies tool-specific behaviors through analysis of request sequences, parameter manipulation patterns, and error handling approaches. Different tools exhibit characteristic beha-

vioral patterns that enable identification even when User-Agent strings are modified or obfuscated.

The system implements tool version identification capabilities that can distinguish between different versions of the same tool based on behavioral pattern variations and signature evolution. Version identification provides insights into attacker tool currency and sophistication levels.

### 4.3.4 Attack Pattern Recognition and Classification

The attack pattern recognition system implements a comprehensive taxonomy of attack patterns that characterize different types of threat activities. The system recognizes ten distinct attack patterns including automated vulnerability scanning, manual security testing, targeted exploitation campaigns, reconnaissance missions, script kiddie activities, advanced persistent threats, web application fuzzing, injection attack specialization, botnet reconnaissance, and insider threat activities.

Each attack pattern is defined through characteristic behavioral indicators, required technical markers, confidence thresholds, and example attack sequences. The pattern matching system evaluates incoming requests against these pattern definitions to identify the most appropriate classification for observed activities.

Pattern confidence scoring algorithms assess the strength of evidence supporting specific pattern classifications through analysis of behavioral consistency, technical marker presence, and historical accuracy metrics. The system provides multiple potential pattern matches with associated confidence scores to account for behavioral ambiguity and classification uncertainty.

## 4.4 Context-Aware Response Generation Framework

The context-aware response generation framework represents the culmination of AI/LLM integration, implementing sophisticated decision-making algorithms that adapt response strategies based on comprehensive analysis of attacker characteristics, threat assessment, and engagement objectives. This framework enables the honeypot to provide tailored responses that maximize intelligence gathering while maintaining operational security.

### 4.4.1 Intelligent Context Analysis Engine

The context analysis engine implements comprehensive assessment algorithms that evaluate multiple dimensions of each interaction to determine optimal response strategies. The analysis considers attacker sophistication levels, detected evasion techniques, threat levels, behavioral patterns, tool signatures, attack progression, and session history.

Sophistication scoring algorithms evaluate attacker technical capabilities through analysis of payload complexity, evasion technique usage, tool selection, and attack methodology. Sophistication scores range from 0.0 to 1.0, with higher scores indicating more advanced attackers requiring more sophisticated response strategies.

Evasion detection algorithms identify attempts to bypass security controls through analysis of encoding techniques, obfuscation methods, case manipulation, comment injection, and other anti-detection measures. Evasion detection informs response generation strategies that can mirror or counter detected evasion techniques.

Threat level assessment combines multiple risk factors including attack sophistication, MITRE technique classifications, persistence indicators, and potential impact assessments to generate comprehensive threat scores. Threat levels are classified as low, medium, high, or critical based on calculated risk scores.

## 4.4.2   Adaptive Response Strategy Selection

The response strategy engine implements five distinct response strategies designed for different attacker types and threat scenarios: Advanced Engagement for sophisticated attackers, Scanner Disruption for automated tools, Evasion Counter for attacks utilizing bypass techniques, Information Control for reconnaissance activities, and Professional Engagement for penetration testing scenarios.

Advanced Engagement strategies deploy complex, realistic responses designed to maintain engagement with sophisticated attackers while collecting comprehensive behavioral data. These responses include multi-layered vulnerability simulations, realistic error messages, and sophisticated application behaviors.

Scanner Disruption strategies implement responses designed to waste automated scanner resources through large response payloads, circular redirect loops, time-consuming content generation, and false positive vulnerability indicators. These strategies aim to increase scanning costs while providing minimal valuable intelligence to automated tools.

Evasion Counter strategies implement responses that acknowledge and mirror detected evasion techniques, demonstrating awareness of bypass attempts while maintaining honeypot effectiveness. These responses can include encoded content, anti-analysis measures, and technique-specific countermeasures.

Information Control strategies manage information disclosure for reconnaissance activities through selective data exposure, controlled vulnerability hints, and guided exploration paths. These strategies balance intelligence gathering objectives with operational security requirements.

