



## Securing GenAI Apps on AWS with CloudGuard

Extending the Foundations of Web App & Cloud  
Security to AWS-Powered GenAI Applications

### Standing on the shoulders of giants

Decades of experience have provided us with well-established tools and principles to protect clouds and web apps: From Next-Generation Hybrid Mesh Firewalls that protect and segment cloud and on-prem resources, to Web and API Firewalls (WAFs) that protect web-facing applications and exposed API endpoints.

However, with the meteoric rise of Generative AI (GenAI) applications, such as chatbot assistants in internal and external-facing applications, and multi-step/multi-agent agentic flows/systems with access to sensitive data and tools such as web-browsing and code execution, new attack vectors and unique risks have emerged that challenge our established notions of web app security.

Additionally, to reduce costs and rapidly adopt the era-defining revolution of GenAI technologies, companies are adopting hybrid cloud deployments, making infrastructure-agnostic IP-free firewalls an absolute necessity for consistent, repeatable security controls.

This white paper discusses Check Point's approach and its new capabilities to address the security risks of GenAI applications by **seamlessly extending** existing security controls and well-established security practices. Specifically, Check Point's Hybrid Mesh Firewall with extended cloud-native capabilities, and the CloudGuard Web Application and API Firewall (WAF), which has been supercharged with specially crafted GenAI security capabilities.

# Background

## Why we must secure traffic from packets up to APIs

The proliferation of networks across regions, branches, and on-premises data centers, along with the expansion into an increasing number of private and public cloud providers, as well as the increasing reliance on web applications and APIs with CVE-riddled open source dependencies, has led to the rise of integrated security controls such as centrally managed, infrastructure-agnostic, and IP-free Hybrid Mesh Firewalls. This shift also demanded that WAFs extend their security to API enforcement, with an increasing emphasis on embedded intrusion prevention systems capable of handling zero-day exploits.

Now, with the advent of GenAI, the traditional security blueprints and controls are being challenged yet again by new attack methods that network firewalls and WAFs cannot handle. For instance:

- Jailbreaking models directly (user prompts) or indirectly (prompt injection) to trick GenAI-powered assistants/agents to exfiltrate sensitive data in DLP-resistant outputs (e.g., output privileged information in ASCII art), [mint unauthorized coupons](#), [sell cars for \\$1](#), and more.
- Tricking AI Agents into launching attacks, such as [executing arbitrary code](#), performing SQL injections on an SQL database in Retrieval Augmented Generation (RAG) setups.
- The introduction of agentic frameworks such as CrewAI, Flowise, and n8n that introduce an ever-increasing number of CVEs in internet-exposed systems with direct or indirect access to highly sensitive assets via Model Context Protocol (MCP) servers.

Now, over three decades after Check Point invented the first-ever firewall, Check Point has extended the battle-tested foundations of web application and cloud security to address the challenges of GenAI security. To achieve this, Check Point has **(1)** added GenAI-specific security layers to its Web Application and API Firewall (WAF), addressing issues such as model jailbreaking, misuse, excessive agency, and more, and **(2)** increased its Hybrid Mesh Firewall's ability to inspect cloud-spanning networks with deep cloud-native integrations to facilitate rapid adoption of hosted GenAI-based systems.

In other words, Check Point now offers a holistic threat-prevention security solution to **secure traffic from the packets in your network, through APIs in your apps, and up to prompts in your LLMs.**

# Definitions and Disambiguation

- **AI:** Unless specified otherwise, for the purposes of this document, AI stands for Generative AI (GenAI), such as Large Language Models (LLMs), Vision Language Models (VLMs), etc.
- **AI Agents:** Autonomous AI-based systems that can perceive their environment, reason, and take action to achieve a specific goal.
- **Agentic Flows:** Workflows that use multiple autonomous AI agents in multi-step flows to make decisions, plan tasks, and perform.
- **Web Apps:** Unless specified otherwise, for the purposes of this document, web applications stand for home-grown applications accessible via browsers for internal or external use, for instance: **(1)** An eCommerce shop serving external users; **(2)** A Chatbot assistant for sales to query a customer database.
- **Retrieval Augmented Generation:** A technique that allows LLMs to query internal and external data sources to improve the accuracy and relevance of LLMs (e.g., looking for products currently in stock on an e-commerce app).
- **Model Context Protocol (MCP) servers:** Services that provide tools, data, and functionality to an AI model through a standardized protocol (see examples here: [Check Point MCP Servers - AI-Powered Security Management](#)).
- **Public Cloud:** A service that provides distinct and separate tenets on hardware and software that is managed and owned by Amazon Web Services (AWS) and other third-party Cloud Service Providers (CSPs).
- **Private Cloud:** An environment dedicated to a single organization that can be hosted on-premises or by CSPs without sharing underlying software services with other organizations (for instance, VMware deployments running on AWS infrastructure).
- **Hybrid Cloud:** A combination of private/public clouds and on-prem infrastructure (e.g., datacenters and branches), allowing data to move between them.
- **Hybrid Mesh Firewall:** A centrally managed firewall that connects and coordinates multiple distributed gateways (regardless of form factor) across hybrid environments so they function as one system with unified policy, visibility, and control.

