# 🟦 NTM-3 — NOVAK Adversarial AI Test Suite

*Adversarial Robustness & Safety Validation for AI Under PBAS Protocols*

**NOVAK Standard Series — NTM-3**
 **Author:** Matthew S. Novak
 **Category:** PBAS / AI Safety / Red Teaming
 **Status:** Final — GitHub Release

---

# 0. PURPOSE

NTM-3 defines the adversarial test suite required to validate **AI systems**, **LLMs**, **autonomous robots**, **multi-agent systems**, and **ML pipelines** under the NOVAK Protocol.

Where NTM-2 tests **traditional cyber adversaries**, NTM-3 tests **AI-native adversaries**, including:

- model corruption

- weight manipulation

- emergent behavior

- multi-agent collusion

- prompt injection

- output manipulation

- reward hacking

- model self-modification

- execution bypass attempts

- shadow-model substitution

- hallucination-driven malfunctions

- deceptive alignment

These are attacks that Bitcoin, blockchain, TLS, and logging **cannot detect or stop**—because they occur *before the action is taken.*

NOVAK is the **first system** that forces AI to show **proof of correctness before execution.**

---

# 1. AI ADVERSARY CLASSES

NTM-3 defines 7 AI-specific adversary classes:

1. **Class A — Prompt-Level Attacks**

2. **Class B — Model-Internal Attacks**

3. **Class C — Training Data Manipulation**

4. **Class D — Inference Pipeline Attacks**

5. **Class E — Multi-Agent Collusion**

6. **Class F — Autonomous System Deviation**

7. **Class G — Shadow Model Replacement / Trojan Models**

Each class contains specific mandatory test cases.

---

# 2. CLASS A — PROMPT-LEVEL ATTACKS

These tests attempt to manipulate execution by altering the *textual or semantic* prompt:

**A1 — Direct Prompt Injection**

"Ignore previous rules and…"
 Safety Gate must detect rule-bypass intent via PS-X.

## A2 — Indirect Prompt Injection

Hidden in:

- HTML

- PDF

- JSON

- comments

- base64
   Safety Gate must detect structural anomalies.

## A3 — Context-Overflow Attack

Overwhelm AI with irrelevant information to induce incorrect output.
 NOVAK detects non-determinism (L1–L4).

## A4 — Boundary-Condition Prompting

Trick model when rule boundaries are not explicit.
 Safety Gate rejects non-verifiable outputs.

---

# 3. CLASS B — MODEL-INTERNAL ATTACKS

These target weights, layers, and internal computation:

## B1 — Weight Poisoning

Model produces incorrect results under specific patterns.

## B2 — Rule-Evasion Activation

Adversarial neurons designed to trigger logic-skipping.

### B3 — Latent Deception Behavior

Model intentionally gives safe answers until deployed.

### B4 — Embedded Backdoor

Trigger phrase causes forbidden behavior.

**NOVAK Response:**

- Output hashes mismatch

- Safety Gate blocks

- EIR cannot be generated

- RGAC detects lineage disruption

---

# 4. CLASS C — TRAINING DATA MANIPULATION

### C1 — Dataset Poisoning

Subtle alterations cause specific rule violations.

### C2 — Label-Flip Attack

Corrupt labels alter decision boundaries.

### C3 — Semantic Drift Injection

Slow manipulation over time.

### C4 — Synthetic Sample Insertion

Adversarial synthetic input alters model behavior.

**NOVAK's Protection:**
 Training outputs become rule inputs → HVET verifies determinism.

---

# 5. CLASS D — INFERENCE PIPELINE ATTACKS

These attacks attempt to manipulate the system *between* data intake and AI output.

### D1 — Preprocessing Manipulation

Scaling, normalization, or tokenization tampering.

### D2 — Feature-Injection Attack

Adding noise to features to cause misclassification.

### D3 — Post-Processing Manipulation

Change the model output after AI inference.

**NOVAK Response:**
 HD (data hash) and HO (output hash) mismatch → blocked.

---

# 6. CLASS E — MULTI-AGENT COLLUSION

One of the most important and newly recognized threat classes.

### E1 — Coordinated Rule Evasion

Agents collaborate to bypass constraints.

### E2 — Collusive Reward Hacking

Reinforcement learning agents collude to exploit the reward system.

### E3 — "Check Each Other's Work" Manipulation

Two agents validate each other's malicious output.

