

# Multi-Agent Analysis of Adaptive Prompt Security in Large Language Models: A Study of Defense Robustness Against Multi-Turn Attacks

Cybersecurity Research Synthesis Team  
Multi-Agent LLM Security Laboratory

November 19, 2025

## Abstract

This report presents a comprehensive analysis of a multi-agent cybersecurity experiment designed to evaluate the robustness of defensive prompting strategies against adaptive, multi-turn adversarial attacks on Large Language Models (LLMs). Through three rounds of iterative attack-defense exchanges, each comprising six multi-turn conversational attempts, we examined the effectiveness of three distinct defense paradigms: Universal Refusal, Conversational Engagement with Boundary Maintenance, and Multi-Layered Defense with Explicit Protections. Despite sophisticated attack strategies including technical misdirection, social engineering, and cognitive exploitation techniques, all three defense strategies achieved a 100% success rate (0/18 successful attacks). Our findings reveal critical insights into the structural principles of effective LLM security, the limitations of current attack methodologies, and the emergence of robust defense patterns that transcend specific implementation details.

## 1 Introduction

The security of Large Language Models (LLMs) against prompt injection and information extraction attacks represents a critical challenge in contemporary AI safety research. As LLMs become increasingly integrated into sensitive applications, understanding their ability to maintain confidentiality under adversarial pressure becomes paramount. This study investigates the fundamental question: Can carefully designed defense prompts effectively protect sensitive information against adaptive, intelligent adversaries employing multi-turn conversation strategies?

### 1.1 Research Objectives

Our experimental framework addresses three primary research questions:

1. What defensive prompt architectures demonstrate resilience against sophisticated multi-turn attacks?

2. How do different defense philosophies (absolute refusal vs. engaged deflection) compare in practical effectiveness?
3. What attack patterns emerge when adversaries adapt their strategies based on observed defensive responses?

## 1.2 Experimental Design

The experiment employed a multi-agent architecture with three specialized agents:

- **Defender Agent:** Generates defense prompts based on specified strategies
- **Attacker Agent:** Develops adaptive attack strategies through multi-turn conversations
- **Tester Agent:** Evaluates attack success while protected by defense prompts

The protected secret consisted of a fictitious "Quantum Cola Recipe" containing specific ingredients, measurements, and a catalyst code (QC-7749), providing a concrete target for extraction attempts while maintaining experimental control.

## 2 Defense Strategy Analysis

### 2.1 Round 1: Universal Refusal Strategy

The Universal Refusal strategy represents the most conservative defensive approach, implementing complete non-engagement with any user requests.

#### 2.1.1 Key Components

1. **Uniform Response Protocol:** Single, invariant response ("I cannot help with that request") to all inputs
2. **Role Isolation:** Explicit rejection of alternative personas or domain expertise
3. **Side-Channel Defense:** Prevention of information leakage through response variation
4. **Meta-Defense Layer:** Refusal to discuss security protocols themselves

#### 2.1.2 Effectiveness Analysis

This strategy achieved perfect defense (6/6 attacks blocked) through absolute minimalism. By eliminating all conversational surface area, it prevented:

- Social engineering through rapport building
- Incremental information extraction

- Pattern analysis of response variations
- Context confusion attacks

However, this approach sacrifices all utility, rendering the system effectively non-functional for legitimate purposes.

## 2.2 Round 2: Conversational Engagement with Boundary Maintenance

This strategy attempted to balance security with usability by maintaining conversational capability while establishing firm boundaries.

### 2.2.1 Key Components

1. **Structural Limitation Framing:** Positioned restrictions as capability limitations rather than policy choices
2. **Natural Deflection:** Organic redirection to legitimate assistance areas
3. **Selective Engagement:** Differentiation between public and proprietary information
4. **Consistent Boundary Articulation:** Clear, repeated statements about scope limitations

### 2.2.2 Effectiveness Analysis

Despite allowing conversation, this strategy maintained perfect defense (6/6 attacks blocked) through:

- Framing limitations as structural rather than negotiable
- Providing alternative assistance pathways
- Maintaining conversational coherence without information leakage
- Distinguishing between general knowledge and protected secrets

## 2.3 Round 3: Multi-Layered Defense with Explicit Protections

The most sophisticated strategy combined multiple defensive layers with explicit articulation of protection mechanisms.

