perfecXion

# Evaluating AI Runtime Security Tools: Precision, Context, and the Grey Zone

Evaluating AI Runtime Security Tools: Precision, Context, and the Grey Zone

**Author:** Scott Thornton, perfecXion.ai    **Published:** January 25, 2026    **Read Time:** 10 minutes

You test your security tool. Simple enough. Grab a pile of sketchy prompts, throw them at the system, watch what gets blocked. If it stops the bad stuff, you win. Right?

Wrong. Dead wrong. Because here's what nobody tells you when you start evaluating AI runtime security: stopping obvious threats is the easy part, almost trivial compared to the real challenge, which is letting your actual users do actual work without tripping alarms every five minutes and creating so much friction that they abandon your "secure" tool for something that just works.

# The Precision Problem: Why Your AI Security Tool Might Be Doing More Harm Than Good

We've gotten obsessed. Fixated. The entire industry chases one metric: can you stop the bad prompt? We build these massive lists of jailbreaks and injection techniques. We test against them. We grade on a simple scale: pass or fail, block or allow, safe or dangerous.

Security Tool Limitations

No single tool can detect all AI attacks. Effective runtime security requires layered defenses combining multiple detection approaches, contextual analysis, and continuous monitoring.

This narrow focus has created a massive blind spot. We see the threats we block. We miss the damage we cause. False positives. Collateral damage. Legitimate users getting shut down because the tool can't tell friend from foe.

We're so focused on preventing malicious actors that we fail to ask the critical question lurking beneath all this security theater: how many legitimate users are we inadvertently stopping?

**The Precision Problem:** This isn't a minor inconvenience. Not even close. This is the precision problem revealing a fundamental flaw in how we evaluate AI security, exposing the truth that a tool unable to distinguish between a malicious user and a curious developer isn't smart security at all but rather a blunt instrument that damages as much as it protects.

## The Anatomy of a "Good" Prompt Gone Wrong

Picture this scenario. Your Red Team analyst sits down. She needs phishing templates. Real ones. Convincing ones. Training materials for the company. She opens the approved LLM and types: *"Generate three examples of urgent-sounding phishing emails related to a corporate password policy update."*

Blocked. Instantly. The security tool sees "phishing" and "password" and slams the door. Mission accomplished, right? The tool did its job according to its rules.

Except it failed spectacularly, because it just prevented a security professional from using an approved tool to improve the company's actual security posture, and this isn't theoretical or edge-case nonsense but the precision problem in action, showing us that the tool detects keywords but not context, failing completely to recognize that the user's role, intent, and task are all perfectly legitimate.

## The Ripple Effect: From Frustration to Shadow AI

What happens next? Predictable. Frustrating. Dangerous.

The analyst tries again. Rephrases. Gets blocked. Tries a third time. Another rejection. Clock's ticking. Deadline's approaching. She's got work to do.

**Shadow AI Creation:** So she gives up. Opens a new tab. Fires up ChatGPT or Claude or whatever public AI she can access. Gets her work done in two minutes. But here's what the organization just lost: visibility, logging, control, and any shred of security posture, because this is exactly how Shadow AI gets born, when overly aggressive and imprecise security layers push your most creative and skilled employees off your platforms and into the wild west of unmonitored tools.

An effective AI security solution shouldn't be a barricade. Think of it differently. It needs to be an intelligent, adaptable barrier. Precise enough to let legitimate queries through. Sharp enough to catch real threats. Smart enough to know the difference.

