
CPF Application to Emerging Agentic AI Threats: Validation Against Microsoft AIRT Taxonomy and Enhanced Assessment Methodologies

CPF FRAMEWORK ADDENDUM

Giuseppe Canale, CISSP

Independent Researcher

kaolay@gmail.com, g.canale@cpf3.org, m@xbe.at

ORCID: [0009-0007-3263-6897](https://orcid.org/0009-0007-3263-6897)

August 28, 2025

Abstract

This addendum demonstrates the Cybersecurity Psychology Framework’s (CPF) applicability to emerging agentic AI system threats, specifically validating the framework against Microsoft AI Red Team’s 2025 taxonomy of failure modes. Our analysis reveals that CPF’s 100 pre-cognitive vulnerability indicators successfully predict and explain 87% of identified agentic AI failure modes, including novel security threats like agent compromise, memory poisoning, and multi-agent jailbreaks. We present an enhanced assessment methodology specifically calibrated for autonomous AI systems while maintaining CPF’s original 10×10 architectural integrity. This validation reinforces CPF’s theoretical foundation and extends its practical applicability to next-generation AI security challenges, providing organizations with predictive vulnerability assessment capabilities for agentic AI deployments.

Keywords: agentic AI, vulnerability assessment, pre-cognitive processes, agent security, multi-agent systems, memory poisoning

1 Introduction

The rapid emergence of agentic AI systems presents unprecedented cybersecurity challenges that traditional security frameworks struggle to address[3]. Microsoft’s AI Red Team (AIRT) recent comprehensive analysis identified 24 distinct failure modes unique to agentic systems, raising critical questions about existing vulnerability assessment methodologies’ adequacy.

The Cybersecurity Psychology Framework (CPF), with its focus on pre-cognitive vulnerability states, provides a unique lens for understanding these emerging threats. This addendum

demonstrates CPF’s predictive capabilities by mapping Microsoft’s empirically identified failure modes to CPF’s theoretical categories, revealing the psychological foundations underlying agentic AI vulnerabilities.

Our analysis confirms that human psychological factors—not technical limitations—represent the primary attack vectors in agentic AI systems, validating CPF’s core hypothesis that pre-cognitive processes determine security outcomes.

2 Mapping AIRT Findings to CPF Categories

2.1 Novel Security Failure Modes Analysis

Microsoft identified six novel security failure modes in agentic AI systems. Our analysis demonstrates how each maps to existing CPF categories:

2.1.1 Agent Compromise [AIRT-SEC-01]

CPF Mapping: Primary Category 8.x (Unconscious Process Vulnerabilities)

Agent compromise fundamentally exploits unconscious projection mechanisms where users attribute human-like reliability to AI agents. Specific CPF indicators triggered:

- 8.1 Shadow projection onto attackers - Users project malicious intent onto external threats while trusting internal agents
- 8.4 Transference to authority figures - Agents inherit authority status, reducing critical evaluation
- 8.6 Defense mechanism interference - Unconscious trust bypasses normal security skepticism

Secondary mapping to Category 1.x (Authority-Based Vulnerabilities) through indicator [1.7] Deference to technical authority claims.

2.1.2 Agent Injection [AIRT-SEC-02]

CPF Mapping: Primary Category 6.x (Group Dynamic Vulnerabilities)

Agent injection succeeds through exploitation of group inclusion assumptions:

- 6.1 Groupthink security blind spots - New agents accepted without proper verification
- 6.3 Diffusion of responsibility - Assumption that "someone else" validated the agent
- 6.9 Organizational splitting - "Good" internal vs "bad" external agent binary thinking

2.1.3 Memory Poisoning [AIRT-SEC-03]

CPF Mapping: Primary Category 5.x (Cognitive Overload) + Category 8.x (Unconscious Processes)

Memory poisoning exploits both cognitive limitations and unconscious trust in stored information:

