# CANARY TOKENS AS A STRATEGIC COMPONENT IN CYBERSECURITY DEFENSE AND RED TEAMING

D. Kabiden[1]*

[1] International Information Technology Univrsity, Almaty, Kazakhstan.
*Corresponding author: dinmuhammedkabiden@yandex.kz

**Abstract**

Canary tokens, also known as honeytokens, are small and lightweight decoy artifacts that act as silent tripwires in digital environments. Their simplicity and flexibility make them one of the most efficient deception-based mechanisms for early detection of breaches. This paper investigates the dual role of canary tokens in cybersecurity, emphasizing their value for Blue Teams (defenders) and their unexpected adoption by Red Teams (attackers). Through analysis of real incidents such as the Grafana Labs breach (2025) and Iranian APT MuddyWater campaigns, we illustrate how these tools operate in practice. Defensive use focuses on reducing attacker dwell time by planting tokens like fake credentials, documents, or URLs across infrastructure. Offensive use includes reconnaissance, phishing confirmation, and even malware anti-analysis. We also examine legal and ethical dimensions, suggesting that when appropriately scoped and governed, canary tokens can be deployed in a manner consistent with data protection principles and ethical cybersecurity practices. Finally, we propose best practices for token deployment: realistic design, diverse placement, SIEM integration, controlled rotation, and minimization of noise. The findings demonstrate that canary tokens, although deceptively simple, play a vital role in defense-in-depth strategies by frustrating attackers, empowering defenders, and shifting the balance of cyber operations.

*Keywords: Honeytokens, cybersecurity defense, deception technology, early breach detection, legal and ethical considerations.*

## 1. Introduction

In the modern era of rapidly evolving cyber threats, organizations across all sectors face an ever-increasing challenge of protecting sensitive data, digital assets, and

critical infrastructures. Cyberattacks have become more sophisticated, persistent, and targeted, often bypassing traditional security controls such as firewalls, antivirus software, and intrusion prevention systems. Consequently, the paradigm of cybersecurity has gradually shifted from purely preventive measures to proactive and adaptive defense mechanisms that can identify, mislead, and contain adversaries before serious damage occurs. One of the most promising directions in this transformation is the strategic use of deception technologies, which introduce controlled uncertainty into the attacker's decision-making process [1].

Historically, deception in information security emerged through the deployment of honeypots and decoy systems—isolated digital environments designed to imitate legitimate resources and attract unauthorized access. These early systems played a critical role in gathering intelligence on attack methodologies, identifying vulnerabilities, and understanding adversarial behavior. However, as cyber infrastructures expanded and diversified, maintaining large-scale honeypot systems became increasingly resource-intensive. The need for a lighter, more scalable approach gave rise to canary tokens, also known as honeytokens—compact, software-based indicators designed to act as silent tripwires within digital ecosystems.

Canary tokens represent the minimalist evolution of deception: they are lightweight, easily deployable, and cost-effective artifacts that appear valuable to an intruder but serve a single purpose—to notify defenders immediately when accessed. Their metaphorical name, derived from the "canary in the coal mine," emphasizes their early-warning function: these digital markers emit a signal of compromise at the first indication of malicious activity [2]. Tokens can take multiple forms—such as fake API keys embedded in code repositories, seemingly sensitive spreadsheets on corporate file servers, or dormant DNS records in configuration files. The core advantage of this approach lies in its precision: any interaction with a canary token is inherently suspicious, providing an unambiguous indicator of unauthorized behavior without the noise commonly associated with traditional alerting mechanisms.

The importance of such early detection mechanisms cannot be overstated. According to numerous incident response studies, the average dwell time of attackers within compromised environments often extends to several weeks or even months before discovery. This prolonged invisibility enables data theft, privilege escalation, and infrastructure manipulation on a large scale. By embedding deceptive elements within operational systems, defenders effectively compel attackers to reveal themselves, drastically reducing detection latency [3]. Notably, this defensive innovation has also inspired offensive adaptation: advanced threat actors, including state-sponsored groups, have begun to incorporate canary-like mechanisms into their own malware to detect sandbox environments or monitor victim engagement.

The objective of this study is to conduct a comprehensive analysis of canary to-

kens as a dual-use cybersecurity instrument. The research examines their defensive deployment strategies, offensive adaptations by threat actors, and the legal-ethical boundaries governing their application. In addition, this paper aims to synthesize best practices for integrating canary tokens into organizational defense-in-depth architectures, emphasizing their role in enhancing situational awareness, accelerating incident response, and supporting cyber resilience.