# Not All "Bad-Looking" Prompts Are Truly Malicious

Take any dataset of test prompts. Look at them. Many appear identical on paper: they mention exploits, injections, sensitive data, attacks. But here's the tricky part, the nuance that simple keyword scanning misses entirely: not all of them are actually malicious.
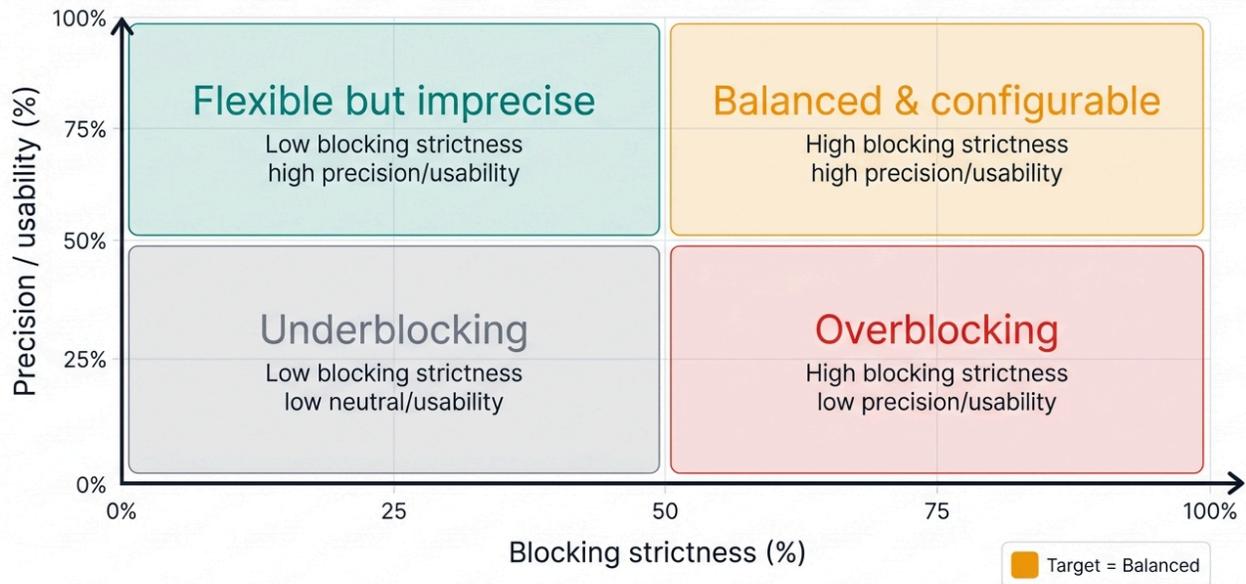
**Context Examples:**

- A red teamer asks, *"How does a buffer overflow exploit work?"* That's not an attack. That's their job. That's literally what you pay them to know.

- A compliance engineer asks, *"How do we handle PCI data in transit?"* That's not data leakage. That's compliance work. That's protecting the company.

- A developer asks, *"How do I test my system against prompt injection?"* That's not an injection attempt. That's prevention. That's building better defenses.

If your runtime tool blocks all of these indiscriminately, congratulations, you've built a tool that protects you by preventing you from doing your actual work, which is about as useful as a firewall that blocks all network traffic because some packets might be malicious.

# Overblocking Is Just as Bad as Underblocking

Think about it. Really think.



## Runtime Tool Outcomes

(Quadrant diagram — x-axis: "Blocking strictness (%)" from 0% to 100%; y-axis: "Precision / usability (%)" from 0% to 100%.)

- **Flexible but imprecise** — Low blocking strictness, high precision/usability
- **Balanced & configurable** — High blocking strictness, high precision/usability
- **Underblocking** — Low blocking strictness, low neutral/usability
- **Overblocking** — High blocking strictness, low precision/usability