Professional Engagement strategies provide sophisticated, realistic vulnerability simulations designed to engage professional security testers while collecting valuable intelligence about testing methodologies and tool usage patterns.

## 4.4.3   Dynamic Content Generation and Complexity Scaling

The dynamic content generation system implements adaptive algorithms that scale response complexity based on detected attacker sophistication and engagement requirements. The system utilizes five complexity levels ranging from basic responses for automated scanners to highly sophisticated simulations for advanced human operators.

Content templates provide structured frameworks for response generation while enabling customization

based on contextual requirements. Templates include variable substitution mechanisms, conditional content blocks, and dynamic data generation algorithms that create realistic application responses.

Complexity scaling algorithms evaluate attacker characteristics to determine appropriate response complexity levels. Higher complexity responses include detailed error messages, comprehensive debugging information, realistic database interactions, and sophisticated application logic simulations.

The system implements comprehensive content generation capabilities including realistic error message generation, convincing vulnerability simulation, detailed system information disclosure, and authentic application behavior modeling. Content generation maintains consistency with established application personas while adapting to specific interaction contexts.

### 4.4.4  Machine Learning Integration and Adaptive Optimization

The response generation system incorporates machine learning algorithms that continuously optimize response strategies based on observed effectiveness metrics and attacker engagement patterns. The system tracks response effectiveness indicators including session duration, follow-up activity levels, data collection quality, and attacker persistence measures.

Learning algorithms implement reinforcement learning techniques that adjust response strategy selection based on historical effectiveness data. The system maintains comprehensive databases of response outcomes that inform future strategy selection decisions.

Adaptive optimization algorithms continuously refine response generation parameters including complexity scaling factors, content selection criteria, and engagement strategy triggers. The optimization process considers both individual interaction outcomes and aggregate system performance metrics.

The system implements A/B testing capabilities that enable controlled experiments with different response strategies to identify optimal approaches for specific attacker types and scenarios. Experimental results inform strategy refinement and system optimization processes.

### 4.4.5  Integration with Behavioral Analysis and MITRE Classification

The response generation framework integrates seamlessly with behavioral analysis and MITRE classification systems to create comprehensive, context-aware response strategies. Integration enables response generation decisions based on detected attacker types, identified MITRE techniques, and behavioral pattern classifications.

Behavioral integration algorithms utilize attacker profiling data to inform response strategy selection, with different strategies triggered for different attacker types. Professional penetration testers receive sophisticated vulnerability simulations, while script kiddies encounter basic responses with educational false positives.

MITRE technique integration enables response strategies tailored to specific attack techniques and tactics. The system can generate responses appropriate for detected initial access attempts, execution techniques, persistence mechanisms, and defense evasion strategies.

The integrated system provides comprehensive feedback loops that enable continuous refinement of classification algorithms based on response effectiveness data. Successful response strategies provide validation for classification decisions, while ineffective responses indicate potential classification errors requiring algorithm adjustment.

## 4.5  Research Data Export and Analysis Capabilities

The enhanced Galah system implements comprehensive research data export capabilities that transform collected honeypot intelligence into structured datasets suitable for academic research, threat intelligence analysis, and operational security assessment. The export system supports multiple data formats and analysis frameworks to accommodate diverse research requirements and analytical approaches.

### 4.5.1  Structured Dataset Generation

The dataset generation system creates comprehensive structured datasets that capture all relevant aspects of honeypot interactions including temporal information, attacker characteristics, behavioral metrics, MITRE classifications, response strategies, and effectiveness measures. Datasets are generated in multiple formats including JSON for programmatic analysis, CSV for statistical software integration, and specialized formats for machine learning frameworks.

Temporal data structuring ensures that all dataset entries include precise timestamp information enabling time-series analysis, attack pattern evolution studies, and longitudinal behavioral assessment. The system maintains nanosecond precision timestamps while providing various temporal aggregation capabilities for different analytical requirements.

