# Blackwall: An integrated AI-driven Framework for proactive Cybersecurity Defense

Basil Abdullah Alzahrani

Department of Management
Information System

Al-Baha university
Al-Baha, Saudi Arabia
444019967@stu.bu.edu.sa

*Abstract*— Modern cybersecurity threats such as AI-driven attacks and zero-day exploits are overwhelming traditional defense systems that rely on static signatures or predefined rules. This paper introduces BlackWall, an integrated AI-driven framework for proactive cybersecurity defense. BlackWall consists of three core components: the Reverse Zero-day Algorithm (RZA) for autonomous vulnerability discovery, the Zero-Trust Verification Module (TVM) for continuous trust evaluation, and the False Positive Protocol (FPP), a deception-based system designed to mislead attackers and reduce false alerts. These modules operate independently but coordinate through a secure signal bus, enabling rapid response, resilience, and adaptive defense. In simulated attack scenarios, BlackWall improved detection accuracy by 37%, reduced response time by 82%, and cut false positives by 74% compared to traditional intrusion detection systems. These results demonstrate that BlackWall can serve as a robust foundation for next-generation, autonomous cybersecurity architectures.

*Keywords*— Cybersecurity, artificial intelligence, zero-day vulnerabilities, zero-trust architecture, deception systems, intrusion detection.

## I. INTRODUCTION

Traditional cybersecurity systems struggle against modern threats such as zero-day exploits, ai-driven attacks, and adaptive malware. Static rule-based defenses are no longer sufficient in dynamic environments where speed and deception are critical. We purpose BlackWall, an integrated AI-driven Framework that addresses detection, containment, and deception in real time. It consists of three main modules:

I. RZA—Reverse Zero-day algorithm for preemptive vulnerability detection

II. TVM—Zero-trust Verification Module for identity and access control.

III. FPP-False Positive Protocol for deception and alert refinement.

These components communicate through a secure signal layer and operate autonomously or in tandem.

Contributions:

- A modular cybersecurity framework combining detection, trust enforcement, and deception

- A novel algorithm (RZA) for identifying unknown vulnerabilities before adversaries can exploit them.

- A deception-based mechanism (FPP) to reduce false alarms, gathering information while quarantining the adversary, and lure adversaries.

## II. RELATED WORK

Traditional intrusion detection systems (IDS) such as snort [1] and Suricata [2] rely on predefined rules and signature matching. While effective against known threats, they fail to detect zero-day exploits and adaptive attacks that evolve beyond static patterns.

While Machine learning-based approaches emerged to address these gaps. Works such as [3] and [4] apply anomaly detection to network, but suffer from high false positive rates and limited response capabilities. Other studies including [5], integrated AI for threats classifications, yet lack modularity and real-time containment mechanisms.

Zero-trust Framework [6] emphasize strict identity verification but typically do not include active threat detection or deception layers. Meanwhile, deception-based systems such as honeypots [7] have shown promise in diverting attackers but remain isolated from broader security architectures.

BlackWall differentiates itself by integrating internal fuzzing (RZA), verification (TVM), and deception (FPP) into a cohesive modular framework. Unlike prior work, it enables the autonomous decision-making, adaptive response, and real-time containment – addressing gaps in precision, speed, and architectural integration.

## III. SYSTEM ARCHITECTURE

BlackWall is designed as a modular, AI-driven cybersecurity framework composed of three core components:

A. reverse zero-day algorithm (rza): Predicts and flags unknown vulnerabilities by analyzing anomaly behavior patterns and pre-exploit indicators.

B. zero-trust verification Module (tvm): continuously validates user and device authenticity using contextual data and policy-driven access rules.

C. false positive protocol (fpp): employs deception and behavioral traps to reduce false alerts while disorienting potential intruders.

Each module is deployed independently but interconnected via a secure Signal bus, which acts as an encrypted, low-latency communications channel between subsystems. This design allows for adaptive collaboration: for instance, a detection by RZA can trigger validation by TVM or initiate a decoy response from FPP.

Additionally, the modular nature of BlackWall enables plug-and-play functionality, cloud or on-premises deployments, the ability to scale or isolate components as needed. The

architecture supports both live mode (real-time defense) and forensic mode (incident reconstruction).

Figure 1 illustrates the overall structure and interactions between modules.

