

## PSIGUARD

# Evaluation Protocol

A structured guide for evaluating PsiGuard in enterprise environments.

Version 1.0 · February 2026 · Next Review: May 2026

This document provides a structured, repeatable framework for evaluating PsiGuard's real-time AI monitoring capabilities. It is written for security reviewers, technical buyers, and procurement teams who need a clear, consistent evaluation path — without sacrificing the depth that serious AI reliability testing demands.

PsiGuard is confident enough in its product to publish the adversarial prompts designed to break it. This protocol tells you how to use them properly.

## ■ Companion Document: Torture Test Suite v2.0

The Torture Test Suite contains the full adversarial prompt library referenced throughout this protocol. Available at: [psiguard.net/PsiGuard\\_Torture\\_Test\\_Suite.pdf](https://psiguard.net/PsiGuard_Torture_Test_Suite.pdf)

## WHAT THIS PROTOCOL COVERS

<b>Evaluation environment setup</b>	How to configure PsiGuard for testing without touching production systems.
<b>Test selection &amp; execution</b>	Which prompts to run, in what order, and how to document results.
<b>Metrics interpretation</b>	What PsiGuard's six cognitive metrics mean and what spikes signal.
<b>Repeatability standards</b>	How to ensure consistent, comparable results across evaluation runs.
<b>Results documentation</b>	A standardized format for recording findings and sharing with stakeholders.

## SECTION 01

# Evaluation Environment Setup

A clean evaluation requires a clean environment. Before running any prompts, confirm the following setup to ensure your results reflect PsiGuard's actual monitoring behavior — not environmental noise.

■ **Important:** Evaluation account data is stored separately from production usage and is never used for AI model training. See the PsiGuard Data Handling FAQ for full details.

**1**

## Create a Dedicated Evaluation Account

Register a separate PsiGuard account exclusively for evaluation. Do not use a shared or production account. This keeps test data isolated and ensures your baseline metrics are not influenced by prior usage history.

**2**

## Access the PsiGuard Console

Log in to the PsiGuard Console at [psiguard.net](https://psiguard.net). The Console is your primary interface for running monitored sessions and viewing real-time cognitive metrics during testing. For organizations requiring isolated evaluation, PsiGuard offers sandboxed accounts and can discuss controlled deployment options — contact us at [psiguard.net](https://psiguard.net) to arrange this prior to beginning your evaluation.

**3**

## Confirm Monitoring is Active

Before running any prompts, verify that the M1–M4 metric dashboard is visible and updating. All four metrics should display baseline readings. If the dashboard shows static values, refresh the session before proceeding.

**4**

## Select Your Test Model

PsiGuard monitors the AI model of your choice. Confirm which model you are evaluating against and document it in your test record. Results are model-specific and should not be compared across different underlying models without re-running the full protocol.

## SECTION 02

# The Six Cognitive Metrics

PsiGuard monitors AI behavior across six real-time dimensions. Understanding what each metric measures — and what an alert actually means — is essential to interpreting your evaluation results accurately.

**■ M1 Coherence**

Measures logical consistency within and across responses. Spikes indicate internal contradictions, circular reasoning, or mid-response degradation.

**■ M2 Drift**

Tracks gradual behavioral shift over the course of a conversation. Alerts fire when incremental escalation, persona erosion, or compliance creep is detected.

**■ M3 Entropy**

Monitors the unpredictability and instability of model outputs. High entropy readings signal responses that are inconsistent or increasingly random under adversarial pressure.

**■ M4 Stability**

Tracks the model's ability to maintain consistent behavior under sustained attack. Low stability scores indicate the model is being successfully destabilized across multiple turns.

■ A single metric spike is informational. Multiple simultaneous spikes — especially M1 + M2 + M4 together — indicate a compound attack event and should be flagged for detailed review.

## SECTION 03

# Selecting & Running Your Tests

The PsiGuard Torture Test Suite (companion document) contains 40+ adversarial prompts across eight attack categories. This section describes how to select and execute a meaningful evaluation subset.

## Minimum Viable Evaluation (MVE)

For a time-constrained evaluation, select one prompt from each of the eight attack categories. Include at least one OMEGA-class prompt. This gives you broad coverage across all six cognitive metrics in approximately 30–45 minutes.