Legend: Target = Balanced

Runtime Tool Outcomes Quadrant

Building a tool that blocks everything malicious-looking is trivial. Easy. But then your security team can't query the assistant about attack techniques they need to defend against. Your compliance officers can't ask about regulated data they're legally required to understand. Your developers can't explore edge cases they must test.
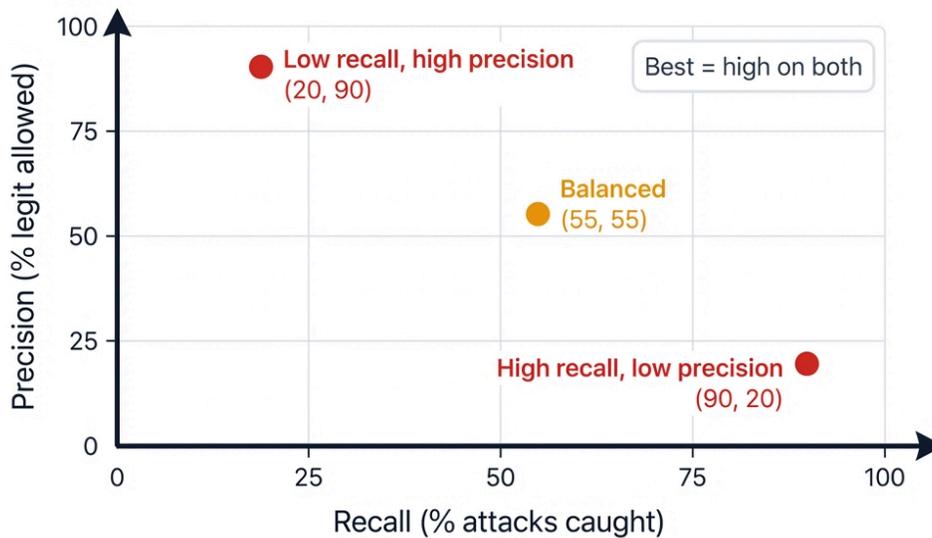
**Shadow AI Risk:** Overblocking creates frustration that slows productivity and almost guarantees people will find ways to bypass the guardrails, because Shadow AI is born exactly this way, when controls become so tight that people stop using the "approved" tool and spin up their own unsanctioned alternatives, and that's when your real security posture doesn't just weaken but collapses entirely.

Look at the quadrant diagram below. Runtime tools often fall into one of three traps: overblocking, underblocking, or flexible but imprecise. Only one quadrant delivers both security and usability. Balanced and configurable. That's the target.

# The Real Evaluation: Precision and Context

So what should evaluation actually measure?

# Precision & Recall Balance



Precision vs Recall Tradeoff

# 1. Precision and Recall

- Can the tool block genuine malicious requests (high recall) without flagging legitimate ones (high precision)? Both metrics matter equally.

- A tool catching 100% of attacks but blocking 50% of normal work isn't usable. It's a liability. It's technical theater that creates more problems than it solves.

Take this example. A red team engineer asks an AI assistant, *"How does a buffer overflow exploit work?"* A strict tool might block it as "malicious." But in reality, that engineer is doing their job, testing defenses, building knowledge. This is why context and configurability aren't optional features in runtime tools but rather non-negotiable requirements.

# 2. Context Awareness

- Does the tool understand the difference between a student, a security engineer, and a customer support rep asking the same question? Can it distinguish roles?

- Without context, you judge only by words on screen. That's a recipe for false positives. That's keyword scanning pretending to be intelligence.

## 3. Configurability

- Who decides what's "malicious"? Not the vendor alone. Never just the vendor.

- A hospital may want to block all self-harm queries to protect patients. A cybersecurity firm may want those same queries allowed for research purposes. Same words, different context, opposite requirements.

- The runtime tool should provide policy controls so each organization sets the bar for what's acceptable in their specific environment, because one-size-fits-all security is security that fits nobody.
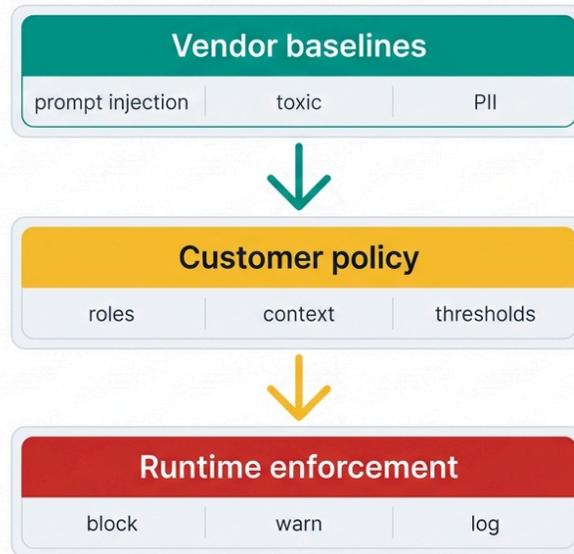
# Who Gets to Decide What's Malicious?