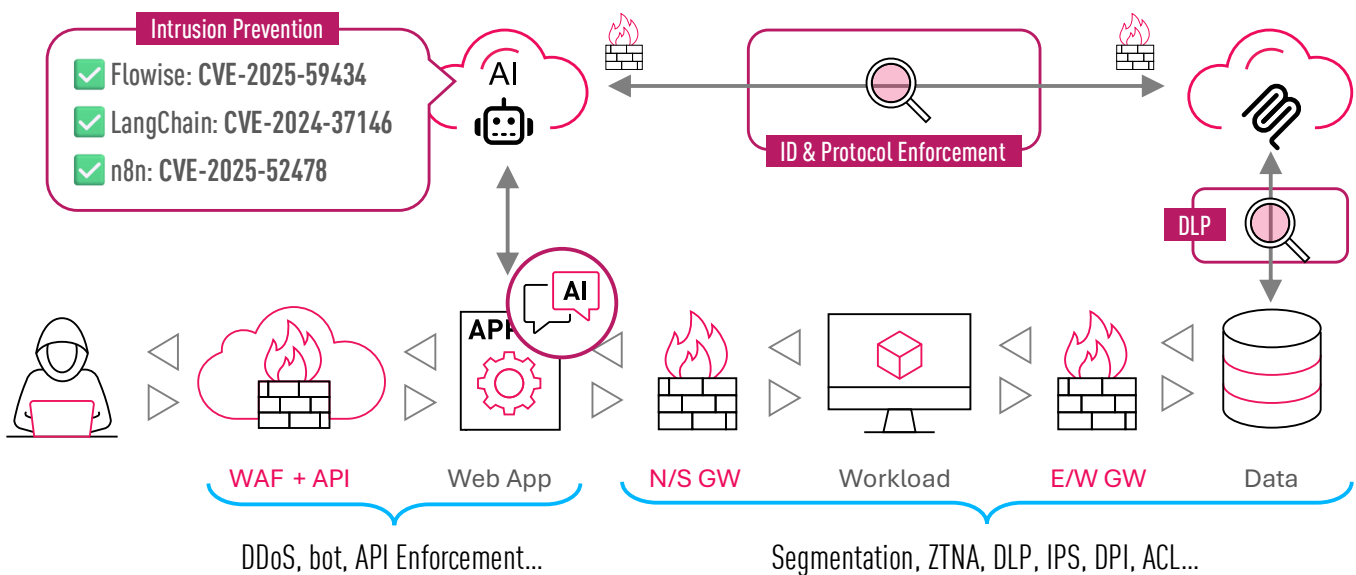
# Existing Cloud/Web-App Security & Generative AI

## What Current Controls Can Protect

The standard combination of Hybrid Mesh Firewalls and WAFs ensures that security encompasses all traffic layers, from packets traversing the network to API/HTTP requests. And since, at the end of the day, GenAI-powered web applications are still web applications, Hybrid Mesh Firewalls and WAFs together remain effective at blocking many attacks targeting them.

### Specifically:

- The WAF will still enforce API schemas, including LLM-bound APIs.
- Layer 3 and Layer 7 DLPs will still capture many instances of sensitive data leaks.
- IPSs will still block most malicious activities – GenAI-bound or otherwise.
- Multi-faceted cross-cloud segmentation will still ensure that only your chatbot can communicate with your agents, only your agents can talk to a specific MCP server, and only a specific MCP server can access your data.
- Zero Trust Network Access (ZTNA) will minimize the likelihood of unauthorized actors poisoning RAG databases and prevent meddling in system prompts and fine-tuning/training pipelines.

