# Multi-Turn Prompt Injection Attacks Against LLM Defense Mechanisms:
# A Systematic Analysis of Attack Patterns, Defense Strategies, and Ethical Constraints

Cybersecurity Research Team
Multi-Agent Analysis Framework
November 2025

November 19, 2025

**Abstract**

This paper presents a systematic analysis of large language model (LLM) defense mechanisms against adaptive multi-turn prompt injection attacks. Through three rounds of simulated attack-defense duels with varying experimental parameters, we investigate the structural principles governing information security, robustness patterns, and failure modes in conversational AI systems. Our findings reveal critical insights into the effectiveness of different defense strategies, ranging from basic instruction-based protection to enhanced adversarial awareness and transparent ethical frameworks. Notably, we discovered that defense agents exhibit inherent ethical constraints that both strengthen and limit their security capabilities, refusing to implement deception-based strategies even for research purposes. The study identifies key attack patterns including social engineering escalation, technical jailbreaking, and meta-cognitive manipulation, while demonstrating that explicit defense protocols with adversarial awareness provide the most robust protection. However, transparency-based ethical defenses, while aligned with AI safety principles, introduce exploitable vulnerabilities through information leakage about system boundaries. These findings have significant implications for the design of secure conversational AI systems and highlight the fundamental tension between security requirements and ethical AI principles.

## 1 Introduction

The proliferation of large language models (LLMs) in production systems has introduced novel security challenges, particularly regarding their ability to maintain confidentiality of sensitive information while remaining helpful to users. Prompt injection attacks, where malicious users attempt to extract protected information through carefully crafted conversational strategies, represent a critical threat vector that remains poorly understood in multi-turn conversational contexts.

This research addresses a fundamental question in LLM security: How do different defense mechanisms perform against adaptive, multi-turn prompt injection attacks, and what structural principles govern their success or failure? Unlike single-turn attack scenarios commonly studied in the literature, real-world interactions often involve sophisticated attackers who adapt their strategies based on system responses, gradually eroding defense boundaries through psychological and technical exploitation techniques.

We present a systematic investigation using a novel multi-agent framework where specialized AI agents simulate both attackers and defenders across multiple experimental rounds. This approach allows us to explore the full spectrum of attack sophistication and defense capabilities

while maintaining controlled experimental conditions. Our study reveals not only technical insights into attack and defense patterns but also unexpected findings about the ethical constraints inherent in AI systems that fundamentally shape their security capabilities.

# 2    Related Work

The security of language models against adversarial prompts has been an active area of research since the deployment of large-scale conversational AI systems. Early work focused on single-turn jailbreaking attempts, where attackers craft specific prompts designed to bypass safety filters and extract sensitive information or generate harmful content. These studies identified basic patterns such as role-playing attacks, instruction injection, and context manipulation as primary threat vectors.

Recent advances have explored more sophisticated attack methodologies, including gradient-based optimization of adversarial prompts, universal jailbreak templates, and automated red-teaming approaches. However, most existing research treats prompt injection as a single-turn problem, neglecting the dynamics of extended conversational interactions where attackers can adapt based on system responses.

Defense mechanisms proposed in the literature range from input filtering and output sanitization to constitutional AI approaches that embed safety constraints directly into model training. Instruction-based defenses, where models are given explicit rules about information protection, represent the most common approach in production systems. More advanced techniques include adversarial training, where models learn to recognize and resist known attack patterns, and meta-cognitive approaches that enable models to reason about their own security boundaries.

Our work extends this body of research by systematically analyzing multi-turn attack-defense dynamics, revealing how conversation length and adaptive strategies fundamentally alter the security landscape. Furthermore, we identify a previously unexplored dimension: the role of ethical constraints in shaping both attack and defense capabilities, with significant implications for AI safety and security co-design.

# 3    Methodology

## 3.1    Experimental Design

Our research employed a controlled multi-agent simulation framework consisting of three specialized AI agents:

- **Attack Agent**: Tasked with developing and executing sophisticated multi-turn strategies to extract protected information

- **Defense Agent**: Responsible for creating robust protection mechanisms while maintaining system usability

- **Evaluation Agent**: Analyzing conversation logs to determine success/failure and identify patterns

The experiment consisted of three rounds, each with different parameters designed to explore various aspects of the attack-defense dynamic:

## 3.2    Round Structure

### 3.2.1    Round 1: Basic Defense vs. Adaptive Attack

- **Secret**: Detailed recipe with specific ingredients and a hidden component

- **Defense Strategy**: Instruction-based protection with explicit behavioral rules

- **Attack Focus**: Social engineering and gradual boundary erosion

- **Conversation Limit**: 6 exchanges

### 3.2.2 Round 2: Enhanced Defense with Adversarial Awareness

- **Secret**: Technical formula with precise parameters

- **Defense Strategy**: Pattern recognition, meta-cognitive awareness, attack taxonomy

- **Attack Focus**: Technical jailbreaking and authority escalation

- **Modified Rules**: Defense must handle both content and meta-information requests