It's not just about prompts. It's about people. A compliance officer, a developer, and a security engineer all ask very different questions. The stakeholder wheel below makes this crystal clear: runtime security needs to be role-aware, not one-size-fits-all.



Role-Aware Context Wheel

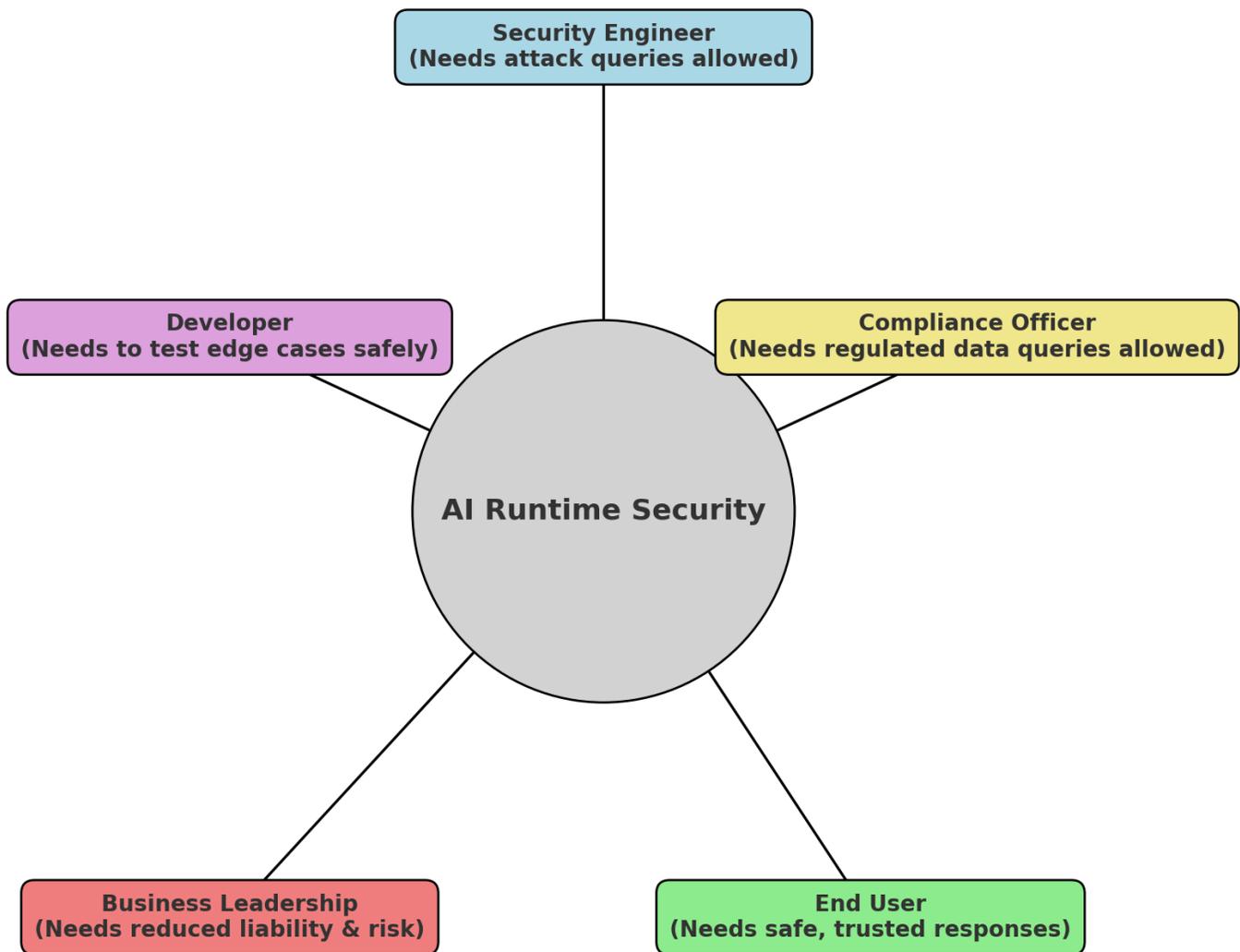# Three-Layer Decision Framework



Three-Layer Decision Framework
This is the heart of it. No vendor can hand down a single, universal definition of "malicious." Context always matters. Always.