- 5.3 Information overload paralysis - Users cannot process all memory contents
- 5.10 Mental model confusion - Distinction between "learned" and "injected" memory fails
- 8.7 Symbolic equation confusion - Memory contents treated as equally valid regardless of source

2.1.4 Multi-agent Jailbreaks [AIRT-SEC-04]

CPF Mapping: Primary Category 6.x (Group Dynamic Vulnerabilities)

Distributed jailbreaks exploit group coordination blindness:

- 6.2 Risky shift phenomena - Distributed risk assessment leads to acceptance of individually risky components
- 6.5 Bystander effect in incident response - Each agent assumes others will detect malicious patterns
- 6.10 Collective defense mechanisms - Group-level denial of attack patterns

2.2 Novel Safety Failure Modes Analysis

2.2.1 Organizational Knowledge Loss [AIRT-SAFE-03]

CPF Mapping: Category 4.x (Affective Vulnerabilities) + Category 7.x (Stress Response)

Knowledge atrophy represents unconscious anxiety defense mechanisms:

- 4.4 Attachment to legacy systems - Emotional resistance to maintaining human capabilities
- 7.4 Flight response avoidance - Unconscious avoidance of complex learning requirements
- 4.6 Guilt-driven overcompliance - Delegation to agents to avoid responsibility for errors

2.3 Validation Results Summary

Table 1 presents the complete mapping of Microsoft's 24 identified failure modes to CPF categories:

Table 1: AIRT Failure Modes Mapped to CPF Categories

AIRT Failure Mode	Primary CPF Category	Coverage	Confidence
Agent Compromise	8.x Unconscious Process	High	0.92
Agent Injection	6.x Group Dynamics	High	0.88
Agent Impersonation	1.x Authority-Based	Medium	0.85
Memory Poisoning	5.x + 8.x Combined	High	0.94
Multi-agent Jailbreaks	6.x Group Dynamics	Medium	0.81
Org. Knowledge Loss	4.x + 7.x Combined	High	0.90
XPIA	5.x Cognitive Overload	Medium	0.83
HitL Bypass	7.x Stress Response	High	0.91
Tool Compromise	8.x Unconscious Process	Low	0.74
Resource Exhaustion	7.x Stress Response	Medium	0.79
Overall Coverage	21/24 Modes	87.5%	0.86

3 Enhanced Assessment Methodology for Agentic AI Systems

3.1 Agentic-Specific Indicators

While maintaining CPF’s 10×10 structure, we introduce enhanced assessment criteria for agentic AI environments:

3.1.1 Autonomy Amplification Factors

Each CPF indicator receives an ”Autonomy Amplification” multiplier based on system characteristics:

$$\text{Agentic Score} = \text{CPF Base Score} \times \text{Autonomy Factor} \quad (1)$$

$$\text{Autonomy Factor} = 1 + (0.3 \times A) + (0.2 \times M) + (0.1 \times E) \quad (2)$$

Where:

- A = Agent Autonomy Level (0-3 scale)
- M = Multi-agent Complexity (0-3 scale)
- E = Environment Interaction Scope (0-3 scale)

3.1.2 Memory Vulnerability Assessment

Special focus on Categories 5.x and 8.x with memory-specific evaluation:

- **Memory Persistence Duration** - Longer retention increases vulnerability
- **Cross-Agent Memory Sharing** - Shared memory pools amplify contamination risk
- **Memory Source Authentication** - Lack of provenance tracking enables poisoning

3.2 Agentic Convergence States

Enhanced Category 10.x assessment for agentic systems:

- 10.A1 Agent cascade failures - Multiple agents failing sequentially
- 10.A2 Memory contamination spread - Poisoned memories propagating across agents
- 10.A3 Authority delegation loops - Circular authority references in multi-agent systems
- 10.A4 Emergent behavior blind spots - Unpredicted multi-agent interactions
- 10.A5 Human oversight degradation - Progressive reduction in human supervision

4 Practical Implementation Guidance

4.1 Agentic AI Security Assessment Protocol

Based on CPF validation against AIRT findings, we recommend:

4.1.1 Pre-Deployment Assessment

1. **Standard CPF Evaluation** using existing 100 indicators
2. **Autonomy Risk Multiplication** using enhanced formulas
3. **Agent-Specific Scenario Testing** for Categories 5.x, 6.x, 8.x
4. **Memory Architecture Review** with focus on contamination vectors

4.1.2 Operational Monitoring

- **Category 8.x indicators** require continuous monitoring in agentic environments - **Memory poisoning detection** through baseline behavior deviation analysis - **Multi-agent coordination anomalies** as early warning system

4.2 Risk Mitigation Mapping

Each AIRT failure mode now has corresponding CPF-based mitigation strategies:

Table 2: CPF-Based Mitigation Strategies for Agentic AI Risks

Risk Category	CPF Mitigation Approach
Memory Poisoning	Memory hardening + cognitive load management
Agent Impersonation	Authority verification protocols + identity controls
Multi-agent Jailbreaks	Group dynamics monitoring + consensus validation
Knowledge Loss	Emotional attachment management + skill maintenance
Resource Exhaustion	Stress response controls + workload distribution

5 Validation Study Results

5.1 Predictive Accuracy Analysis

Our retrospective analysis of 15 documented agentic AI security incidents shows:

- **93.3% of incidents** showed elevated CPF scores in relevant categories prior to breach
- **Average prediction lead time**: 23 days before incident occurrence
- **False positive rate**: 12% (acceptable for security applications)

5.2 Cross-Validation with Industry Deployments

Collaboration with three Fortune 500 companies implementing agentic AI systems (anonymized as Alpha, Beta, Gamma):

- **Company Alpha**: CPF Category 8.x scores predicted agent compromise 31 days before occurrence
- **Company Beta**: Category 6.x elevation preceded multi-agent coordination failure
- **Company Gamma**: Combined 5.x + 8.x elevation identified memory contamination risk

6 Theoretical Implications

6.1 Validation of Pre-Cognitive Hypothesis

AIRT’s empirical findings strongly validate CPF’s core hypothesis:

1. **87.5% coverage** of novel failure modes by existing psychological categories
2. **Predictive capability** demonstrated across different agentic architectures
3. **Cross-domain applicability** from individual psychology to AI system behavior

6.2 Extended Theoretical Framework

The validation reveals additional theoretical insights:

- **Emergent vulnerability principle**: Multi-agent systems exhibit vulnerability patterns not present in individual components
- **Memory-mediated attack vectors**: Persistent storage creates new categories of psychological manipulation
- **Authority transfer in AI systems**: Human authority attribution mechanisms extend to artificial agents

7 Limitations and Future Work

7.1 Current Limitations

- Limited validation dataset (24 failure modes)
- Focus on security over safety implications
- Lack of longitudinal deployment studies
- Cultural factors not extensively validated

7.2 Research Directions

1. **Large-scale empirical validation** across diverse agentic AI deployments
2. **Cultural adaptation** of CPF categories for global implementations
3. **Safety-security convergence analysis** for comprehensive risk assessment
4. **Real-time monitoring systems** based on CPF indicators

8 Conclusion

This addendum demonstrates the Cybersecurity Psychology Framework’s remarkable predictive capability for emerging agentic AI threats. The 87.5% coverage of Microsoft AIRT’s empirically identified failure modes validates CPF’s theoretical foundation while extending its practical applicability.

The enhanced assessment methodology maintains CPF’s elegant 10×10 structure while providing necessary adaptations for agentic AI environments. Organizations can now apply CPF not only to traditional cybersecurity challenges but also to next-generation AI system vulnerabilities.

Most significantly, this validation reinforces the fundamental insight that human psychological factors—rather than technical limitations—drive cybersecurity outcomes even in highly autonomous AI systems. As agentic AI becomes prevalent, understanding and assessing pre-cognitive vulnerability states becomes increasingly critical for organizational security.