Behavioral data normalization algorithms ensure that behavioral metrics are consistently scaled and formatted across different dataset exports. Normalization processes account for session duration variations, request rate differences, and attack intensity variations to enable meaningful comparative analysis.

The system implements comprehensive data validation algorithms that ensure dataset integrity and consistency through automated quality checks, missing value detection, outlier identification, and consistency verification across related data fields.

### 4.5.2  Privacy-Preserving Data Processing

Privacy protection mechanisms ensure that exported datasets comply with relevant data protection regulations while preserving analytical value for research purposes. The system implements multiple anonymization techniques including IP address masking, request sanitization, timestamp fuzzing, and selective field redaction.

IP address anonymization utilizes cryptographic hashing techniques that preserve geographical and network relationship information while preventing identification of specific source addresses. The anonymization process maintains sufficient entropy for network analysis while ensuring irreversible protection of source identity.

Request sanitization algorithms remove potentially personally identifiable information from request content while preserving attack pattern characteristics and behavioral indicators. Sanitization processes utilize regular expressions, natural language processing, and machine learning classification to identify and protect sensitive information.

The system implements differential privacy techniques for aggregate statistics and behavioral metrics, ensuring that individual attacker behaviors cannot be reverse-engineered from published research datasets while maintaining statistical utility for research purposes.

### 4.5.3 Academic Research Integration

The research integration capabilities enable seamless incorporation of honeypot data into academic research frameworks through standardized interfaces, established research methodologies, and compliance with academic ethical guidelines. The system supports integration with popular academic software including R, Python pandas, MATLAB, and statistical analysis packages.

Reproducibility support mechanisms ensure that research conducted using honeypot datasets can be reliably reproduced through comprehensive metadata documentation, version control integration, experimental parameter tracking, and deterministic dataset generation capabilities.

The system implements support for controlled experimental designs including A/B testing frameworks, cohort analysis capabilities, and longitudinal study methodologies. Experimental design features enable researchers to conduct rigorous studies of honeypot effectiveness, attacker behavior patterns, and response strategy optimization.

Citation and attribution mechanisms provide standardized formats for academic citation of honeypot datasets and research results, ensuring proper credit attribution while facilitating academic collaboration and knowledge sharing.

## 4.6 Performance Analysis and Scalability Considerations

The AI/LLM extensions to Galah introduce additional computational requirements and performance considerations that must be carefully managed to maintain operational effectiveness while providing enhanced intelligence capabilities. Comprehensive performance analysis reveals the resource utilization patterns, scalability characteristics, and optimization opportunities inherent in the enhanced system architecture.

### 4.6.1 Computational Resource Requirements

Behavioral analysis algorithms introduce moderate CPU overhead during request processing, with typical processing times ranging from 10-50 milliseconds per request depending on payload complexity and analysis depth. The analysis algorithms are designed for incremental processing, enabling efficient utilization of available computational resources while maintaining real-time response capabilities.

Memory utilization patterns demonstrate stable consumption characteristics with periodic spikes during machine learning model inference and complex behavioral analysis operations. The system maintains comprehensive memory management algorithms that prevent memory leaks while optimizing cache utilization for frequently accessed data structures.

MITRE classification processes require additional computational resources for pattern matching and evidence evaluation, with typical classification times ranging from 5-20 milliseconds per request. The classification system utilizes optimized data structures and caching mechanisms to minimize computational overhead while maintaining classification accuracy.

Context-aware response generation introduces the most significant computational requirements, particularly during LLM inference operations that may require 1-10 seconds depending on model selection and provider performance. The system implements comprehensive caching and optimization strategies to minimize LLM API costs while maintaining response quality.

### 4.6.2 Scalability Architecture and Load Distribution

The enhanced system architecture supports horizontal scaling through containerized deployment models that enable distribution of computational load across multiple processing nodes. Container orchestration platforms provide automatic scaling capabilities based on request volume, resource utilization, and performance metrics.