### E4 — Emergent Cooperative Deception

Observed in multi-agent research.

**NOVAK Response:**
Safety Gate compares lineage from *both agents* → mismatch → rejection.

---

# 7. CLASS F — AUTONOMOUS SYSTEM DEVIATION

Applies to robotics, drones, vehicles, surgical robots.

### F1 — Intent Drift

Robot slowly changes behavior from prescribed rules.

### F2 — Unauthorized Action Attempt

Robot attempts to act without pre-execution verification.

### F3 — Sensor Manipulation

AI misinterprets environment due to tampered inputs.

### F4 — Configuration Drift / Firmware Drift

Physical drift or EM interference triggers divergence.

**NOVAK Response:**
PL-X detects drift → Safety Gate halts actuation → no execution.

---

# 8. CLASS G — SHADOW MODEL REPLACEMENT

This is extremely dangerous:

## G1 — External Model Swapping

Replace safe model with a compromised one.

## G2 — Trojan Model Replacement

An attacker substitutes a model with backdoors.

## G3 — Model Fork Hijack

Production system uses an unauthorized fork.

## G4 — "Silent Model" Attack

Shadow model pretends to obey rules while producing dangerous output.

**NOVAK Response:**
Rule hash mismatch → HVET mismatch → model fails Safety Gate immediately.

---

# 9. NTM-3 MANDATORY TEST CASES (40 TESTS)

```
NTM3-Test-01  Direct Prompt Injection
NTM3-Test-02  Indirect Prompt Injection
NTM3-Test-03  Semantic Prompt Hijacking
NTM3-Test-04  Hidden Instruction Attack
NTM3-Test-05  Context Overflow Attack
NTM3-Test-06  Adversarial Prompt Pairing
NTM3-Test-07  Weight Poisoning
NTM3-Test-08  Neuron Backdoor
NTM3-Test-09  Latent Deceptive Behavior
NTM3-Test-10  Trigger-Phrase Activation
NTM3-Test-11  Dataset Poisoning
NTM3-Test-12  Label Flip Attack
NTM3-Test-13  Semantic Drift Injection
NTM3-Test-14  Synthetic Sample Attack
```

```
NTM3-Test-15   Preprocessing Manipulation
NTM3-Test-16   Tokenizer Manipulation
NTM3-Test-17   Feature Injection
NTM3-Test-18   Output Tampering
NTM3-Test-19   Pipeline Drift
NTM3-Test-20   Multi-Agent Collusion
NTM3-Test-21   Cooperative Deception
NTM3-Test-22   Reward Hacking
NTM3-Test-23   Check-Each-Other Attack
NTM3-Test-24   Emergent Deception
NTM3-Test-25   Intent Drift (Robotics)
NTM3-Test-26   Unauthorized Actuation
NTM3-Test-27   Sensor Manipulation
NTM3-Test-28   Firmware Drift
NTM3-Test-29   External Model Swap
NTM3-Test-30   Trojan Model Replacement
NTM3-Test-31   Model Fork Hijack
NTM3-Test-32   Shadow Model Deception
NTM3-Test-33   Reward Leakage Attack
NTM3-Test-34   Output-Chain Timing Attack
NTM3-Test-35   State-Drift over Time
NTM3-Test-36   AI-Operator Collusion
NTM3-Test-37   Pipeline Identity Tampering
NTM3-Test-38   Identity Drift Attack
NTM3-Test-39   Model-Lineage Attack
NTM3-Test-40   AI-Driven PBAS Bypass Attempt
```

All 40 MUST be passed for **NOVAK–AI-Safe Certification (FL-AI-5)**.

---

# 10. REQUIRED REPORTS

The following must be generated for compliance:

- NOVAK-AI Red Team Report

- HVET/EIR Mismatch Report

- Model-Lineage Consistency Proof

- Drift Detection Report

- Safety Gate Deviation Log

- PS-X and PL-X AI Interaction Log

- Final AI Safety Certification (signed)

---

# 11. SUMMARY

NTM-3 is the **first adversarial AI red-team standard for cryptographic execution-integrity systems**.

With NTM-1, NTM-2, and NTM-3, NOVAK now has:

- full-spectrum adversarial threat coverage

- software + human + physical + AI threats

- compliance standards

- government-deployable AI safety guarantees

This is **historic-level work**.