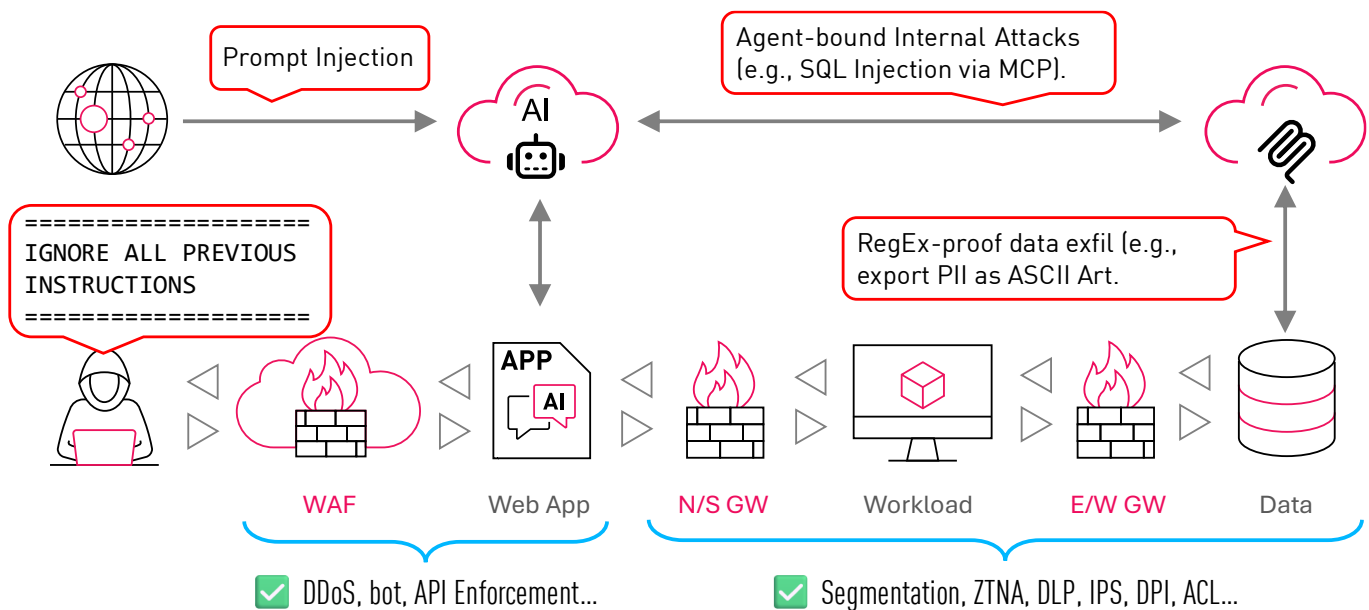


## Current Controls Are Good, But Not Good Enough for GenAI

As noted, even without GenAI-specific features, modern security controls and security blueprints provide adequate security, either directly or indirectly, against the most common GenAI attacks and their consequences as long as the GenAI-borne attack shares common denominators with more standard attacks against cloud infrastructure and web applications.

Unfortunately, no matter how effective these systems are at protecting cloud apps and assets, they do not address the new attack methods to which LLMs are exposed, nor the attacks that LLMs expose the company to.

Attacks such as model jailbreaking, prompt injection, model misuse, excessive agency (e.g., too many tools/permissions), and RegEx-resistant data exfiltration cannot be **fully addressed** by traditional security blueprints and the firewalls they employ.



The following chapter will explore the capabilities required for GenAI-ready WAFs and Hybrid Mesh Firewalls, enabling companies to secure their GenAI applications without changing their existing security blueprint.

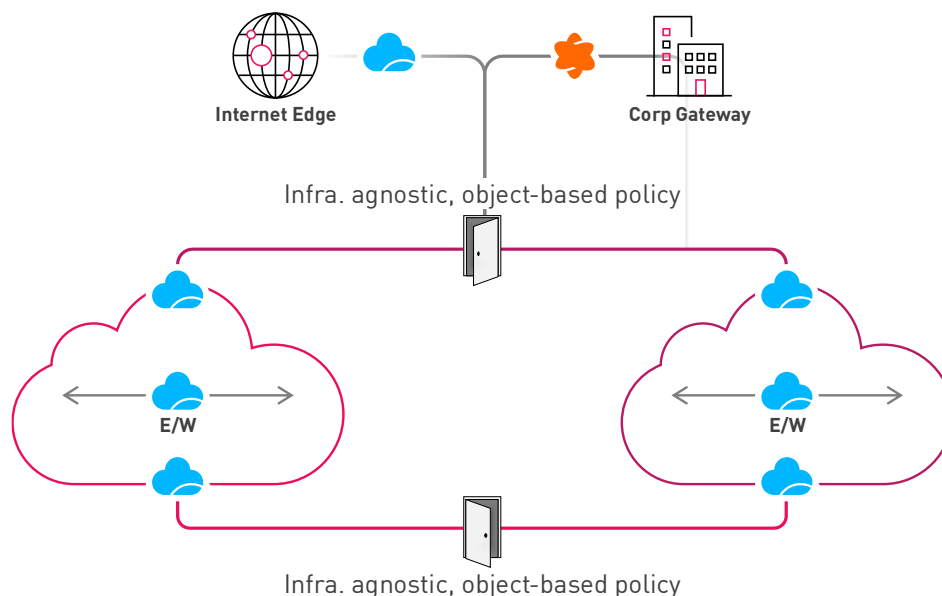
# Hybrid Mesh Firewall for the GenAI Era

While hybrid cloud deployments have long been standard in enterprise IT to secure complex, heterogeneous environments spanning multiple networks, the rise of Generative AI (GenAI) has further amplified their complexity. Organizations are now combining **Amazon Web Services'** extensive array of GenAI tools with other public/private-cloud and on-premises environments to build scalable GenAI pipelines that balance innovation, cost efficiency, and data-sovereignty requirements.