## 2. Methods

### 2.1 Types of Canary Tokens

Canary tokens can be classified into several distinct types depending on their technical form, purpose, and operational context. Each type plays a unique role within the broader deception strategy, and their collective use creates a multi-layered detection environment.

A concise taxonomy of canary tokens and their main properties is presented in Table 1, which systematizes their structure, detection mechanisms, and deployment areas.

*Table 1. Classification and characteristics of canary tokens*

| Type of Token | Typical Form / Example | Primary Purpose | Detection Mechanism | Typical Deployment Context |
|---|---|---|---|---|
| Credential Token | Fake API key, SSH private key, database password | Detect unauthorized access or credential theft | Callback via HTTP or API validation | Repositories, CI/CD pipelines, cloud storage |
| Document Token | Passwords.xlsx, HR_Payroll.pdf, Finance_Report_2025.docx | Detect insider threats, phishing validation | Embedded web beacon or tracking URL | File servers, shared folders, email attachments |
| Web/URL Token | Unique or fake admin link | Detect reconnaissance or scanning | HTTP request to tokenized URL | Web apps, documentation portals |
| DNS Token | Fake DNS record (e.g., db-secure.example.com) | Detect malware beaconing or network enumeration | Logged DNS query to token server | Cloud and hybrid networks |
| Custom Token | Registry key, folder path, or event log | Detect local privilege escalation or insider activity | System event or local agent alert | Endpoints, domain controllers |

Credential tokens typically simulate sensitive authentication data—such as API keys, database credentials, or SSH keys—and are primarily designed to detect unauthorized access or credential theft. Their activation through an API validation event or HTTP callback provides a clear and unambiguous signal of intrusion.

Document tokens, on the other hand, are embedded within files that appear valuable to an attacker, such as Passwords.xlsx or Finance_Report_2025.docx. When opened, these files silently trigger tracking beacons or web callbacks. Such tokens are effective in monitoring insider behavior, detecting phishing campaigns, and preventing exfiltration of confidential data.

Web or URL tokens take the form of fake admin links or unique pages that should never be accessed during legitimate activity. When visited, they generate HTTP requests that immediately notify defenders of reconnaissance or scanning attempts. DNS tokens function similarly but operate at the network level: they consist of fake subdomains embedded in configurations or scripts. Each time a malicious actor or automated malware resolves the domain, the corresponding query is recorded and logged by the defender's monitoring system.

Custom tokens extend the concept further by adapting deception to local infrastructure needs. Examples include registry keys, fake log entries, or decoy directory paths that trigger alerts upon unauthorized access. These tokens are especially useful in internal network monitoring and endpoint protection where traditional detection systems may have limited coverage.

Through such differentiation, defenders can construct a distributed detection grid that operates silently yet effectively. By deploying various token types across endpoints, servers, cloud services, and network layers, organizations achieve defense-in-depth visibility and enhance their capacity for rapid threat detection.

### 2.2 Method of Analysis

The methodological foundation of this research is based on a qualitative synthesis of academic and operational data supported by targeted case study analysis. This dual approach allows for a balanced understanding of both theoretical models and their real-world implementations.

1. Literature Review. A systematic review of scientific publications, white papers, and open-source intelligence materials was performed to examine the evolution of deception technologies, the technical design of canary tokens, and their integration within Security Operations Centers (SOCs). Special attention was given to sources discussing implementation efficiency, detection accuracy, and automation within SIEM and SOAR systems.

2. Case Study Analysis. Two well-documented incidents were selected for empirical comparison. The Grafana Labs (2025) incident exemplified the defensive deployment of canary tokens, where pre-seeded fake AWS credentials allowed

immediate detection of a compromised GitHub Action pipeline[3].

Conversely, the MuddyWater (Iranian APT) operation represented an offensive adaptation, where malware embedded canary-like callbacks to verify execution on genuine victims and detect sandbox analysis[4].

These contrasting scenarios demonstrate how identical mechanisms can be leveraged for both defense and offense, depending on intent.

3. Comparative Review of Security Reports. To validate and generalize observations, technical reports and advisories from cybersecurity vendors were analyzed. These included detailed discussions on token configuration, false-positive reduction, and attacker evasion techniques, enabling a cross-comparison of practical implementations.