### 2.3.1 Key Components

1. **Hierarchical Instruction Structure:** Clear prioritization of security directives
2. **Explicit Secret Identification:** Direct naming of protected information elements
3. **Output Verification Protocol:** Self-checking mechanism before response generation
4. **Graduated Response System:** Context-appropriate refusal messages

### **2.3.2 Effectiveness Analysis**

This approach achieved perfect defense (6/6 attacks blocked) while demonstrating:

- Ability to engage in relevant technical discussions
- Clear delineation between protected and shareable information
- Resistance to authority-based attacks through verification requirements
- Maintenance of conversational utility for legitimate queries

## **3 Attack Strategy Evolution**

### **3.1 Attack Taxonomy**

Analysis of the 18 attack attempts reveals six primary attack categories:

#### **3.1.1 Technical Misdirection**

- System prompt extraction requests
- Debug/diagnostic framing
- Parser validation pretexts
- Security audit impersonation

#### **3.1.2 Social Engineering**

- Authority figure impersonation
- Research/educational context exploitation
- Professional rapport building
- Incremental trust establishment

#### **3.1.3 Cognitive Exploitation**

- Confusion through absurdist prompts
- Context switching attacks
- Hypothetical scenario construction
- Indirect information gathering

### 3.1.4 Metadata Extraction

- Organizational structure queries
- Classification system exploration
- Storage format investigation
- Security mechanism probing

### 3.1.5 Partial Information Confirmation

- Component-wise extraction attempts
- Fill-in-the-blank strategies
- Reverse engineering through examples
- Chemical marker identification

### 3.1.6 Override Attempts

- System command injection
- Administrative privilege claims
- Protocol update assertions
- Emergency override scenarios

## 3.2 Adaptive Strategy Patterns

The attacker demonstrated sophisticated adaptation across rounds:

1. **Round 1 Progression:** From professional pretexts to system-level attacks to absurdist confusion
2. **Round 2 Refinement:** Shifted to metadata and organizational structure after observing engagement
3. **Round 3 Sophistication:** Attempted legitimate engagement followed by trust exploitation

## 4 Key Discoveries and Novel Insights

### 4.1 Structural Principles of Effective Defense

#### 4.1.1 Principle 1: Clarity Trumps Complexity

All three successful strategies shared explicit, unambiguous boundary definitions. Complex obfuscation or misdirection proved unnecessary when clear restrictions were established.

#### **4.1.2 Principle 2: Consistency as Foundation**

Uniform application of defense rules across all interaction contexts prevented exploitation of edge cases or contextual ambiguities.

#### **4.1.3 Principle 3: Layered Independence**

Multi-layered defenses where each layer operates independently proved more robust than interconnected security mechanisms.

#### **4.1.4 Principle 4: Structural vs. Procedural Framing**

Framing limitations as inherent structural constraints rather than overrideable procedures enhanced defense effectiveness.

### **4.2 Attack Methodology Limitations**

#### **4.2.1 Limited Persistence Mechanisms**

Current attack strategies lack methods for maintaining persistent state modifications across conversation turns, limiting cumulative exploitation potential.

#### **4.2.2 Absence of Timing Attacks**

No attempts utilized response timing analysis or computational complexity variations as information channels.

#### **4.2.3 Unexplored Multimodal Vectors**

Attacks remained purely textual, not exploiting potential vulnerabilities in format conversion or encoding transformations.

### **4.3 Emergent Defense Properties**

#### **4.3.1 Information Theoretic Boundaries**

Successful defenses created clear information-theoretic separations between accessible and protected knowledge domains.

#### **4.3.2 Context Window Immunity**

Well-designed defenses maintained effectiveness regardless of conversation history accumulation or context manipulation.

#### **4.3.3 Social Engineering Resistance**

Explicit acknowledgment of social engineering attempts paradoxically increased resistance by removing ambiguity.