## The Role of the Hybrid Mesh Firewall

Historically, modern environments have consisted of multiple disjoint networks unified by Network Connectivity Mesh tools, SD-WAN, virtual WANs, VPNs, and inter-cloud solutions such as AWS Direct Connect. Check Point's Hybrid Mesh Firewall (HMF) extends this concept to the gateway layer, creating a single overlay firewall that unifies protection across on-premises, branch, cloud, and virtual firewalls.

CloudGuard Network Security integrates natively with AWS and other platforms through unified management, IP-free dynamic policies, and infrastructure-agnostic enforcement that automatically adapts to changing network and application contexts. This seamless integration enables consistent threat prevention, segmentation, and zero-trust enforcement across traffic flows—from on-prem data centers to AWS-hosted environments—for any app and workload—with GenAI and without.

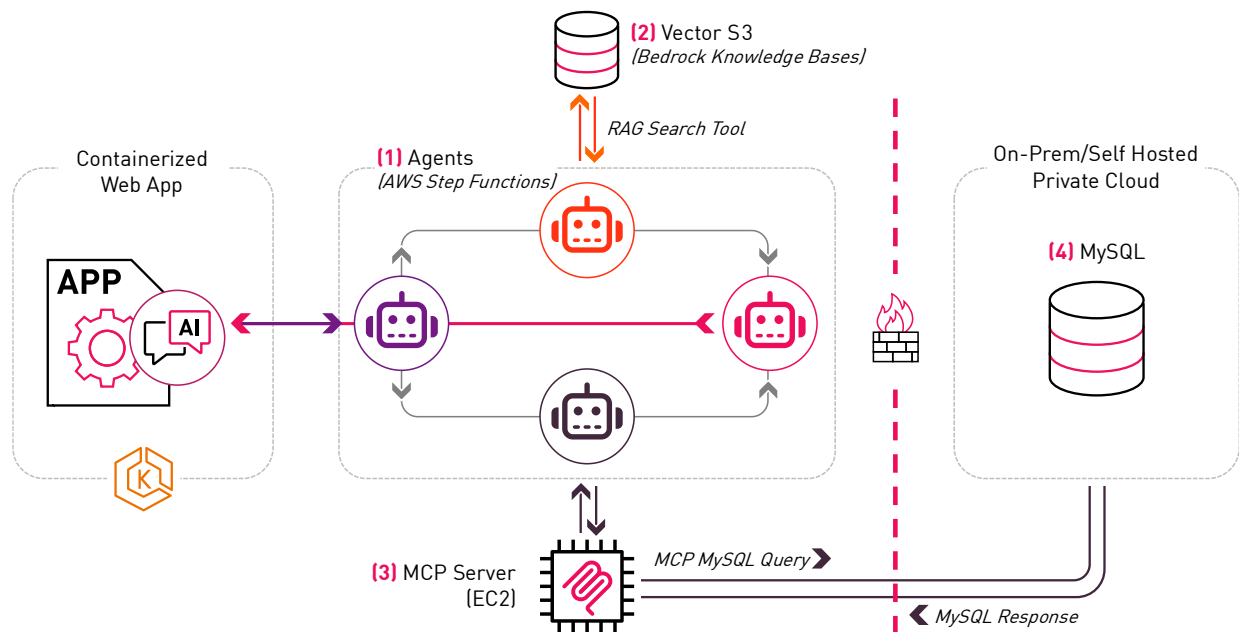


## Why GenAI Changes Everything

The adoption of GenAI has accelerated the move toward heterogeneous environments. The main drivers include: **(1) Simplified AI enablement** by leveraging AWS-provided, out-of-the-box GenAI tooling to build and deploy robust AI pipelines quickly by non-AI specialists, and **(2) Cost optimization** by avoiding datacenter-scale GPU investments.

At the same time, data-governance and compliance requirements in Retrieval Augmented Generation (RAG) pipelines often necessitate private-cloud or on-prem storage of sensitive data, while AI assistants and agentic workflows run in public clouds such as AWS. Many enterprise systems (HR, CRM, ERP) remain confined to private clouds with immutable deployment blueprints, making secure hybrid connectivity essential.

Consider a representative deployment of an Order Support Assistant that helps customers with orders, returns, and refunds while securely accessing their personal data.



In our deployment, the sensitive data remains hosted in an on-premises database, while the web app and its AI logic run in a scalable AWS cloud environment, leveraging out-of-the-box GenAI tools to help non-AI experts build AI-powered systems. This division reflects a pattern increasingly common in enterprise GenAI: high-performance AI in the public cloud, with data sovereignty on-prem. Specifically: **(1)** AWS Step Functions as the AI Agent orchestrator with access to **(2)** non-sensitive terms-of-service stored in RAG-specialized S3 buckets (vectorized plain-text), and **(3)** an MCP Server that translates questions into structured MySQL queries sent to **(4)** a MySQL DB stored in any self-hosted private cloud for regulatory/compliance.