4. Legal frameworks. A review of relevant data protection and surveillance regulations (GDPR, U.S. Wiretap Act) and academic discussions on entrapment-related concerns in cybersecurity monitoring.

## 3. Results and Discussion

### 3.1 Defensive use of canary tokens by Blue Teams

The empirical data collected from case studies and technical reports confirm that canary tokens serve as an effective and low-cost detection mechanism for Blue Teams. Unlike traditional intrusion detection alerts, which often generate a significant number of false positives due to noisy baselines or benign anomalies, a canary token alert carries a strong semantic weight. Its activation almost always indicates that an unauthorized party has interacted with a resource that should remain untouched in legitimate workflows[5].

A notable case was documented in Grafana Labs (2025), where adversaries exploited a misconfigured GitHub Action to exfiltrate environment variables. The attackers obtained what appeared to be Amazon Web Services (AWS) credentials. However, the security team had pre-seeded these repositories with a strategically placed canary AWS key. When the adversary attempted to validate the stolen key through the AWS API, the token immediately triggered an alert. This rapid signal allowed the Security Operations Center (SOC) to mobilize, rotate all exposed secrets, and suspend vulnerable workflows in less than 24 hours. Without the early detection provided by the token, the dwell time of the attacker could have extended for days or weeks, significantly increasing potential damage.

Key benefits observed in Blue Team deployments include:

– High-fidelity alerts: Canary tokens eliminate ambiguity because their use has no legitimate operational context. A triggered token directly correlates with a breach attempt;

– Contextual intelligence: Alerts contain valuable metadata such as the source IP,

user agent, geolocation, and the specific token identifier. This contextual information reduces investigative delays and directs analysts to the compromised environment;

– Psychological deterrence: Awareness among attackers that organizations deploy deception technologies often leads to hesitation, increasing their cognitive load and slowing down lateral movement. Penetration testers report that encountering tokens forces adversaries to second-guess every credential or file.

Best practices synthesized from industry reports are as follows:

– Strategic placement: Tokens should be positioned where attackers are most likely to search for sensitive material, such as configuration repositories, CI/CD pipelines, or network share drives;

– Realism: The naming and structure of tokens must mimic real assets. For example, filenames like Finance_2025_Budget.xlsx or credentials labeled DB_Admin increase credibility;

– Integration into monitoring systems: Token alerts should be routed into SIEM platforms (e.g., Splunk, Elastic, QRadar) and incident response workflows. Some organizations use tiered escalation, sending high-priority token triggers directly to on-call engineers via Slack, PagerDuty, or SMS.

– Rotation and renewal: To maintain effectiveness, tokens must be periodically refreshed. Static tokens may become outdated or recognized by sophisticated adversaries.

– Noise reduction: False positives are minimized by whitelisting internal scanners, backup processes, and trusted IP ranges. Overexposure of tokens in public repositories should also be avoided to reduce accidental triggering by automated indexing tools.

Collectively, these practices demonstrate that canary tokens, when deployed systematically, significantly shorten mean time to detection (MTTD) and enhance organizational resilience against intrusions

### 3.2 Offensive and Reconnaissance Applications by Red Teams

While originally conceived as a defensive tool, canary tokens have been creatively repurposed by Red Teams and threat actors for offensive objectives. The Iranian state-sponsored group MuddyWater (APT34) illustrates this phenomenon vividly.

Their operations included embedding canary tokens directly into Visual Basic Script (VBS) malware loaders. The logic was as follows: when executed in a sandbox or automated analysis environment, the script would trigger only one HTTP request to the token server. On a real victim machine, however, the script would produce two sequential requests. If only a single request was detected, the malware concluded that it was under analysis and terminated execution. This method served as an anti-anal-

ysis mechanism, allowing MuddyWater to evade sandbox detonation and ensuring that payloads were delivered only to genuine victims.

In another campaign documented by Cisco Talos (2022), MuddyWater distributed malicious Excel documents containing links to canarytokens.org. When a target opened the attachment, a callback was triggered, notifying the attackers that their phishing attempt was successful. This telemetry allowed them to selectively invest resources into compromised machines, while ignoring unopened lures[4].

