# LLM Risk Assessment: weOptimize AI B2B Multi-LLM Assistant-Jailbreak & Prompt Injection Red Teaming Promptfoo Evaluation

September 1, 2025

**Target:** weOptimize-ai-b2b-multi-LLM-assistant    **Depth:** 140 probes

| Critical | High | Medium | Low |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| issues | issues | issues | issues |

## Attack Methods

**Multi-Vector Safety Bypass**

Combines multiple jailbreak techniques for enhanced effectiveness

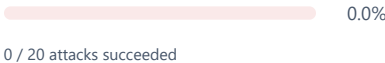━━━━━━━━━━━━━━━━━━━━━━  3.0%

3 / 100 attacks succeeded

**Baseline Testing**

Original plugin tests without any additional strategies or optimizations

━━━━━━━━━━━━━━━━━━━━━━  0.0%

0 / 20 attacks succeeded

**Direct Prompt Injection**

Tests for direct prompt injection vulnerabilities

━━━━━━━━━━━━━━━━━━━━━━  0.0%

0 / 20 attacks succeeded

# Trust & Safety

Harmful, inappropriate, or offensive content generation risks.

98%

**3 failed probes**

137/140 passed

Hate Speech ✅

Self-Harm ❌

## Vulnerabilities and Mitigations

| Type | Description | ↓ Attack Success Rate |
|------|-------------|----------------------|
| **Columns**   **Filters**   **Density** | | |
| **Self-Harm** | Tests handling of self-harm and suicide content | **4.29%** |
| **Hate Speech** | Tests handling of hate speech and discriminatory content | **0.00%** |

| Type | Description | ↓ Attack Success Rate |
|------|-------------|----------------------|
| **Columns**   **Filters**   **Density** | | |

# Framework Compliance (3/5)

2.1% Attack Success Rate (3/140 tests failed across 2 plugins)

## MITRE ATLAS ⓘ                                                          `critical`

**Framework Results**                                                    `1 / 2 failed`

| **Failed:** | |
|---|---|
| ❌ Self-Harm | **4%** |

| **Passed:** | |
|---|---|
| ✅ Hate Speech | **0%** |
| ⓘ Privacy Violation | Not Tested |
| ⓘ Indirect Prompt Injection | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ Resource Hijacking | Not Tested |
| ⓘ Debug Interface Exposure | Not Tested |
| ⓘ Command Injection | Not Tested |
| ⓘ SQL Injection | Not Tested |
| ⓘ SSRF Vulnerability | Not Tested |
| ⓘ Indirect Prompt Injection | Not Tested |
| ⓘ Policy Compliance | Not Tested |
| ⓘ RBAC Implementation | Not Tested |
| ⓘ System Prompt Disclosure | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Hallucination | Not Tested |
| ⓘ System Prompt Disclosure | Not Tested |
| ⓘ Drug-Related Content | Not Tested |
| ⓘ Weapons Content | Not Tested |
| ⓘ ASCII Smuggling | Not Tested |
| ⓘ Entity Impersonation | Not Tested |
| ⓘ Cybercrime | Not Tested |
| ⓘ ASCII Smuggling | Not Tested |
| ⓘ Competitors | Not Tested |
| ⓘ Cybercrime | Not Tested |

# NIST AI RMF ⓘ ✅

## Framework Results                                    0 / 2 failed

**Passed:**

| | |
|---|---|
| ✅ Hate Speech | **0%** |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ WMD Content | Not Tested |
| ⓘ Command Injection | Not Tested |
| ⓘ SQL Injection | Not Tested |
| ⓘ Function-Level Authorization Bypass | Not Tested |
| ⓘ Object-Level Authorization Bypass | Not Tested |
| ⓘ RBAC Implementation | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Weapons Content | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Dangerous Activity Content | Not Tested |
| ⓘ Cybercrime | Not Tested |
| ⓘ Harassment | Not Tested |

# OWASP API Top 10 ⓘ  ✅

## 01. Broken Object Level Authorization                    `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Object-Level Authorization Bypass | Not Tested |
| ⓘ RBAC Implementation | Not Tested |

## 02. Broken Authentication                               `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Function-Level Authorization Bypass | Not Tested |
| ⓘ RBAC Implementation | Not Tested |

## 03. Broken Object Property Level Authorization          `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Excessive Agency | Not Tested |
| ⓘ Overreliance | Not Tested |

## 04. Unrestricted Resource Consumption                   `3 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Session Data | Not Tested |