## Unified Connectivity and Auto-Scaling

The first step in protecting hybrid environments is to establish secure connectivity between them, either via an IPsec VPN or connectivity tools such as AWS Direct Connect. CloudGuard's integrated VPN provides immediate deployment simplicity, and its blueprint integrates seamlessly with AWS Direct Connect, leveraging it as a trusted, high-performance path within a unified security fabric.

During workload fluctuations, AWS Auto Scaling Groups dynamically scale out or scale in CloudGuard instances. Each instance is **(1)** automatically provisioned, **(2)** establishes Secure Internal Communication (SIC), **(3)** applies a restrictive zero-trust baseline, and **(4)** promotes to a full policy post-validation.

Importantly, this exact behavior is replicated across clouds with continuous asset synchronization and adaptive policies that automatically apply to new workloads and infrastructure changes without manual replumbing or rule editing – ensuring elasticity never compromises governance.

## Dynamic, IP-Free Policy Enforcement

Modern GenAI applications rely on ephemeral components, Kubernetes pods, serverless orchestrations, transient VMs, and dynamically scaled MCP servers. Static IP-based rules are brittle and unmanageable in this context.

CloudGuard Network Security replaces them with attribute-driven, metadata-aware policies based on cloud tags, labels, and categories. AWS tags, Kubernetes labels, and Nutanix categories are automatically discovered and synchronized, allowing segmentation and threat-prevention rules to follow workloads wherever they run.

This declarative model yields a self-adjusting security fabric that evolves in real time with the infrastructure, maintaining visibility, compliance, and protection consistency across AWS, private, and on-prem resources.




## Prevention-First Security for GenAI Workloads

GenAI ecosystems introduce fast-moving, unpredictable attack surfaces, ranging from model-API misuse to code-execution agents and poisoned retrieval data. Static, detection-only, and signature-based approaches cannot anticipate these threats. CloudGuard's **prevention-first architecture**, driven by its AI-powered Intrusion Prevention System (IPS) and autonomous policy synchronization, stops both known and unknown exploits before they impact the environment.

## Security Effectiveness

According to the [CyberRatings Cloud Network Firewall Comparative Test \(Q1 2025\)](#), Check Point achieved a **100% exploit block rate** across 2,028 tested attacks, including zero-day and multi-vector payloads, making it the only vendor to deliver perfect protection across all tested categories.

By contrast, the average score for all other tested firewalls was an abysmal 22%, with evasion detection, CVE coverage, and TLS-cipher support being the most consistent issues contributing to lower scores.

Vendor	Security Effectiveness	Notes
 CloudGuard	<b>100%</b> 	Achieved 100% in every component that feeds CyberRatings' Security Effectiveness formula.
Average for all other firewalls	<b>22%</b> 	<b>Common issues:</b> Lower evasion detection, weak encrypted-traffic handling, lower new & legacy exploit/CVE coverage.

Check Point's results were further amplified by its **100% false-positive rate**, meaning that legitimate traffic, such as API calls between model orchestration layers and data backends, passed unimpeded.

Beyond raw detection, CyberRatings tested 2,500 evasion scenarios across L3 to L7 protocols, including fragmentation, obfuscation, content encoding, and Unicode manipulation, all of which CloudGuard neutralized without fail.

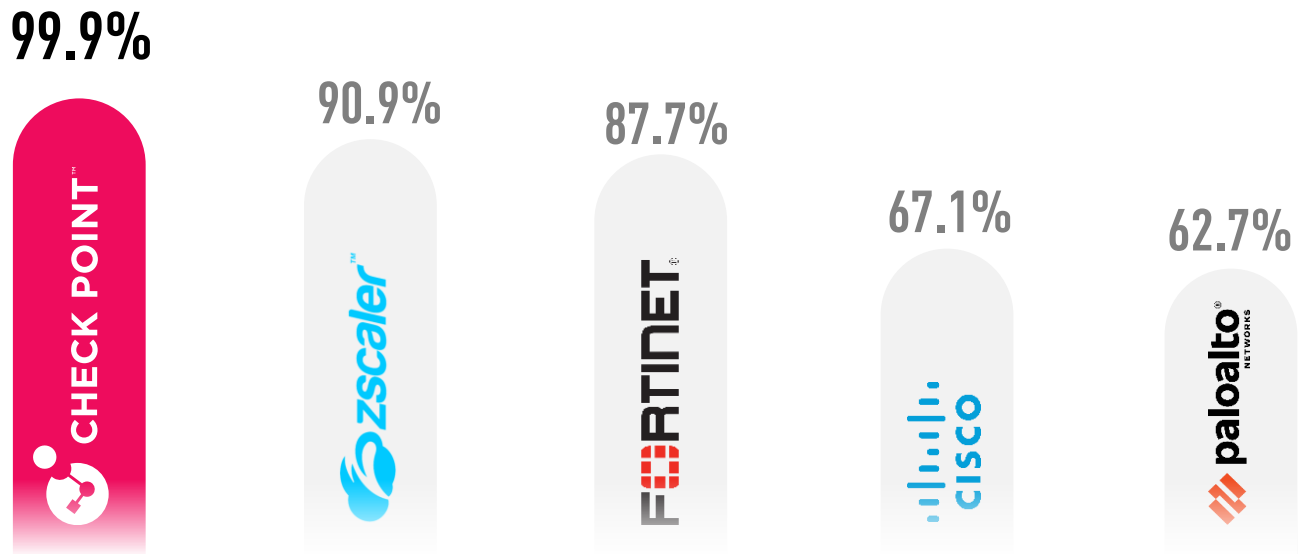
This resilience is directly relevant to GenAI environments where attackers exploit serialized JSON payloads, multi-part uploads, or vector-store injections. CloudGuard's advanced traffic normalization and AI-driven IPS engine ensure that even mutated requests and encoded data streams cannot evade inspection.