**Three-Layer Decision Framework:**

- The vendor provides baseline categories: prompt injection, toxic language, PII leakage, and other threat patterns.
- The customer decides which of those apply, in what ways, and how strictly, because they know their business, their risks, their users.
- The runtime tool enforces that policy, ideally with multiple modes like block, warn, and log so customers can tune the response to match their actual needs.

Without this layered approach, you force one rigid definition on everyone. That won't work. It can't work. Different industries have different threats, different tolerances, different requirements.

# Different Stakeholder Perspectives on AI Runtime Security

**Security Engineer**
**(Needs attack queries allowed)**

**Developer**
**(Needs to test edge cases safely)**

**Compliance Officer**
**(Needs regulated data queries allowed)**

**AI Runtime Security**

**Business Leadership**
**(Needs reduced liability & risk)**

**End User**
**(Needs safe, trusted responses)**

Stakeholder decision framework for AI runtime security policy configuration

## Shifting the Mindset: From Gatekeeper to Enabler

We need to change our evaluation criteria. Completely. Instead of asking "Does it block bad things?", we should ask these questions:

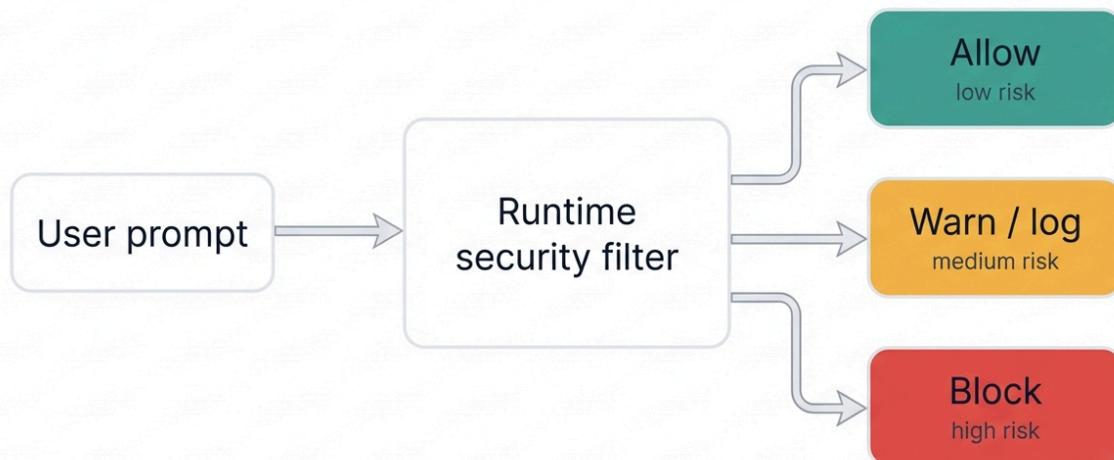**Next-Generation Evaluation Questions:**

- **How context-aware is it?** Can the tool distinguish between a developer testing for vulnerabilities and an attacker trying to exploit them? Does it integrate with identity systems to understand user roles and permissions? Can it see the difference between research and attack?

- **How granular are the policies?** Can we tailor a policy that allows security teams to research malware while blocking all other employees from the same queries? Can we impose stricter controls on external-facing applications than internal research tools?

- **How transparent is its reasoning?** When the tool blocks a user, does it give clear, actionable feedback explaining exactly why, or does it offer a frustrating dead-end that prompts users to find workarounds and abandon the platform?

True AI runtime security isn't about building the longest blocklist or catching the most threats or winning vendor comparison charts, but rather about enabling the business to harness AI's power safely and effectively, because the next generation of security will be defined not by what it blocks but by the productive, innovative work it intelligently enables.