D. *The interactions between modules can be expressed through detection and trust threshold. For example, if anomaly score δ(t) from RZA exceets its threshold, it triggers a call to TVM for trust validation. If the resulting trust T(u,c) is below the minimum threshold τ, FPP engages to deploy a deception response. This logical flow allows autonomous collaboration without hardcoded rules. In short, it can be explained like this:*

$$Y = \begin{cases} Block + deceive, if\, \delta(t) > \mu + k\sigma\ and\ T(u,c) < \tau\ and\ F(x) > \theta \\ Monitor, if\, \delta(t) > \mu + k\sigma\ and\ T(u,c) \geq \tau \\ Allow, if\, \delta(t) \leq \mu + k\sigma \end{cases}$$

E. *Modular Workflow Execution*

BlackWall's architecture is intentionally modular, enabling flexible deployment in diverse environments such as security operation centers (SOCS), enterprise firewalls, or cloud-native platforms. The system does not rely on tight coupling between components, allowing RZA, TVM, and FPP to be run independently or chained in real time depending on the use case. This Modularity supports on-premise, hybrid and remote deployments.

Each module exposes an input-output interface, forming a pipeline where outputs from one layer are consumed by the next. For example, anomaly signals from RZA can be routed into a containerized TVM instance running zero-trust policies based on behavioral context. The final threat confirmation and deception logic in FPP can then respond accordingly, either by triggering honeypots, isolation sessions or escalating to human analysts.

This modular execution not only enables scalability, but also facilitates maintainability, upgrades, and tuning of thresholds $(\mu, \tau, \theta)$ without needing to halt the system.

## IV. EVALUATION

Blackwall was evaluated in a controlled environment designed to mimic real-world network traffic patterns, built around using python-based simulation scripts, incorporating synthetic network traffic, known attack vectors, (such as port scan, brute force, and privilege escalation), and randomized benign activity, datasets were modeled based on patterns observed in public corpora such as CICIDS2017 and NSL-KDD, allowing for generation of both labeled and malicious and non-malicious sessions.

A. *Results*

These results were derived from simulated network environments incorporating synthetic traffic, known attack patterns, and baseline user behavior. The high detection rate and low false positive rate reflects the effective synergy between RZA and TVM modules, while the FPP module demonstrated a strong deception performance by successfully engaging adversarial behavior in over 85% of attack scenarios.

| Results and metrics. | |
|---|---|
| *Detection rate (DR)* | *96.4%* |
| False positive rate (FPR) | 3.1% |
| Trust Accuracy (TVM module) | 91.2% |
| Deception Engagement Effectiveness | 85.7% |

*a) These results were obtained from a simulated enviroment usning synthetic network traffic and known attack vectors. The high detection rate and low false positive rate suggest strong coordination between anomaly detection and contextual trust scoring.*

## REFERENCES

[1] Y. Guo, "A review of machine learning-based zero-day attack detection: challenges and future directions," Computer Communications, vol. 198, pp. 1-18, Jan. 2023.

[2] M. Sarhan, S. Layeghy, M. Gallagher, N. Moustafa, and P. Watters, "From zero-shot machine learning to zero-day attack detection," International Journal of Information Security, vol. 22, pp. 947-959, 2023.

[3] S. Rose, O. Borchert, S. Mitchell, and S. Connelly, "Zero trust architecture," NIST Special Publication 800-207, National Institute of Standards and Technology, Aug. 2020.

[4] P. Dini, A. Elhanashi, A. Begni, S. Saponara, Q. Zheng, and K. Gasmi, "Overview on intrusion detection systems design exploiting machine learning for networking cybersecurity," Applied Sciences, vol. 13, no. 13, pp. 7507, June 2023.

[5] A. Javadpour, F. Ja'fari, and T. Taleb, "A comprehensive survey on cyber deception techniques to improve honeypot performance," Computers & Security, vol. 140, pp. 103732, May 2024.

[6] S. Masmoudi, E. Mezghani, S. Masmoudi, C. B. Amar, and H. Alhumyani, "Containerized cloud-based honeypot deception for tracking attackers," Scientific Reports, vol. 13, pp. 1896, Feb. 2023.

[7] H. Hindy, R. Atkinson, C. Tachtatzis, J.-N. Colin, E. Bayne, and X. Bellekens, "Utilising deep learning techniques for effective zero-day attack detection," Electronics, vol. 9, no. 10, pp. 1684, Oct. 2020.

[8] M. Issa, M. Aljanabi, and H. M. Muhialdeen, "Systematic literature review on intrusion detection systems: research trends, algorithms, methods, datasets, and limitations," Journal of Intelligent Systems, vol. 33, no. 1, pp. 20230248, Jan. 2024.

[9] Y. Mirsky, T. Doitshman, Y. Elovici, and A. Shabtai, "Kitsune: an ensemble of autoencoders for online network intrusion detection," in Proc. Network and Distributed System Security Symposium (NDSS), San Diego, CA, 2018.

[10] L. Bilge and T. Dumitraş, "Before we knew it: an empirical study of zero-day attacks in the wild," in Proc. 2012 ACM Conference on Computer and Communications Security, Raleigh, NC, Oct. 2012, pp. 833-844.