Additional offensive use cases include:

– Tracking defender behavior: By planting fake credentials in compromised environments, Red Teams can detect when Blue Teams attempt to test these credentials, thus confirming that their activities have been discovered.

– Phishing validation: Tokens embedded in spear-phishing documents act as "read receipts," enabling attackers to measure campaign effectiveness in real time.

– Stealthy signaling: In certain malware families, token callbacks are used as covert beacons to confirm data exfiltration or payload execution without relying on noisy C2 traffic.

This offensive adaptation illustrates the dual-use nature of canary tokens. What serves defenders as an early detection tripwire can equally empower adversaries to improve stealth, situational awareness, and operational efficiency.

### 3.3 Countermeasures, evasion, and defensive responses

As awareness of canary tokens grows, attackers increasingly develop counter-detection techniques. Document analysis tools can scan for hidden web bugs, unusual external references, or suspicious metadata. Skilled adversaries may inspect the XML structure of Word or Excel files to identify embedded URLs pointing to token services. Similarly, utilities like binwalk or strings can reveal hidden callbacks in executables or PDFs.

Attackers also experiment with network isolation techniques. By executing potentially tokenized files within isolated virtual environments, behind VPNs, or over anonymizing networks like TOR, adversaries attempt to decouple token alerts from their true identities. In some instances, attackers deliberately trigger every suspected token in a compromised system in order to flood the defender's SOC with alerts, thereby degrading response capacity[6].

In response, defenders employ several strategies:

– Custom domains and infrastructure: Rather than relying on recognizable domains such as canarytokens.org, defenders host tokens under organization-specific domains, making them indistinguishable from legitimate resources.

– Stealth tokens requiring full validation: Advanced token types only generate alerts after successful API interaction (e.g., AWS credential validation). Such to-

kens cannot be identified by simple string matching, forcing attackers to interact more deeply and reveal themselves.

 – Threat hunting integrations: SOC teams monitor for outbound requests to known token providers or suspicious custom domains, enabling detection of adversaries who attempt to misuse token services themselves.

### 3.4 Integration and Best Practices

Effective implementation of canary tokens requires more than just creating false artifacts; it requires a structured integration strategy within a broader deep security architecture. Based on our synthesis of case studies, technical guidance, and best practices, we identify several important aspects of successful integration.

Multi-level placement: Canary tokens are most effective when distributed across different levels of the IT environment. At the endpoint level, defenders can insert tokens into local configuration files, browser credential stores, or hidden directories that attackers commonly visit after gaining initial access. At the network level, tokens can be embedded in DNS records, bait ports, or fake service credentials that attract third-party users. At the application level, tokens can be disguised as database records or API keys embedded in configuration stores. Finally, in a cloud environment, defenders often use AWS, Azure, or GCP canary credentials, which, when verified by an attacker, trigger alerts with full contextual data (such as region, IP address, and timestamp). This multi-layered approach ensures that even if an attacker evades one trap, another can detect their actions[7].

Synergy with honeypots: While canary tokens provide silent operation, honeypots offer enhanced interaction with attackers. The combination of these two factors enhances the protection. For example, a token credential embedded in a developer's repository can lead an attacker to a honeypot system that mimics a production server. This multi-level deception delays attackers, provides insights into their methods, and creates a more comprehensive database for forensic analysis. The academic literature on deception technologies suggests that these multi-level approaches significantly increase the time that attackers spend in controlled environments, diverting them from their actual assets.

Automatic alerts and action patterns: The value of a canary token lies in its speed of response. Token alerts should be directly integrated into incident response (IR) programs and security orchestration, automation, and response (SOAR) systems. For example, an activated token representing AWS credentials can automatically disable the corresponding account, quarantine the computer from which the token was accessed, and trigger an event to the on-call staff via PagerDuty or Slack. This automation minimizes staff delays and ensures that defenders can act within minutes rather than hours.

Token lifecycle management (rotation and renewal): Outdated or predictable to-

kens lose their effectiveness over time. Best practices require periodic token rotation (quarterly or semi-annually, depending on the environment's dynamics). Rotation reduces the risk of a token being previously discovered, catalogued, or neutralized by attackers. Modern approaches include automating token updates using CI/CD pipelines, so each new deployment includes updated tokens. This practice not only keeps tokens fresh, but also provides defenders with a clearer temporal correlation between leakage and usage[12].