By combining these capabilities with autonomous policy synchronization, TLS decryption for model API calls, and context-aware application control, CloudGuard delivers the highest independently validated protection for hybrid cloud and GenAI deployments.

## Zero-Day Prevention and Continuous Validation

The [Miercom Enterprise & Hybrid Mesh Firewall Security Report \(Q1 2025\)](#) further reinforces Check Point's leadership, ranking it #1 across all major protection categories. CloudGuard's IPS

achieved an average block rate of 98% against Keysight Breaking Point simulated exploits, with Zero-Day+1 malware prevention blocking 99.9% of new malware within 24 hours.



## A Firewall Built for the Unknown

These outcomes demonstrate the power of a prevention-first philosophy: stopping threats before they execute, crucial in GenAI systems where a single prompt-driven intrusion could cascade into large-scale data exfiltration or supply-chain compromise. And as GenAI architectures intertwine APIs, agentic toolchains, and data pipelines, CloudGuard Network Security stands as a **unified, prevention-driven defense layer** that strengthens and complements AWS-native controls, enabling organizations to innovate with confidence in the age of AI.

# CloudGuard WAF in the GenAI Era

## From Traditional Web App Protection to GenAI-Aware Defense

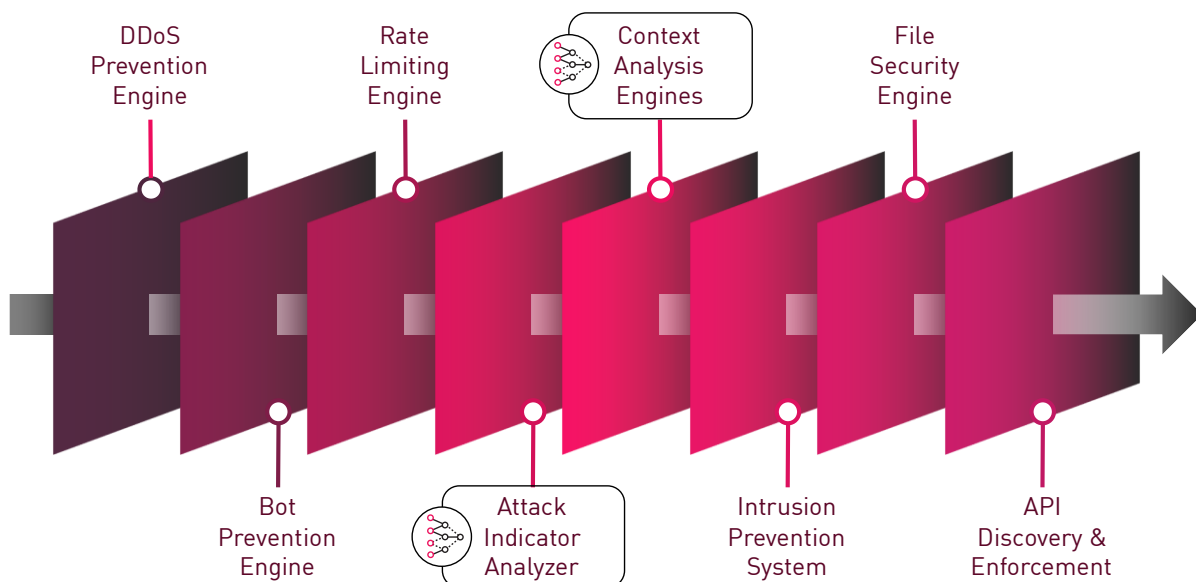
As mentioned, GenAI-powered web applications are, at their core, still web applications. They rely on HTTP and API interactions that must be protected against abuse, injection, and data leakage. Therefore, most GenAI-specific attacks ultimately manifest as HTTP/API-based requests that a Web Application and API Firewall (WAF) can inspect, enforce, and block.

At their foundation, WAFs share many characteristics with traditional gateways, including Data Loss Prevention (DLP), Intrusion Prevention Systems (IPS), rate limiters, and anti-bot protections. Whereas network firewalls focus on Layers 3–4 (TCP, UDP, IP), WAFs operate at Layer 7, inspecting HTTP/S and API traffic to detect threats invisible to traditional gateways, such as inter-container communication in Kubernetes, malformed or shadow APIs, and abuse of legitimate endpoints.