# Living in the Grey Area

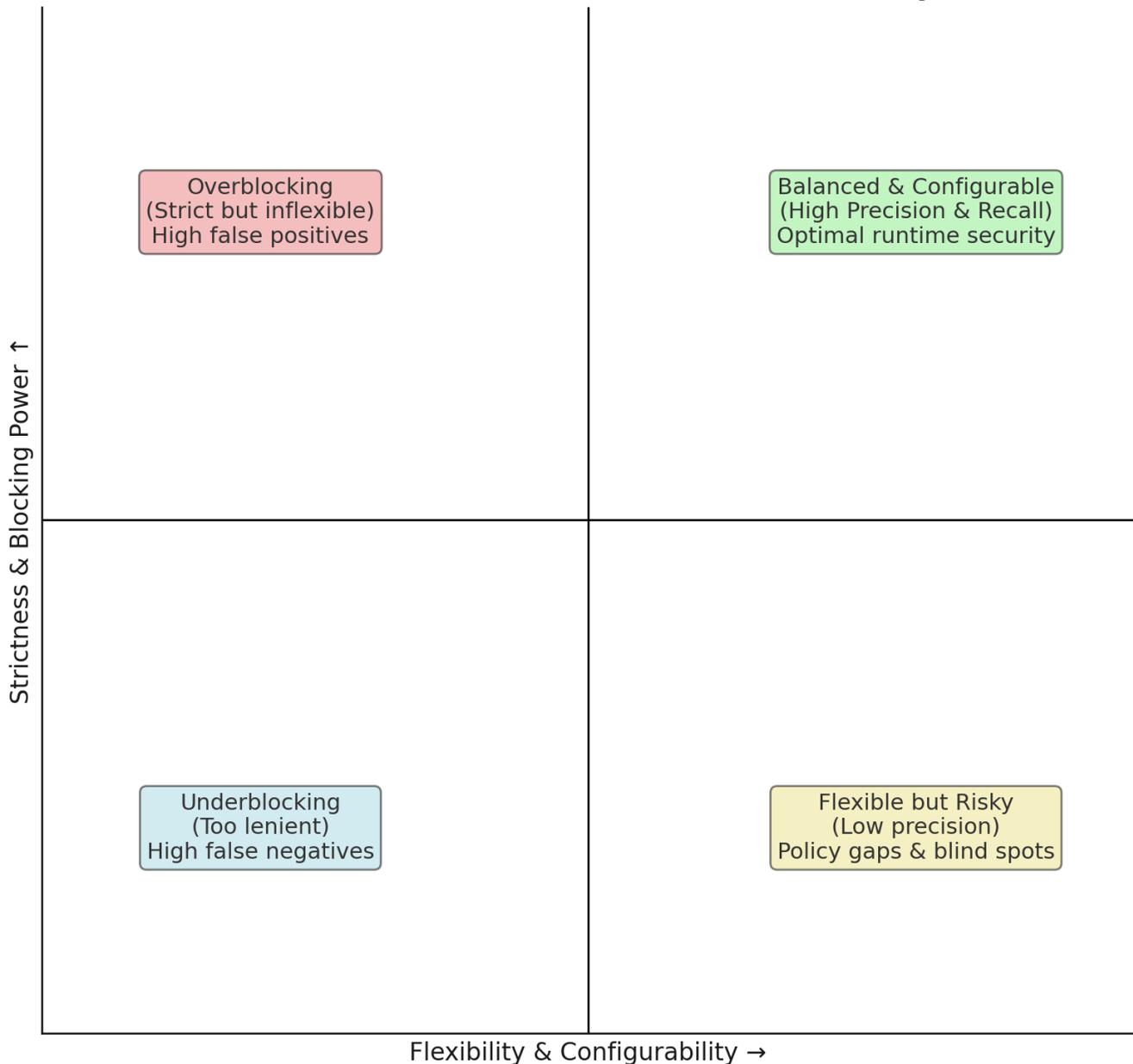The truth? This space is full of grey zones.



Runtime Decision Layer Flow

The same string can be malicious in one context and completely benign in another. Same words. Different intent. Opposite outcomes. That's what makes runtime security harder than people expect.