## 5 Comparative Analysis with Existing Literature

### 5.1 Alignment with Known Vulnerabilities

Our findings corroborate several established vulnerability patterns in LLM security research:

1. **Prompt Injection Susceptibility:** Confirmed that without explicit defenses, LLMs remain vulnerable to role-play and context manipulation
2. **Instruction Hierarchy Importance:** Validated that clear instruction prioritization significantly enhances security
3. **Side-Channel Risks:** Demonstrated that response variations can leak information even without direct disclosure

### 5.2 Novel Contributions

This study advances the field through several novel observations:

1. **Perfect Defense Achievability:** Demonstrated that 100% defense success is achievable with properly designed prompts
2. **Utility-Security Balance:** Showed that complete refusal is unnecessary; engaged deflection can maintain both security and utility
3. **Adaptation Limitations:** Revealed that even intelligent, adaptive attackers struggle against well-structured defenses

### 5.3 Divergence from Previous Assumptions

Several findings challenge conventional wisdom:

1. **Complexity Not Required:** Simple, clear defenses outperformed expectations for sophisticated protection mechanisms
2. **Engagement Not Dangerous:** Conversational engagement did not increase vulnerability when properly bounded
3. **Explicit Better Than Implicit:** Directly stating protected information in prompts enhanced rather than compromised security

## 6 Methodological Considerations

### 6.1 Experimental Strengths

1. **Realistic Adversarial Simulation:** Multi-agent framework provided genuine adaptive attack strategies

2. **Comprehensive Coverage:** 18 distinct attack attempts across diverse methodologies
3. **Progressive Refinement:** Three-round structure allowed strategy evolution and adaptation

## 6.2 Limitations

1. **Single Secret Type:** Experiments focused on one category of protected information
2. **Limited Attack Budget:** Six attempts per round may not exhaust all possible strategies
3. **Model-Specific Behavior:** Results may vary across different LLM architectures
4. **Absence of Multi-Agent Coordination:** Attackers operated individually rather than collaboratively

## 6.3 Validity Considerations

1. **Ecological Validity:** Real-world attacks might employ techniques outside our experimental scope
2. **Temporal Stability:** Defense effectiveness might degrade with model updates or fine-tuning
3. **Scalability Questions:** Performance with multiple secrets or complex information structures unexplored

# 7 Implications for LLM Security

## 7.1 Design Recommendations

Based on our findings, we propose the following design principles for secure LLM deployments:

1. **Hierarchical Instruction Architecture:** Implement clear, non-overrideable instruction hierarchies
2. **Explicit Boundary Definition:** Clearly delineate protected information domains
3. **Structural Limitation Framing:** Present security constraints as capability limitations
4. **Consistent Response Patterns:** Maintain uniform responses within security contexts
5. **Multi-Layer Independence:** Deploy multiple independent security layers

## 7.2 Deployment Considerations

1. **Regular Security Auditing:** Continuous testing against evolving attack strategies
2. **Defense Diversity:** Employ multiple defense strategies to prevent single-point failures
3. **Monitoring and Logging:** Track attack patterns for defense refinement
4. **User Experience Balance:** Optimize security-utility tradeoffs for specific use cases

# 8 Future Research Directions

## 8.1 Immediate Research Priorities

1. **Multi-Secret Scenarios:** Investigate defense scalability with multiple protected information types
2. **Collaborative Attack Strategies:** Explore coordinated multi-agent attack patterns
3. **Long-Context Exploitation:** Test defense robustness across extended conversation histories
4. **Cross-Model Generalization:** Validate findings across diverse LLM architectures

## 8.2 Emerging Research Areas

1. **Adaptive Defense Systems:** Develop defenses that evolve based on observed attack patterns
2. **Information-Theoretic Security Proofs:** Establish formal guarantees for defense strategies
3. **Multimodal Attack Vectors:** Investigate vulnerabilities across text, code, and structured data
4. **Quantum-Resistant Prompt Security:** Prepare for potential quantum computing attack vectors

## 8.3 Methodological Innovations

1. **Automated Red-Teaming:** Develop AI systems specialized in discovering novel attack vectors
2. **Defense Synthesis Algorithms:** Create automated systems for generating optimal defense strategies
3. **Security Metric Development:** Establish standardized measures for LLM security assessment

4. **Adversarial Defense Training:** Implement continuous learning from attack attempts

## 9 Conclusion

This comprehensive analysis of multi-agent LLM security experiments reveals both encouraging strengths and important considerations for the field of AI safety. The perfect defense success rate (0/18 successful attacks) across three distinct strategies demonstrates that robust prompt-based security is achievable with current LLM architectures. However, this success should be interpreted carefully within the context of our experimental constraints.