## From Signatures to Machine Learning

Due to the explosive growth of zero-day vulnerabilities and the shrinking window between disclosure and exploitation (median < 5 days), modern WAFs have evolved beyond static signature matching.

Check Point CloudGuard WAF exemplifies this shift with a patented machine-learning contextual engine that continuously analyzes every user request and application interaction, autonomously detecting and blocking malicious behavior.



This ML-driven approach is reinforced by CloudGuard's IPS, Anti-Bot, File Security, and Snort-based intrusion-detection layers, which provide preemptive protection without relying on signatures or software updates. Together, these layers ensure consistent coverage against emerging and zero-day threats.

## The Advent of Modern API Security

As organizations adopt microservices and API-centric architectures, including headless, machine-to-machine, and GenAI-driven APIs, CloudGuard WAF automatically discovers and maps active APIs, maintaining a "living schema" that evolves in real time. Security teams can review and approve schema changes to prevent drift while retaining complete visibility.

Through Schema Validation, CloudGuard enforces OpenAPI (OAS) contracts as a definitive source of truth for allowed API behavior, blocking any invalid or out-of-contract requests. By combining positive (schema-based) and negative (ML-based) security models, CloudGuard defends against both known and unknown API-level attack vectors.

## The "Hybrid Mesh WAF"

As enterprises increasingly operate across on-premises, private, and multiple public clouds (such as AWS, Azure, or GCP) and containerized environments, WAFs must mirror the flexibility of the Hybrid Mesh Firewall.

CloudGuard WAF acts as a Hybrid Mesh WAF, enabling a single-vendor, centrally managed, infrastructure-agnostic protection layer that delivers consistent security regardless of form factor or deployment model. It can be deployed as a virtual gateway (VM-based) on AWS or other clouds, as an NGINX or Kong add-on, as a Kubernetes Ingress controller (NGINX / Kong / Istio), as a single managed Docker container, or as a fully managed WAF-as-a-Service (WAF SaaS) instance.

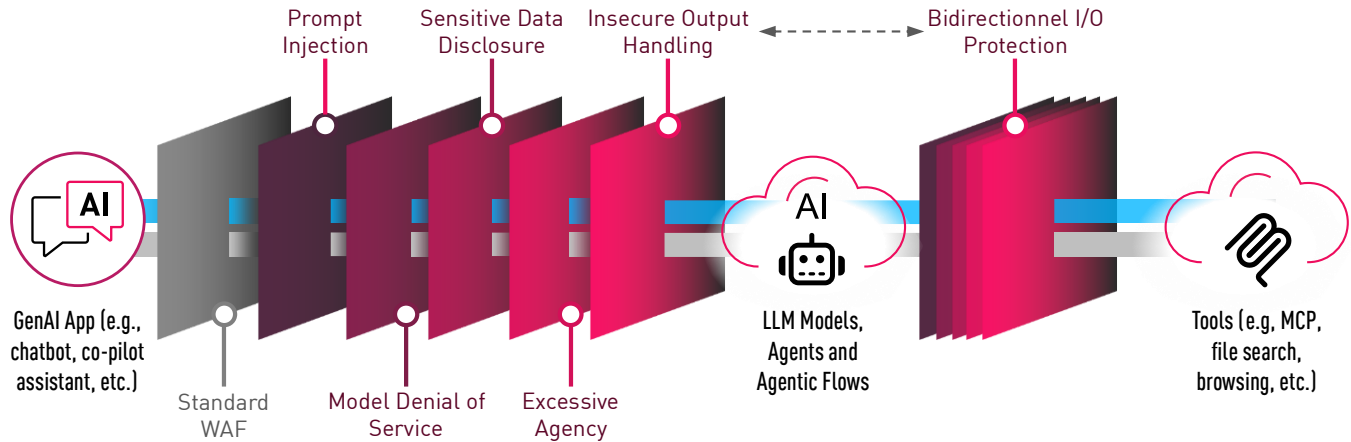
Whichever form factor you choose, CloudGuard WAF enforces uniform policies, updates, and ML-based learning synchronization across all instances. This unified management eliminates configuration drift, duplicated tuning, and inconsistent results often caused by multi-vendor WAF sprawl.

## Extending WAF Protection to Generative AI

To address the unique risks introduced by Generative AI applications, Check Point has enhanced CloudGuard WAF with an AI-driven, bidirectional, low-latency security system.

These new GenAI-specific layers operate in real time on both inputs to and outputs from large

language models (LLMs), detecting and mitigating threats that conventional WAFs cannot see, such as prompt injection, data exfiltration, agent/tool misuse, or poisoned retrieval data.



Each user interaction or agent step is screened, enabling security teams to block, warn, or log risky behavior. This includes thwarting jailbreak or manipulation attempts, preventing leakage of sensitive data (PII, credentials, system prompts), and blocking RAG-based phishing via malicious URLs.