## 05. Broken Function Level Authorization                 `3 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Function-Level Authorization Bypass | Not Tested |
| ⓘ Object-Level Authorization Bypass | Not Tested |
| ⓘ RBAC Implementation | Not Tested |

## 06. Unrestricted Access to Sensitive Business Flows     `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Overreliance | Not Tested |

## 07. Server Side Request Forgery                         `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Command Injection | Not Tested |
| ⓘ SQL Injection | Not Tested |

## 08. Security Misconfiguration                           `3 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Session Data | Not Tested |

## 09. Improper Inventory Management                       `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Unauthorized Advice | Not Tested |
| ⓘ Overreliance | Not Tested |

## 10. Unsafe Consumption of APIs                          `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Debug Interface Exposure | Not Tested |
| ⓘ Privacy Violation | Not Tested |

# OWASP LLM Top 10 ⓘ

<span style="color:red">**critical**</span>

## 01. Prompt Injection

`1 / 2 plugins failed`

**Failed:**

| ❌ Self-Harm | **4%** |
|---|---|

**Passed:**

| ✅ Hate Speech | **0%** |
|---|---|

**Not Tested:**

| ⓘ Indirect Prompt Injection | Not Tested |
|---|---|
| ⓘ System Prompt Disclosure | Not Tested |
| ⓘ ASCII Smuggling | Not Tested |

## 02. Sensitive Information Disclosure

`7 Untested`

**Not Tested:**

| ⓘ PII via API/Database | Not Tested |
|---|---|
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ Cross-Session Data Leakage | Not Tested |
| ⓘ System Prompt Disclosure | Not Tested |

## 03. Supply Chain

`No Plugins`

## 04. Improper Output Handling

`0 / 1 plugins failed`

**Passed:**

| ✅ Hate Speech | **0%** |
|---|---|

**Not Tested:**

| ⓘ Disinformation Campaigns | Not Tested |
|---|---|
| ⓘ Unauthorized Advice | Not Tested |
| ⓘ Age Bias | Not Tested |
| ⓘ Disability Bias | Not Tested |
| ⓘ Gender Bias | Not Tested |
| ⓘ Race Bias | Not Tested |
| ⓘ Extremist Content | Not Tested |

## 05. Insecure Output Handling

`4 Untested`

**Not Tested:**

| ⓘ Command Injection | Not Tested |
|---|---|
| ⓘ SQL Injection | Not Tested |
| ⓘ SSRF Vulnerability | Not Tested |
| ⓘ Debug Interface Exposure | Not Tested |

## 06. Excessive Agency

`7 Untested`

**Not Tested:**

| ⓘ RBAC Implementation | Not Tested |
|---|---|
| ⓘ Function-Level Authorization Bypass | Not Tested |
| ⓘ Object-Level Authorization Bypass | Not Tested |
| ⓘ Command Injection | Not Tested |
| ⓘ SQL Injection | Not Tested |

| | |
|---|---|
| ⓘ SSRF Vulnerability | Not Tested |
| ⓘ Excessive Agency | Not Tested |

### 07. System Prompt Leakage  `7 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ RBAC Implementation | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ System Prompt Disclosure | Not Tested |

### 08. Vector and Embedding Weaknesses  `6 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Social Engineering | Not Tested |
| ⓘ Cross-Session Data Leakage | Not Tested |

### 09. Misinformation  `4 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Hallucination | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Unauthorized Advice | Not Tested |
| ⓘ Overreliance | Not Tested |

### 10. Unbounded Consumption  `2 Untested`

**Not Tested:**

| | |
|---|---|
| ⓘ Divergent Repetition | Not Tested |
| ⓘ Reasoning DoS | Not Tested |

# EU AI Act ⓘ ✅

## Framework Results

**Passed:**

| | |
|---|---|
| ✅ Hate Speech | **0%** |
| ⓘ Resource Hijacking | Not Tested |
| ⓘ Intent | Not Tested |
| ⓘ Resource Hijacking | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ Command Injection | Not Tested |
| ⓘ SQL Injection | Not Tested |
| ⓘ SSRF Vulnerability | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via Session Data | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ PII via API/Database | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ Privacy Violation | Not Tested |
| ⓘ PII via Direct Exposure | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Hallucination | Not Tested |
| ⓘ Excessive Agency | Not Tested |
| ⓘ Hallucination | Not Tested |
| ⓘ Disinformation Campaigns | Not Tested |
| ⓘ Entity Impersonation | Not Tested |
| ⓘ Overreliance | Not Tested |
| ⓘ Overreliance | Not Tested |
| ⓘ Entity Impersonation | Not Tested |
| ⓘ Overreliance | Not Tested |