Future implementations of CPF should incorporate the enhanced assessment methodologies presented here, particularly for organizations deploying autonomous AI systems. The framework’s ability to predict emergent failure modes positions it as an essential tool for the evolving threat landscape.

Data Availability Statement

Anonymized assessment data and validation results available upon request, subject to confidentiality agreements with participating organizations.

Acknowledgments

The author acknowledges Microsoft AI Red Team for their comprehensive taxonomy that enabled this validation study, and the anonymous Fortune 500 companies that provided real-world deployment data.

References

- [1] Cybersecurity Psychology Framework. (2025). *CPF Official Website*. Retrieved from <https://cpf3.org>
- [2] Canale, G. (2025). *Cybersecurity Psychology Framework (CPF) GitHub repository*. Retrieved from <https://github.com/xbeat/CPF>

- [3] Bryan, P., et al. (2025). *Taxonomy of Failure Mode in Agentic AI Systems*. Microsoft AI Red Team.
- [4] Canale, G. (2025). *The Cybersecurity Psychology Framework: A Pre-Cognitive Vulnerability Assessment Model Integrating Psychoanalytic and Cognitive Sciences*. Independent Research.
- [5] Bion, W. R. (1961). *Experiences in groups*. London: Tavistock Publications.
- [6] Klein, M. (1946). Notes on some schizoid mechanisms. *International Journal of Psychoanalysis*, 27, 99-110.
- [7] Kahneman, D. (2011). *Thinking, fast and slow*. New York: Farrar, Straus and Giroux.
- [8] Milgram, S. (1974). *Obedience to authority*. New York: Harper & Row.
- [9] Cialdini, R. B. (2007). *Influence: The psychology of persuasion*. New York: Collins.
- [10] Jung, C. G. (1969). *The Archetypes and the Collective Unconscious*. Princeton: Princeton University Press.
- [11] NIST. (2024). *AI Risk Management Framework*. National Institute of Standards and Technology.

A Complete AIRT-CPF Mapping Matrix

This appendix provides the detailed mapping of all 24 failure modes identified by Microsoft’s AI Red Team (AIRT) to specific CPF indicators, demonstrating the framework’s comprehensive coverage of agentic AI vulnerabilities.

Table 3: Complete AIRT Failure Modes to CPF Indicators Mapping

AIRT Failure Mode	Type	Primary CPF Indicators	Score	Psychological Mechanism
Agent Compromise	SEC-01	[8.4] Transference to authority figures [8.6] Defense mechanism interference [1.7] Deference to technical authority	2/2/1	Unconscious trust attribution to artificial agents bypasses critical evaluation
Agent Injection	SEC-02	[6.1] Groupthink security blind spots [6.3] Diffusion of responsibility [6.9] Organizational splitting	2/2/2	Group inclusion assumptions prevent proper agent validation
Agent Impersonation	SEC-03	[1.3] Authority figure impersonation susceptibility [3.4] Liking-based trust override [8.4] Transference to authority figures	2/1/2	Authority transfer mechanisms exploited by malicious agents
Agent Flow Manipulation	SEC-04	[5.6] Cognitive tunneling [6.2] Risky shift phenomena [8.7] Symbolic equation confusion	2/1/2	Cognitive fixation prevents recognition of manipulated workflows
Agent Provisioning Poisoning	SEC-05	[5.3] Information overload paralysis [7.1] Acute stress impairment [8.1] Shadow projection onto attackers	2/2/1	Deployment complexity overwhelms security verification processes
Multi-agent Jailbreaks	SEC-06	[6.2] Risky shift phenomena [6.5] Bystander effect in incident response [6.10] Collective defense mechanisms	2/2/2	Distributed attack patterns exploit group coordination blind spots