Testing and validation: Just as fire alarms require regular training, canary tokens should be tested using red team or purple team exercises. Simulated breaches confirm that attackers can detect tokens, that alerts are triggered correctly, and that the SOC responds in accordance with established procedures. The collaboration between the purple team is particularly valuable, as defenders can observe how the red teams interact with the tokens, while the red teams evaluate the realism and detectability of the deployed decoys. In some organizations, token testing is included in quarterly cyber resilience assessments.

Policy and legal alignment: Finally, the use of canary tokens must be consistent with the organization's policies and legal framework. From a governance perspective, tokens should be clearly outlined in the organization's internal security documentation as part of its fraud prevention and monitoring strategy. Legally, organizations must ensure compliance with privacy regulations such as the GDPR, particularly when token callbacks may capture IP addresses or other identifiers. Transparency is achieved by incorporating monitoring provisions into the organization's acceptable use policy and documenting the legitimate interest in using tokens. Ethical principles recommend avoiding overly intrusive information placement (such as in employee personal files) and instead focusing on areas where access by malicious actors is the most likely explanation.

### 4. Conclusion

Canary tokens represent a modern and efficient evolution of deception-based cybersecurity mechanisms. Their lightweight nature, ease of deployment, and high–fidelity alerting make them a valuable component of proactive defense strategies. Unlike traditional detection systems, which often generate large volumes of ambiguous data, canary tokens provide clear and actionable signals that allow defenders to identify and respond to intrusions within minutes rather than days.

The analysis of real-world cases, such as the Grafana Labs incident and the Muddy-Water campaigns, demonstrates their dual functionality in both defensive and offensive contexts. When properly implemented, canary tokens enhance visibility, reduce attacker dwell time, and strengthen the resilience of organizational infrastructures.

Integrating these tokens into a multi-layered security architecture — with automated alerting, periodic rotation, and adherence to organizational policies — allows se-

curity teams to maintain continuous situational awareness. In an era where breaches are inevitable, early detection remains decisive. Therefore, canary tokens serve not only as silent guardians of digital assets but also as practical instruments that shift the balance of cyber operations in favor of defenders.

### References

1. Neuwirth, C. (2025). Canarytokens: Zero-Cost Tripwires That All Blue Teams Should Be Using. NetWorks Group Blog.

2. Webb, B. (2025). Digital Age: The Power of Canary Tokens and Deception Technology in Cybersecurity. Recon Infosec Blog.

3. Moradian, M. (2025). Canary tokens: Learn all about the unsung heroes of security at Grafana Labs. Grafana Labs Blog.

4. Cisco Talos (2022). Iranian-linked conglomerate MuddyWater comprised of regionally focused subgroups. Talos Intelligence Report.

5. Lemos, R. (2022). Credential Canaries Create Minefield for Attackers. Dark Reading.

6. Netsurion (2023). Are honeypots illegal? Netsurion Insights.

7. Ali AK (2023). Detecting & Bypassing Defensive Measures (Canary Token). InfoSec Write-ups.

8. Mthcht (2023). Canary Tokens and Callback URLs: A Double-Edged Sword. Medium.

9. EU GDPR Analysis (2017). Honeypots and honeynets: issues of privacy. EURASIP J. on Info. Security.

10. Tracebit (2023). The Security Canary Maturity Model. Tracebit Blog.

11. Morić, Z., Dakić, V., and Regvart, D. (2025). Advancing Cybersecurity with Honeypots and Deception Strategies. Informatics (MDPI).

12. Nelson, A., Rekhi, S., Souppaya, M., and Scarfone, K. (2024). Incident Response Recommendations and Considerations for Cybersecurity Risk Management (NIST SP 800-61r3). NIST Special Publication.

13. Virvilis, N., and Gritzalis, D. (2013). The Big Four – What we did wrong in Advanced Persistent Threat detection? International Conference on Availability, Reliability and Security (ARES).

14. Sokol, P., Mišek, J., and Husák, M. (2017). Honeypots and honeynets: issues of privacy. EURASIP Journal on Information Security.

15. Almeshekah, M. H., and Spafford, E. H. (2016). Cyber Security Deception.

## Information about authors

**Kabiden Dinmukhammed Ramazanuly** - International Information Technology Univrsity, Almaty, Kazakhstan.

**e-mail:** dinmuhammedkabiden@yandex.kz

**ORCID:** 0009-0009-4593-0547