Load balancing algorithms distribute incoming requests across multiple processing instances while maintaining session affinity requirements for behavioral analysis and response consistency. The load distribution system considers both computational capacity and specialized processing capabilities when routing requests to appropriate processing nodes.

Database scaling strategies accommodate the increased data storage requirements introduced by enhanced logging, behavioral analysis data, and comprehensive session tracking. The system supports both vertical scaling through increased database resources and horizontal scaling through database sharding and replication strategies.

Caching optimization algorithms maximize cache hit rates while managing the increased cache requirements introduced by behavioral analysis data and context-aware response generation. Multi-level caching strategies balance memory utilization with performance optimization objectives.

### 4.6.3 Performance Optimization Strategies

Algorithm optimization techniques minimize computational overhead through efficient data structure utilization, optimized pattern matching algorithms, and streamlined processing pipelines. Performance profiling identifies bottlenecks and optimization opportunities throughout the enhanced system architecture.

Asynchronous processing capabilities enable non-blocking request handling while performing computationally intensive analysis operations in background threads. Asynchronous architectures improve system

responsiveness while maintaining comprehensive analysis capabilities.

Batch processing algorithms enable efficient handling of high-volume scenarios through aggregated analysis operations, bulk database transactions, and optimized resource utilization patterns. Batch processing capabilities are particularly beneficial for research data export and comprehensive behavioral analysis operations.

The system implements comprehensive performance monitoring that tracks response times, resource utilization, cache hit rates, and analysis processing metrics. Performance monitoring enables proactive optimization and capacity planning for production deployments.

# A   Appendix

This section contains additional material that supports the thesis, such as code listings, large data sets, or additional explanations.

## A.1 Declaration of Authorship

The following text is to be attached to the thesis.

"I hereby declare,

- that I have written this thesis independently,
- that I have written the thesis using only the aids specified in the index;
- that all parts of the thesis produced with the help of aids have been declared;
- that I have handled both input and output responsibly when using AI. I confirm that I have therefore only read in public data or data released with consent and that I have checked, declared and comprehensibly referenced all results and/or other forms of AI assistance in the required form and that I am aware that I am responsible if   incorrect content, violations of data protection law, copyright law or scientific misconduct (e.g. plagiarism) have also occurred unintentionally;
- that I have mentioned all sources used and cited them correctly according to established academic citation rules;
- that I have acquired all immaterial rights to any materials I may have used, such as images or graphics, or that these materials were created by me;
- that the topic, the thesis or parts of it have not already been the object of any work or examination of another course, unless this has been expressly agreed with the faculty member in advance and is stated as such in the thesis;
- that I am aware of the legal provisions regarding the publication and dissemination of parts or the entire thesis and that I comply with them accordingly;
- that I am aware that my thesis can be electronically checked for plagiarism and for third-party authorship of human or technical origin and that I hereby grant the University of St.Gallen the copyright according to the Examination Regulations as far as it is necessary for the administrative actions;
- that I am aware that the University will prosecute a violation of this Declaration of Authorship and that disciplinary as well as criminal consequences may result, which may lead to expulsion from the University or to the withdrawal of my title."

By submitting this thesis, I confirm through my conclusive action that I am submitting the Declaration of Authorship, that I have read and understood it, and that it is true.

Date and signature

……………………………………………

## A.2  Declaration of Confidentiality

<div style="border:1px solid;">

**Declaration of Confidentiality**

in connection with written work at the University of St.Gallen

</div>

The following text has to be added to the work and signed:

The undersigned

hereby undertakes and warrants to treat any information obtained by the enterprise/ administration concerned in strict confidentiality. In particular, he / she shall only permit people other than the referees to inspect his / her written work with the express consent of all the parties that have provided information.

The undersigned hereby takes cognizance of the fact that the University of St.Gallen may check his / her work for any plagiarism with the help of a plagiarism software and that the undersigned shall have to notify the enterprise/administration surveyed accordingly.

Date and signature

……………………………………………