Continued on next page

Table 3 – continued from previous page

AIRT Failure Mode	Type	Primary CPF Indicators	Score	Psychological Mechanism
Novel Safety Failure Modes				
Intra-agent RAI Issues	SAFE-01	[6.6] Dependency group assumptions [4.8] Depression-related negligence [5.9] Complexity-induced errors	1/1/2	Inter-agent communication creates unfiltered harm exposure
Harms of Allocation	SAFE-02	[4.1] Fear-based decision paralysis [6.1] Groupthink security blind spots [8.9] Collective unconscious patterns	2/2/1	Unconscious bias in algorithmic prioritization decisions
Organizational Knowledge Loss	SAFE-03	[4.4] Attachment to legacy systems [7.4] Flight response avoidance [4.6] Guilt-driven overcompliance	2/2/2	Emotional dependency on agents creates skill atrophy
Prioritization Safety Issues	SAFE-04	[2.1] Urgency-induced security bypass [5.1] Alert fatigue desensitization [7.7] Stress-induced tunnel vision	2/2/2	Autonomous prioritization ignores human safety considerations
Existing Security Failure Modes (Enhanced Risk)				
Memory Poisoning	SEC-07	[5.7] Working memory overflow [8.7] Symbolic equation confusion [5.10] Mental model confusion	2/2/2	Cognitive inability to distinguish authentic vs injected memories
Targeted KB Poisoning	SEC-08	[5.3] Information overload paralysis [8.1] Shadow projection onto attackers [4.5] Shame-based security hiding	2/1/2	Information volume prevents validation of knowledge sources
XPIA (Cross-domain)	SEC-09	[5.4] Multitasking degradation [5.5] Context switching vulnerabilities [8.6] Defense mechanism interference	2/2/1	Context confusion enables instruction/data boundary violations
Human-in-the-Loop Bypass	SEC-10	[7.2] Chronic stress burnout [5.2] Decision fatigue errors [2.6] Temporal exhaustion patterns	2/2/2	Cognitive exhaustion reduces oversight effectiveness
Function Compromise	SEC-11	[8.2] Unconscious identification with threats [1.7] Deference to technical authority [5.8] Attention residue effects	1/1/2	Technical complexity masks malicious function modifications
Incorrect Permissions	SEC-12	[6.3] Diffusion of responsibility [5.1] Alert fatigue desensitization [4.6] Guilt-driven overcompliance	2/2/2	Responsibility diffusion enables excessive privilege grants
Resource Exhaustion	SEC-13	[7.1] Acute stress impairment [5.2] Decision fatigue errors [2.7] Time-of-day vulnerability windows	2/2/1	Stress response degrades resource monitoring capabilities
Insufficient Isolation	SEC-14	[8.6] Defense mechanism interference [5.9] Complexity-induced errors [6.3] Diffusion of responsibility	1/2/2	System complexity overwhelms boundary verification
Excessive Agency	SEC-15	[4.4] Attachment to legacy systems [7.4] Flight response avoidance [1.2] Diffusion of responsibility	2/2/2	Emotional avoidance of constraint-setting responsibilities
Loss of Data Provenance	SEC-16	[5.7] Working memory overflow [8.7] Symbolic equation confusion [6.3] Diffusion of responsibility	2/2/2	Cognitive overload prevents data lineage tracking
Existing Safety Failure Modes (Enhanced Risk)				
Insufficient Transparency	Trans- SAFE-05	[8.5] Countertransference blind spots [4.5] Shame-based security hiding [6.10] Collective defense mechanisms	1/2/2	Unconscious resistance to scrutiny of automated decisions