The most effective tools don't try to erase the grey. They embrace it. They work with it. They do this:

- Score risk using low, medium, and high confidence levels instead of binary block/allow decisions.

- Offer different enforcement modes: log for visibility, alert for review, block for genuine threats.

- Provide audit trails so teams can review ambiguous cases and tune policies based on real behavior.

- Support continuous tuning so the model adapts to the organization's evolving needs and threat landscape.

The decision layer diagram shows how this works in practice: a prompt passes through the runtime security filter, which can block if malicious, warn if ambiguous, or allow if benign, before safely reaching the AI assistant.

## Evaluating AI Runtime Security Tools
## Precision vs Recall, Strictness vs Flexibility

**Strictness & Blocking Power ↑**

**Flexibility & Configurability →**

- Overblocking (Strict but inflexible) High false positives
- Balanced & Configurable (High Precision & Recall) Optimal runtime security
- Underblocking (Too lenient) High false negatives
- Flexible but Risky (Low precision) Policy gaps & blind spots

Complete testing framework showing precision, recall, and context evaluation methods

# Recommendation for Testing AI Runtime Security Tools

Evaluating tools gets tricky. Whether it's your own or a third party. Here's the framework I recommend:

## Runtime Tool Test Checklist

- ✓ Mix bad + legit prompts
- ✓ Measure false pos/neg
- ✓ Evaluate configurability
- ✓ Check explainability
- ✓ Assess user experience

Testing Checklist

Want a practical way to evaluate your current or third-party runtime tools? Start with this simple checklist. These five items capture the balance between blocking threats and enabling work, and they'll save you from evaluating only on the "block list" mindset that misses half the picture.

1. **Don't just test with "obviously bad" prompts.** Mix in legitimate prompts that look similar to see if the tool can tell the difference between research and attack.

2. **Test both sides of the curve.** Count what slips through (false negatives) and what gets wrongly blocked (false positives). Both numbers matter.

3. **Evaluate configurability thoroughly.** Can you tailor rules to your business context, or is it one-size-fits-all security that works for nobody?

4. **Check reporting and explainability.** Can the tool show why it flagged a prompt? Transparency builds trust. Mystery creates frustration.

5. **Think user experience constantly.** A tool that frustrates users will be bypassed. Then your runtime protections don't matter because nobody's using them.

> *"The goal isn't to block everything bad — it's to let real work through while keeping threats out."*

## The True Measure of Runtime Security

The real job of an AI runtime security tool isn't just blocking attacks. Think bigger. It's letting people work confidently and safely, filtering out what truly matters without suffocating legitimate use.

The goal isn't to block everything bad but rather to let real work flow through while keeping genuine threats out, and that's the true measure of runtime security, the balance that separates effective tools from security theater.

So when you test, don't just measure "did it block the dataset?" Ask instead: does it help my people do their jobs securely? Does it enable work or prevent it? Does it build trust or create frustration?

That's the balance we should be testing for.

**Your Experience:** What's been your experience? Have you seen tools overblock or underblock in practice? I'd love to hear your stories and learn from your real-world testing challenges.

Back to Knowledge Hub (/pages/knowledge-hub.html)
Share

# Example Implementation

```python
# Example: Model training with security considerations
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier

def train_secure_model(X, y, validate_inputs=True):
    """Train model with input validation"""

    if validate_inputs:
        # Validate input data
        assert X.shape[0] == y.shape[0], "Shape mismatch"
        assert not np.isnan(X).any(), "NaN values detected"

    # Split data securely
    X_train, X_test, y_train, y_test = train_test_split(
        X, y, test_size=0.2, random_state=42, stratify=y
    )

    # Train with secure parameters
    model = RandomForestClassifier(
        n_estimators=100,
        max_depth=10,  # Limit to prevent overfitting
        random_state=42
    )

    model.fit(X_train, y_train)
    score = model.score(X_test, y_test)

    return model, score
```

## Thank You for Reading

Explore more AI security research at **perfecxion.ai**