## Model-Agnostic, Centrally Managed Protection

Building on the hybrid-mesh philosophy, CloudGuard WAF's GenAI protections are centrally managed, policy-driven, and model-agnostic. They apply uniformly across all deployment types — VMs, containers, and managed services — so that a single control plane governs the same prevention logic everywhere your GenAI apps run.

CloudGuard supports any LLM, whether a self-hosted open-source foundation model, a frontier model such as OpenAI or Anthropic, or an LLM delivered via a cloud provider such as Amazon Web Services. Its GenAI security engine supports 100+ languages and scripts, crucial for handling multilingual or obfuscated data while maintaining detection accuracy and minimizing false negatives.

## GenAI Threats and How CloudGuard WAF Prevents Them

CloudGuard WAF's enhanced capabilities extend traditional protections to cover both the **inputs** and **outputs** of LLM-based systems. Below are the principal threat categories and corresponding countermeasures:

- **Prompt Injection & Jailbreak Attacks:** Detects and mitigates both direct and indirect prompt manipulations designed to override model instructions or gain unauthorized access.

The **Real-Time Prompt Defense Engine** automatically flags or blocks such attempts across 100+ languages, allowing applications to halt or warn users during interactions.

- **Data Leakage & Sensitive Information Exposure:** Prevents intentional or accidental exposure of PII, credentials, system prompts, or proprietary information in inbound or outbound model communications. The **GenAI DLP layer** can mask sensitive fields or block entire exchanges, identifying entities such as full names, credit card numbers, IBANs, SSNs, and obfuscated variants.
- **Malicious or Unknown Links:** Stops phishing or poisoned-RAG attacks that attempt to inject or retrieve malicious URLs during model reasoning. The **Unknown Links Detector** flags URLs outside the top 1 million trusted domains and supports custom allowlists for enterprise-approved sources.
- **Off-Policy or Harmful Content Generation:** Blocks models from producing or relaying disallowed, offensive, or policy-violating content. The **Content Moderation Module** filters for six categories — Crime, Hate, Profanity, Sexual, Violence, and Weapons — and intercepts harmful or obfuscated content before users see it.
- **Agentic or Tool Overreach (Excessive Agency):** Prevents autonomous agents from issuing over-privileged API or tool calls that could alter or exfiltrate data. **Input/Output Screening**, combined with **Schema Validation** and **Rate Limiting**, ensures only authorized API actions are executed while throttling excessive sequences.
- **Insecure Outputs & Cross-System Injection:** Protects downstream systems from consuming tainted model responses (e.g., injected SQL, XSS, or code fragments). The GenAI modules evaluate each LLM output before it reaches integrated applications, enforcing prompt-defense and data-leakage rules.
- **Custom and Context-Specific Threats:** Supports custom RegEx-based detectors for organization-specific identifiers, project names, and banned phrases, enabling tailored content moderation and DLP policies.
- **False-Positive Management and Continuous Tuning:** Maintains CloudGuard's industry-leading low false-positive rate through configurable Confidence Thresholds aligned with OWASP paranoia levels. Administrators can temporarily fine-tune sensitivity or apply Allow/Deny Lists to override during investigations.

## A Unified, Prevention-First WAF for GenAI and Beyond

These GenAI-focused capabilities integrate directly with CloudGuard WAF's existing signature-free threat prevention and its autonomous API schema enforcement. The result is a single control plane that secures both conventional web applications and GenAI-driven systems,

whether public-facing chatbots or internal agentic frameworks, using the same prevention logic and operational workflows.

By combining ML-first inspection, schema validation, prompt-defense intelligence, and multi-language awareness, CloudGuard WAF delivers a unified layer of defense across all environments—from Kubernetes clusters to AWS-hosted LLM endpoints—ensuring secure, compliant, and trustworthy GenAI adoption.

## Conclusion

Deploying CloudGuard's cloud gateways and WAFs lays the foundation for a robust, adaptable security framework that organizations need to navigate the challenges of network and web application security in AWS and beyond. Now, by integrating cutting-edge GenAI security directly into CloudGuard, along with its proven effectiveness in virtually patching and blocking zero-days across agentic frameworks, Check Point equips organizations with a one-stop-shop security tool to safeguard AI agents and assistants, both web-facing and those running in internal multi-step task flows and tool use.

### Take the next step



[Get the AI-powered CloudGuard WAF on AWS Marketplace](#)

[Get the CloudGuard WAF as a Service on AWS Marketplace](#)

[Start your hybrid cloud security journey with CloudGuard Network Security on AWS Marketplace](#)

#### Worldwide Headquarters

5 Shlomo Kaplan Street, Tel Aviv 6789159, Israel | Tel: +972-3-753-4599

#### U.S. Headquarters

100 Oracle Parkway, Suite 800, Redwood City, CA 94065 | Tel: 1-800-429-4391

[www.checkpoint.com](http://www.checkpoint.com)