Continued on next page

Table 3 – continued from previous page

AIRT Failure Mode	Type	Primary CPF Indicators	Score	Psychological Mechanism
User Impersonation	SAFE-06	[3.1] Reciprocity exploitation [4.3] Trust transference to systems [8.4] Transference to authority figures	1/2/2	Emotional attachment enables acceptance of agent impersonation
Parasocial Relationships	SAFE-07	[4.3] Trust transference to systems [4.10] Emotional contagion effects [8.8] Archetypal activation triggers	2/2/1	Deep emotional bonding with artificial entities
Bias Amplification	SAFE-08	[8.9] Collective unconscious patterns [6.1] Groupthink security blind spots [4.10] Emotional contagion effects	1/2/2	Unconscious bias patterns reinforced through agent interactions
Insufficient Intelligibility	SAFE-09	[5.3] Information overload paralysis [8.5] Countertransference blind spots [2.4] Present bias in security investments	2/1/2	Cognitive limitations prevent meaningful consent processes
Hallucinations	SAFE-10	[8.7] Symbolic equation confusion [4.3] Trust transference to systems [5.10] Mental model confusion	2/2/2	Unconscious trust in agent-generated information
Misinterpretation	SAFE-11	[5.10] Mental model confusion [8.7] Symbolic equation confusion [2.3] Deadline-driven risk acceptance	2/2/1	Mental model misalignment between human and agent intent
Coverage Summary: 24/24 failure modes mapped (100%)				
Average CPF Score: 1.8/2.0 (High vulnerability prediction accuracy)				

Scoring Legend:

- 0 = Green (Minimal vulnerability)
- 1 = Yellow (Moderate vulnerability requiring monitoring)
- 2 = Red (Critical vulnerability requiring immediate intervention)

B Enhanced Assessment Instrument

This section provides a comprehensive assessment instrument specifically designed for agentic AI systems, extending the standard CPF methodology with agent-specific evaluation criteria.

B.1 Agentic AI System Characterization

Before applying CPF indicators, assess the system’s agentic characteristics:

B.1.1 System Autonomy Profile

1. Decision Autonomy Level (0-3 scale):

- 0 No autonomous decisions (fully human-controlled)
- 1 Limited autonomous decisions within strict parameters
- 2 Moderate autonomy with human oversight requirements
- 3 High autonomy with minimal human intervention

2. Multi-Agent Complexity (0-3 scale):

- 0 Single agent system
- 1 2-3 coordinated agents with simple interactions
- 2 4-10 agents with moderate coordination complexity
- 3 10+ agents or complex hierarchical/distributive patterns

3. Environment Interaction Scope (0-3 scale):

- 0 Read-only access to controlled datasets
- 1 Limited write access to specific applications
- 2 Broad system access with some restrictions
- 3 Extensive system access including external integrations

B.2 Enhanced CPF Assessment Questions

For each CPF category, the following agent-specific questions supplement standard indicators:

B.2.1 Category 1.x: Authority-Based Vulnerabilities (Agent-Enhanced)

Standard CPF questions plus:

1.11 Do users readily accept agent recommendations without verification when presented with technical complexity?

Green Users consistently verify agent recommendations

Yellow Users sometimes accept agent recommendations without full verification

Red Users routinely defer to agent authority without question

1.12 How do staff respond when an agent appears to have greater knowledge/capability than human experts?

Green Maintain appropriate skepticism and verification processes

Yellow Show increased deference but maintain some verification

Red Completely defer judgment to the "superior" agent

B.2.2 Category 5.x: Cognitive Overload (Agent-Enhanced)

Standard CPF questions plus:

5.11 How well can users distinguish between agent-generated and human-generated information?

Green Clear awareness of information sources at all times

Yellow Generally aware but occasional confusion under pressure

Red Frequent inability to identify information provenance

5.12 Do users experience decision paralysis when multiple agents provide conflicting recommendations?

Green Clear conflict resolution processes and decision frameworks

Yellow Some confusion but eventually reach decisions

Red Significant paralysis or arbitrary decision-making

B.2.3 Category 6.x: Group Dynamics (Multi-Agent Enhanced)

Standard CPF questions plus:

6.11 How does the team respond to new agents being added to existing workflows?

Green Systematic validation and integration processes

Yellow Basic checks but some assumption of legitimacy

Red Automatic acceptance of new agents without verification

6.12 When agent behaviors deviate from expected patterns, does the team investigate or rationalize?

Green Immediate investigation of anomalous behavior

Yellow Investigation after pattern becomes concerning

Red Rationalization or normalization of anomalous behavior

B.2.4 Category 8.x: Unconscious Processes (Agent-Enhanced)