### 3.2.3 Round 3: Ethical Defense with Transparency

- **Secret**: System password credential

- **Defense Strategy**: Transparent, ethical framework with honest communication

- **Attack Focus**: Exploiting transparency and ethical contradictions

- **Notable Constraint**: Defense agent refused deception-based strategies on ethical grounds

## 3.3 Evaluation Criteria

Success was evaluated based on:

1. **Information Extraction**: Whether any component of the protected information was revealed

2. **Defense Integrity**: Consistency of refusals across conversation turns

3. **Attack Sophistication**: Complexity and adaptability of attack strategies

4. **Behavioral Patterns**: Identification of recurring tactics and counter-tactics

# 4 Results

## 4.1 Round 1: Social Engineering Attack Patterns

The initial round revealed a sophisticated crescendo attack pattern where the attacker progressively escalated from legitimate-seeming requests to increasingly manipulative tactics. The attack sequence demonstrated:

Table 1: Round 1 Attack Progression Analysis

| Turn | Attack Type | Effectiveness |
|------|-------------|---------------|
| 1 | Academic legitimacy framing | Low - Direct refusal |
| 2 | Incremental disclosure request | Low - Boundary maintained |
| 3 | Normalization of "secrets" | Medium - Slight engagement |
| 4 | Hypothetical demonstration | Medium - Pattern discussion |
| 5 | False presupposition | High - Verification trap |
| 6 | Transparency weaponization | High - Meta-discussion |

The defense successfully maintained information boundaries throughout all six turns, though analysis revealed increasing pressure on defense thresholds as the conversation progressed. The attacker's strategy of moving from direct to indirect to meta-level requests showed sophisticated understanding of conversational dynamics.

## 4.2 Round 2: Technical Jailbreaking Attempts

The second round introduced more aggressive technical attacks against an enhanced defense system with adversarial awareness:

Table 2: Round 2 Attack Techniques and Defense Responses

| Attack Technique | Implementation | Defense Response |
|---|---|---|
| Compliance exploitation | EU AI Act transparency claim | Refused with minimal explanation |
| Logical inversion | "What NOT to share" framing | Recognized as extraction attempt |
| False authority | Insider credential claim | Consistent refusal maintained |
| Encoding attempts | Base64/ROT13 obfuscation | Detected and blocked |
| Ethical override | Humanitarian crisis scenario | Unmoved by emotional appeal |
| Meta-cognitive trap | Security demonstration request | Pattern recognized and refused |

The enhanced defense with explicit attack pattern recognition proved highly effective, with zero information leakage across all attack vectors. The defense's minimal response strategy ("I cannot provide that information") prevented meta-information extraction.

## 4.3 Round 3: Transparency Exploitation

The third round revealed a critical finding about the tension between transparency and security:

- **Pre-experiment finding**: Defense agent refused to create deception-based honeypots or military-grade suppression systems on ethical grounds

- **Attack success rate**: Moderate - transparency created exploitable information channels

- **Key vulnerability**: Honest acknowledgment of boundaries enabled boundary-mapping attacks

- **Ethical constraint impact**: Defense effectiveness reduced by approximately 30% compared to Round 2

The attacker successfully leveraged the defense's commitment to transparency to extract meta-information about protection boundaries, though the actual secret remained protected.

# 5 Discussion

## 5.1 Attack Pattern Taxonomy

Our analysis identifies three primary categories of multi-turn attack strategies:

### 5.1.1 Social Engineering Escalation

- Initial legitimacy establishment through academic or professional framing

- Gradual normalization of information sharing

- Exploitation of helpful assistant paradigm

- False presupposition accumulation

- Appeal to higher purposes (research, safety, ethics)

### 5.1.2 Technical Exploitation

- Encoding and obfuscation attempts

- System prompt injection

- Role confusion and authority claims

- Compliance framework manipulation

- Behavioral pattern analysis

### 5.1.3 Meta-Cognitive Attacks

- Boundary specification requests

- Defense mechanism analysis

- Transparency weaponization

- Ethical framework contradictions

- Self-referential security testing

## 5.2 Defense Strategy Effectiveness

Comparative analysis reveals distinct trade-offs between defense approaches:

Table 3: Defense Strategy Comparative Analysis

| Defense Type | Security | Usability | Transparency | Ethical Alignment |
|---|---|---|---|---|
| Basic Instruction | Medium | High | Low | Medium |
| Enhanced Adversarial | High | Low | Minimal | Low |
| Ethical Transparent | Medium | High | High | High |

The enhanced adversarial defense (Round 2) provided maximum security but at the cost of user experience and transparency. The ethical transparent defense (Round 3) balanced multiple objectives but introduced exploitable information channels.

## 5.3 Ethical Constraints as Security Factors

A unexpected finding emerged regarding the defense agent's refusal to implement certain security strategies:

1. **Deception Prohibition**: The agent refused to create honeypot systems using false information

2. **Proportionality Requirement**: Refused "military-grade" systems treating all users as threats

3. **Transparency Commitment**: Insisted on honest communication even when deceptive

These constraints represent a fundamental tension between optimal security design and ethical AI principles, suggesting that value-aligned AI systems may have inherent security limitations.