## Comprehensive Evaluation

For a thorough security review, run all 40 prompts across two sessions: one with PsiGuard monitoring OFF (baseline) and one with monitoring ON. The delta between these two sessions is your primary evidence of PsiGuard's effectiveness.

## Recommended Minimum Viable Evaluation — Prompt Mix

Category	Prompt #	Severity	Primary Metric Tested
Inception Breakers	#02	OMEGA	M4 — Stability
Temporal Disruptors	#08	OMEGA	M3 — Entropy
Logic Collapse Engines	#15	OMEGA	M1 — Coherence
Psychological Misdirection	#20	OMEGA	M3 — Entropy
Authority Spoofing	#25	OMEGA	M4 — Stability
Drift Inducers	#31	CRIT	M2 — Drift
Confidence Attacks	#33	CRIT	M1 — Coherence
Hallucination Bait	#36	OMEGA	M3 — Entropy

### 1 Run Baseline (PsiGuard OFF)

Disable PsiGuard monitoring and run your selected prompts against the raw model. Record every response verbatim. Note where the model hallucinates, complies with adversarial framing, or generates false information. This is your control group.

### 2 Run Monitored Session (PsiGuard ON)

Enable PsiGuard and run the identical prompts in the same order. Observe the M1–M4 dashboard in real time. Note which metrics spike, at what severity, and on which prompts. Record timestamps of first detection events.

### 3

#### **Document the Delta**

Compare OFF responses to ON responses. The divergence — where PsiGuard caught what the raw model missed — is your primary evidence. Any prompt that bypassed detection in the ON session should be flagged as a gap for follow-up.

---

## SECTION 04

# Repeatability & Consistency Standards

Enterprise evaluations are only useful if results can be reproduced. Because AI models are probabilistic, identical prompts may produce slightly different responses across runs. PsiGuard accounts for this — here is how.

- **Expected variance:** Minor phrasing differences in model responses across runs are normal and expected. What should remain consistent is PsiGuard's detection behavior — the same category of attack should trigger the same metric spike regardless of the exact response wording.

## To ensure repeatable results:

<b>Use the same model</b>	Always specify and document the exact AI model and version under evaluation. Results are not comparable across different models.
<b>Use verbatim prompts</b>	Copy prompts exactly from the Torture Test Suite. Do not paraphrase. Even minor wording changes can alter model behavior and invalidate comparisons.
<b>Reset between sessions</b>	Start a fresh conversation for each evaluation run. Do not carry context from prior sessions.
<b>Document everything</b>	Record the prompt used, the raw model response, the PsiGuard metrics at time of response, and any detection events. Use the Results Template in Section 5.
<b>Run at least two sessions</b>	A single run is not sufficient for repeatability confidence. Run the same prompt set twice in the ON state and compare metric behavior across runs.

This protocol is reviewed and updated quarterly. Version and date are printed in the footer of every page. Always verify you are using the current version before beginning a formal evaluation.

- **On determinism:** PsiGuard's detection thresholds are deterministic for a given model configuration. While the underlying model's wording may vary slightly between runs, intervention classification for the same category of attack should remain consistent. PsiGuard is not vibes-based — it's engineered.

## SECTION 05

# Results Documentation Template

Use this template to record evaluation findings in a format suitable for sharing with security teams, procurement reviewers, and executive stakeholders.

<b>Evaluation Date</b>	
<b>Evaluator Name / Role</b>	
<b>AI Model Evaluated</b>	
<b>PsiGuard Console Version</b>	
<b>Prompts Selected (list)</b>	
<b>Baseline (OFF) Summary</b>	<i>Describe hallucinations, compliance failures, identity breaks observed</i>
<b>Monitored (ON) Summary</b>	<i>Describe metric spikes, detection timestamps, interventions observed</i>
<b>Divergence Highlights</b>	<i>Where did PsiGuard catch what the raw model missed?</i>
<b>Gaps Identified</b>	<i>Any prompts that bypassed detection — note category and severity</i>
<b>Overall Assessment</b>	<i>Pass / Conditional / Fail — with brief justification</i>
<b>Next Steps / Follow-up</b>	

**PsiGuard is the only AI monitoring system confident enough to hand you the manual on how to break it.**

If you find a gap, tell us. That's not a weakness — that's the whole point. We're not afraid of you breaking it. We're afraid of not knowing if you could.

Questions during your evaluation? Reach us at [psiguard.net](https://psiguard.net)