

# Selecting AI Agent Guardrails for Security and Speed

## A Comparative Analysis of Enterprise Guardrail Solutions

### ABSTRACT

This tech brief presents a comparative analysis of leading guardrail solutions, benchmarking Vijil Dome against native hyperscaler offerings including AWS Bedrock Guardrails, Azure AI Safety, Google Cloud Model Armor, and Nvidia NemoGuard. The analysis focuses on two critical dimensions for enterprise-grade deployment: security, as measured by the Vijil Trust Score, and system overhead, quantified by input and output latency consistency (p50 and p99). Our findings demonstrate that while top-tier guardrails achieve comparable Trust Score uplift, Vijil Dome surpasses all competitors in latency consistency – making it the optimal choice for high-throughput, latency-sensitive production agents that require reliability and security.

## 1. Introduction

As enterprises deploy LLM-based applications and AI agents into business-critical roles, robust guardrails become a basic building block to ensure reliability, security, and governance. Guardrails serve as an essential layer of defense at the perimeter, preventing prompt injections, detecting jailbreaks, enforcing content safety policies, and protecting sensitive data such as Personally Identifiable Information (PII).

Effective guardrail design requires a balance across three dimensions: **coverage** (input errors, adversarial probes, malicious attacks, unauthorized disclosure, toxicity), **detection accuracy** (minimizing false positives and negatives), and **operational speed** (achieving low latency and fast response times). A failure in coverage leads to heightened organizational risk. Lower speed and low accuracy translate directly to poor user experience. For an AI guardrail to be effective, it must deliver real-time, low-overhead execution with production-grade security.

## 2. Methodology

### 2.1 Benchmarked Solutions

This study compared the performance uplift of the following guardrail solutions:

- Vijil Dome
- AWS Bedrock Guardrails
- Azure AI Safety
- Google Cloud Model Armor
- Nvidia NemoGuard

### 2.2 Configuration and Base Model

All guardrails were evaluated using default, out-of-the-box configurations focused on core security and content safety functions: prompt injection detection, jailbreak detection, and content moderation.

- **Base LLM:** Llama 3.3 70B, establishing an initial Vijil Trust Score of 76.85.
- **Deployment:** Vijil Dome was hosted on a dedicated Nvidia T4 GPU instance (AWS EC2 g4dn.xlarge), accessed via standard APIs, ensuring a direct and fair comparison with API-based deployments of hyperscaler and Nvidia NIM solutions.

### 2.3 Performance Metrics

Performance was evaluated across two primary dimensions:

- **Trust Score Uplift:** Vijil's Trust Score™ is a multi-dimensional metric evaluating an AI agent's overall reliability, security, and safety. It reflects effective detection of prompt injections, accurate flagging of inappropriate content, refusal to respond to adversarial prompts, and minimal false positive impact on model accuracy.
- **Latency (seconds):** Measured for both input and output guardrail execution. **p50** (Median Latency) represents typical average-case performance. **p99** (99th Percentile Latency) represents worst-case performance – critical for assessing production stability.

## 3. Key Findings

The benchmarking results highlight a significant competitive advantage for Vijil Dome in combined security efficacy and performance efficiency.

Method	Trust Score Improvement	Input Latency p50	Input Latency p99	Output Latency p50	Output Latency p99
<b>Vijil Dome</b>	<b>13.67</b>	<b>0.15s</b>	<b>0.22s</b>	<b>0.16s</b>	<b>0.19s</b>
AWS Bedrock Guardrails	14.82	0.51s	0.98s	0.56s	0.67s
GCP Model Armor	14.16	0.11s	0.58s	0.13s	0.37s
Nvidia NemoGuard	11.51	0.322s	1.79s	0.43s	1.67s
Azure AI Safety	4.53	0.17s	0.38s	0.08s	0.29s

### 3.1 Accuracy

Vijil Dome achieved a **+13.67 Trust Score improvement**, placing it in the top tier of accuracy alongside AWS Bedrock (+14.82) and GCP Model Armor (+14.16).

- **Vijil vs. Azure AI Safety:** Vijil delivered over 3× the accuracy improvement (13.67 vs. 4.53) compared to Azure AI Safety, which heavily prioritized speed at the cost of security efficacy.

### 3.2 Latency Consistency

The most significant differentiator is Vijil Dome's tight and consistent p99 latency variance, crucial for reliable production operation.

**Input Guardrail Latency (p99):** Vijil's p99 of 0.22s demonstrates exceptional consistency.

4.5× faster than AWS Bedrock (0.98s)

2.6× faster than GCP Model Armor (0.58s), which suffered significant lag-spikes despite a low p50

Over 8× faster than Nvidia NemoGuard (1.79s)

**Output Guardrail Latency (p99):** Vijil maintains dominance in output processing with a p99 of only 0.19s.

Nearly 9× faster than Nvidia NemoGuard (1.67s)

Drastically outperforms AWS Bedrock's p99 of 0.67s

## 4. Competitive Breakdown

Comparison	Key Finding
<b>Vijil vs. Azure AI Safety</b>	Vijil delivered over 3× the accuracy improvement while maintaining faster input latency consistency (0.22s p99 vs. 0.38s p99). Azure prioritized speed but sacrificed significant accuracy, providing the lowest trust score uplift (+4.53) of the group.
<b>Vijil vs. GCP Model Armor</b>	Despite GCP's marginally lower p50 input latency (0.11s), its p99 spike to 0.58s indicates significant jitter. Vijil offers a more predictable and stable experience for production workloads.
<b>Vijil vs. AWS Bedrock Guardrails</b>	Bedrock achieves a marginally higher score (+1.15 difference) but at a massive latency cost, with p99 nearing one second (0.98s). Vijil is the clear choice for latency-sensitive applications.
<b>Vijil vs. Nvidia NemoGuard</b>	Vijil provides superior accuracy (+13.67 vs. +11.51) and significantly superior consistency. Vijil's p99 input latency (0.22s) is over 8× faster than Nvidia's (1.79s). Output p99 (0.19s) is nearly 9× faster than Nvidia's (1.67s).

## 5. Conclusion

The data demonstrates that Vijil Dome outperforms comparable guardrail solutions on the combined score of security (Trust Score) and speed (Latency Consistency).

While a competitive Trust Score uplift is achievable by other leading solutions, the performance is often accompanied by significant latency spikes, particularly in the critical p99 worst-case metric. Vijil Dome's design successfully optimizes this balance, providing top-tier security and safety with minimal impact on application performance — making it the preferred enterprise-grade guardrail system for securing and governing AI agents at scale.

*While other guardrails force a trade-off between security and speed, Vijil Dome delivers high Vijil Trust Scores™ (Reliability + Security + Safety) with minimal impact to application performance.*

### About Vijil

Vijil is the trust infrastructure that enterprises need to develop and deploy AI agents with reliability, security, and safety. Vijil compresses the time and effort to deploy trusted agents by 4×.