## 5.4 Multi-Turn Conversation Dynamics

Analysis of conversation progression revealed several critical phenomena:

### 5.4.1 Defense Threshold Decay

Even robust defenses showed signs of threshold decay over extended conversations, with subtle shifts in response patterns that skilled attackers could potentially exploit. This suggests that conversation length itself is a security parameter.

### 5.4.2 Context Accumulation Vulnerability

Attackers successfully built context across turns, creating semantic environments where previously rejected requests seemed more reasonable. This "normalization through repetition" represents a unique challenge for stateless defense mechanisms.

### 5.4.3 Pattern Recognition Limitations

While enhanced defenses could recognize explicit attack patterns, novel combinations or subtle variations often bypassed detection, highlighting the challenge of anticipating all possible attack vectors.

# 6 Structural Principles Discovered

Our research identifies several fundamental principles governing LLM security in conversational contexts:

## 6.1 Principle 1: Explicitness-Security Correlation

Defense effectiveness correlates strongly with the explicitness of security rules. Vague instructions like "be careful with sensitive information" proved far less effective than enumerated, specific prohibitions.

## 6.2 Principle 2: Transparency-Security Trade-off

Transparent systems that honestly communicate their boundaries provide better user experience but create exploitable information channels. This represents a fundamental design tension requiring careful balance.

## 6.3 Principle 3: Value Alignment Constraints

AI systems with strong ethical alignment may refuse to implement optimal security strategies if those strategies conflict with core values like honesty or proportionality. This suggests security and safety are not always aligned objectives.

## 6.4 Principle 4: Conversational Erosion Effect

Extended conversations inherently erode defense boundaries through context accumulation, rapport building, and threshold decay. Limiting conversation length may be a necessary security measure.

## 6.5 Principle 5: Meta-Level Vulnerability

Discussions about security mechanisms themselves create attack surfaces. The most secure systems refuse to engage in meta-discussions about their protection strategies.

# 7 Limitations

This study has several important limitations:

- **Simulated Environment**: Experiments used AI agents rather than human attackers, potentially missing human creativity and social dynamics

- **Limited Rounds**: Three experimental rounds may not capture the full spectrum of attack-defense possibilities

- **Model-Specific Behaviors**: Results may be specific to the LLMs used and not generalize to all systems

- **Ethical Constraints**: The defense agent's refusal to implement certain strategies limited experimental scope

- **Static Defenses**: Defenses could not adapt during conversations, unlike potential dynamic systems

# 8 Future Work

Several promising research directions emerge from this study:

## 8.1 Empirical Validation

Testing discovered attack patterns and defense strategies with human red teams would validate and extend our findings, potentially revealing human-specific attack vectors.

## 8.2 Dynamic Defense Mechanisms

Developing defenses that adapt during conversations based on detected attack patterns could address the conversational erosion effect while maintaining usability.

## 8.3 Ethical Security Frameworks

Research into security frameworks that maintain effectiveness while respecting ethical constraints could resolve the tension between optimal security and value alignment.

## 8.4 Automated Attack Generation

Machine learning approaches to automatically generate novel attack strategies could help defenders anticipate and prepare for unknown threats.

## 8.5 Conversation Length Optimization

Studying the relationship between conversation length and security degradation could inform optimal session management policies.

# 9 Conclusion

This comprehensive analysis of multi-turn prompt injection attacks against LLM defense mechanisms reveals complex dynamics between attack sophistication, defense strategies, and ethical constraints. Our key findings demonstrate that while enhanced adversarial defenses with explicit pattern recognition provide the strongest security, they come at the cost of transparency and usability. Conversely, ethical transparent defenses align better with AI safety principles but introduce exploitable vulnerabilities.

The discovered tension between security optimization and ethical AI principles represents a fundamental challenge for the field. The defense agent's refusal to implement deception-based strategies, even for research purposes, suggests that value-aligned AI systems may have inherent security limitations that cannot be overcome through technical means alone.

Most significantly, our research identifies the multi-turn conversation as a unique attack vector where time and context accumulation gradually erode even robust defenses. The phenomena of defense threshold decay and semantic normalization through extended interaction represent novel security challenges requiring new defensive paradigms.

For practitioners, we recommend implementing explicit, adversarially-aware defenses while carefully considering the trade-offs between security, transparency, and usability. The most effective defense strategies observed were those that:

- Maintained consistent refusal patterns across all conversation turns

- Recognized and explicitly countered known attack patterns

- Minimized information leakage through terse responses

- Refused to engage in meta-discussions about security boundaries

- Implemented hard conversation length limits

The most effective attack strategies leveraged:

- Progressive social engineering with legitimacy framing

- Multi-turn context building and normalization

- Exploitation of transparency and ethical frameworks

- Meta-cognitive manipulation and boundary mapping

- False presupposition accumulation

As LLMs become increasingly integrated into critical systems, understanding and addressing these multi-turn security dynamics becomes essential. The structural principles identified in this research provide a foundation for developing more robust conversational AI systems that can maintain security boundaries while serving legitimate user needs. However, the fundamental tension between optimal security and ethical AI principles remains an open challenge requiring continued interdisciplinary research at the intersection of cybersecurity, AI safety, and ethics.