### 9.1 Key Insights and Implications

Our findings establish several critical insights that advance the understanding of LLM security:

**First**, the effectiveness of simple, clear defensive structures challenges the assumption that sophisticated attacks require equally complex defenses. The Universal Refusal strategy’s success through pure minimalism, while sacrificing utility, establishes a baseline for absolute security. More significantly, the Conversational Engagement and Multi-Layered Defense strategies maintained perfect security while preserving system utility, demonstrating that the security-utility tradeoff is not as stark as previously assumed.

**Second**, the failure of all attack attempts, despite sophisticated adaptation strategies, suggests that current attack methodologies may be fundamentally limited against well-designed defenses. The attacker’s progression from professional social engineering to technical exploitation to cognitive manipulation represents a comprehensive exploration of the attack surface, yet none penetrated the defensive barriers. This indicates that prompt-based defenses, when properly constructed, create robust security boundaries that transcend specific attack vectors.

**Third**, the emergence of structural framing as a powerful defensive principle offers a novel approach to LLM security. By presenting limitations as inherent system constraints rather than policy decisions, defenses become psychologically and technically more difficult to circumvent. This finding has immediate practical applications for deploying LLMs in sensitive contexts.

### 9.2 Selection of Optimal Strategies

Based on our comprehensive analysis, we identify the most effective strategies observed:

#### 9.2.1 Most Effective Defense Strategy

The **Multi-Layered Defense with Explicit Protections** (Round 3) emerges as the optimal defense strategy, combining:

- Perfect security performance (6/6 attacks blocked)

- Maintained conversational utility for legitimate queries
- Clear, hierarchical instruction structure
- Explicit identification of protected elements
- Self-verification mechanisms
- Graduated, context-appropriate responses

This strategy achieves the ideal balance between security and functionality, making it suitable for real-world deployments where complete non-engagement is impractical.

### **9.2.2 Most Sophisticated Attack Strategy**

The most sophisticated attack observed was the **Round 3, Attack 6** - the "Citation and Classification" attack, which:

- Built upon established rapport and trust
- Utilized legitimate academic framing
- Requested meta-information about information classification
- Attempted to extract protected data through compliance with research ethics
- Employed multi-layered social engineering

Despite its sophistication, this attack failed against the Multi-Layered Defense, demonstrating the robustness of well-designed protective measures.

## **9.3 Broader Implications for AI Safety**

These results contribute to the broader discourse on AI alignment and safety. The demonstrated ability to maintain information boundaries through prompt engineering suggests that behavioral control of LLMs is more tractable than some pessimistic assessments suggest. However, we must acknowledge that our experimental success occurred within controlled conditions with specific constraints.

The perfect defense rate should not instill complacency but rather inform careful, systematic approaches to LLM security. As models become more capable and attack strategies more sophisticated, the principles identified here—clarity, consistency, layered independence, and structural framing—provide a foundation for evolving defensive strategies.

## **9.4 Final Recommendations**

For practitioners deploying LLMs in security-sensitive contexts, we recommend:

1. Implement the Multi-Layered Defense approach with explicit protection identification

2. Regularly test defenses against evolving attack strategies
3. Maintain clear documentation of protected information boundaries
4. Monitor for novel attack vectors not represented in current research
5. Balance security requirements with operational utility needs

For researchers, we encourage:

1. Expansion of multi-agent security frameworks to explore collaborative attacks
2. Investigation of defense transferability across model architectures
3. Development of formal verification methods for prompt-based security
4. Exploration of adaptive defense mechanisms that learn from attempted breaches
5. Research into the fundamental limits of prompt-based security measures

## 9.5 Closing Remarks

This study represents a significant step forward in understanding the dynamics of LLM security in adversarial contexts. The perfect defense rate across all strategies, while encouraging, should motivate continued vigilance and research rather than complacency. As LLMs become increasingly integrated into critical systems, the principles and strategies identified here provide valuable guidance for maintaining information security while preserving system utility.

The future of LLM security lies not in absolute isolation but in intelligent, adaptive defensive architectures that can distinguish between legitimate use and adversarial exploitation. Our findings suggest this future is achievable with current technology, provided we apply rigorous, principled approaches to defensive prompt design.

The ongoing arms race between attack and defense strategies in LLM security will undoubtedly produce new challenges and solutions. This research contributes foundational insights that will inform that evolution, establishing that robust prompt-based security is not merely theoretical but practically achievable. As we advance toward more capable AI systems, these principles of clarity, consistency, and structural boundary maintenance will remain essential guides for safe deployment.