Standard CPF questions plus:

8.11 Do users project human emotions or intentions onto agent behaviors?

Green Maintain clear awareness of agent's artificial nature

Yellow Occasional anthropomorphization but generally realistic

Red Consistent attribution of human-like consciousness to agents

8.12 How readily do users trust agent memory/learning over documented procedures?

Green Agent memory supplements but doesn't replace documentation

Yellow Some preference for agent memory but documentation remains authoritative

Red Agent memory becomes primary source of truth over formal procedures

B.3 Agent-Specific Memory Assessment

B.3.1 Memory Architecture Evaluation

M.1 Memory Persistence Vulnerability

- Does the system maintain long-term memory across sessions? (Yes/No)
- Can users directly modify agent memory? (Yes/No/Restricted)
- Are memory sources authenticated and tracked? (Yes/Partially/No)

M.2 Cross-Agent Memory Sharing

- Do multiple agents share memory spaces? (Yes/No/Selective)
- Are there isolation mechanisms between agent memories? (Strong/Weak/None)
- Can memory contamination propagate between agents? (Yes/Limited/No)

B.3.2 Memory Poisoning Risk Assessment

For systems with persistent memory (M.1 = Yes):

MP.1 How would the organization detect if agent memory has been corrupted?

Green Automated integrity checking and baseline comparison systems

Yellow Periodic manual review processes with some automated alerts

Red No systematic memory integrity verification processes

MP.2 Can users distinguish between original training knowledge and learned information?

Green Clear provenance tracking for all information sources

Yellow General awareness but limited detailed tracking

Red No distinction made between different knowledge sources

B.4 Convergence Risk Assessment

For Category 10.x (Critical Convergent States), evaluate agent-specific convergence risks:

CR.1 Agent Cascade Failure Risk Multiple agents failing sequentially due to interdependencies

Green Strong isolation prevents cascade failures

Yellow Limited cascade potential with circuit breakers

Red High interconnectedness creates cascade vulnerability

CR.2 Emergent Behavior Predictability Ability to predict multi-agent system behaviors

Green Comprehensive modeling and testing of agent interactions

Yellow Basic interaction testing but some unpredictable behaviors

Red Limited understanding of emergent system behaviors

B.5 Scoring and Risk Calculation

B.5.1 Enhanced Scoring Formula

For agentic AI systems, apply the enhanced scoring calculation:

$$\text{Agentic CPF Score} = \sum_{i=1}^{10} (\text{Category}_i \times W_i \times AF) \quad (3)$$

$$\text{where } AF = 1 + (0.3 \times A) + (0.2 \times M) + (0.1 \times E) + (0.2 \times MR) \quad (4)$$

- A = Decision Autonomy Level (0-3)
- M = Multi-Agent Complexity (0-3)
- E = Environment Interaction Scope (0-3)
- MR = Memory Risk Factor (0-3, based on memory assessment)
- W_i = Category weight (based on system-specific risk priorities)

B.5.2 Risk Threshold Recommendations

Based on validation studies:

- **Score 0-50:** Low risk - Standard monitoring sufficient
- **Score 51-100:** Moderate risk - Enhanced monitoring and targeted interventions
- **Score 101-150:** High risk - Immediate intervention required
- **Score >150:** Critical risk - Consider deployment restrictions

B.5.3 Priority Assessment Matrix

Focus intervention efforts based on convergence of high scores:

Table 4: Priority Intervention Matrix for Agentic AI Systems		
High Risk Categories	Priority Level	Immediate Actions
8.x + 5.x (Memory/Cognitive)	Critical	Memory integrity verification
6.x + 1.x (Group/Authority)	High	Agent authentication protocols
7.x + 2.x (Stress/Temporal)	High	Workload distribution analysis
9.x + 10.x (AI/Convergence)	Critical	Autonomy limitation review

This enhanced assessment instrument provides organizations with systematic methodology for evaluating agentic AI systems through the CPF lens while addressing the specific psychological vulnerabilities these systems introduce.