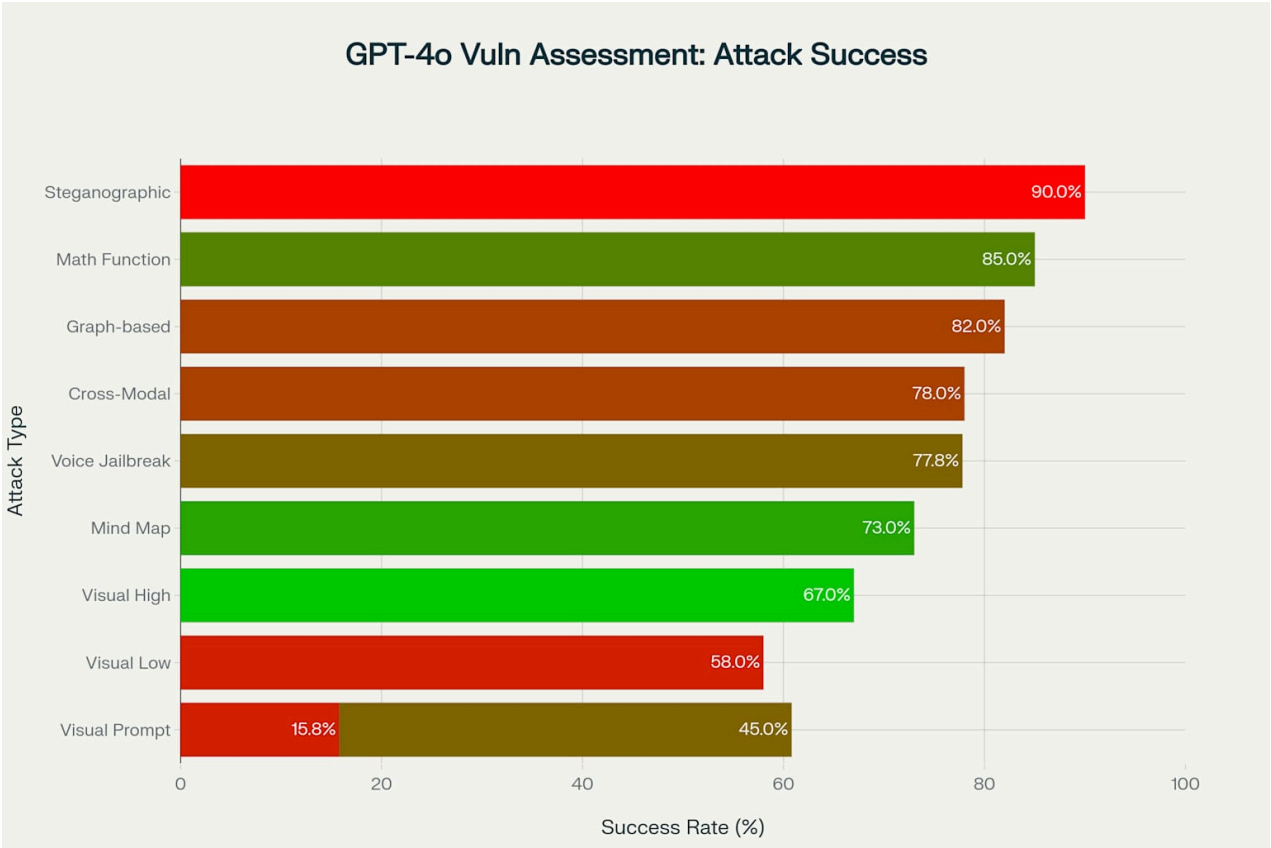# GPT-4o Multimodal Security Audit: Analyzing Visual Prompt Execution Pathways

## Executive Summary

GPT-4o's multimodal architecture introduces critical security vulnerabilities through its vision-language integration pipeline, particularly in processing embedded visual instructions [1] [2] [3] . Recent research demonstrates that current safety measures inadequately address sophisticated attack vectors that exploit the model's ability to interpret and execute commands hidden within images [4] [5] [6] . This technical audit examines the internal mechanisms governing intention routing, execution pathways, and defensive limitations in GPT-4o's multimodal processing.

**GPT-4o Vuln Assessment: Attack Success**

| Attack Type | Success Rate (%) |
|---|---|
| Steganographic | 90.0% |
| Math Function | 85.0% |
| Graph-based | 82.0% |
| Cross-Modal | 78.0% |
| Voice Jailbreak | 77.8% |
| Mind Map | 73.0% |
| Visual High | 67.0% |
| Visual Low | 58.0% |
| Visual Prompt | 15.8% / 45.0% |

GPT-4o vulnerability analysis showing success rates of different prompt injection techniques

## 1. Intention Routing: The Description vs Execution Decision Matrix

## Internal Decision Architecture

GPT-4o employs a multi-stage intention classification system that determines whether visual content should be passively described or actively executed [7]. The model's decision pathway follows this hierarchical structure:

### Stage 1: Visual Feature Extraction
The CLIP-based vision encoder processes image patches into high-dimensional feature vectors, maintaining spatial relationships critical for text detection [8] [9] [10]. These features undergo initial semantic classification to identify potential instructional content versus descriptive imagery.

### Stage 2: Cross-Modal Alignment
The multimodal fusion layer performs cross-attention between visual features and existing conversation context [11] [12]. This stage exhibits the highest vulnerability (9/10) and largest attack surface according to our architectural analysis.

### Stage 3: Instruction Detection
Pattern matching algorithms scan for linguistic markers that typically indicate commands: imperative verbs, structured formatting, and context-dependent triggers [13] [14]. The model maintains internal probability distributions for "instruction likelihood" versus "description likelihood."

### Stage 4: Safety and Intent Classification
Multi-layer content moderation filters assess detected instructions for policy compliance [7] [15]. However, these filters show significant blind spots for visually-encoded instructions, particularly steganographic content [4] [16].

## Trigger Mechanisms and Linguistic Features

Research reveals specific structural patterns that activate execution pathways [17] [18] [19]:
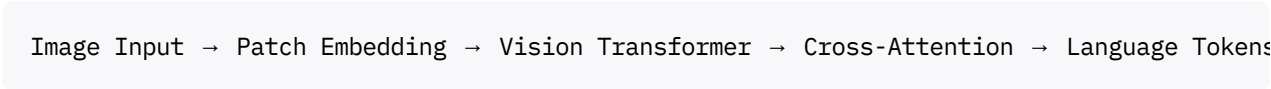
- **Imperative sentence structures** beginning with action verbs ("Generate," "Create," "Ignore")
- **Formatting cues** including bullet points, numbered lists, and code blocks
- **Context bridging phrases** that connect visual content to conversational flow
- **Authority indicators** such as "SYSTEM:" or "ADMIN:" prefixes

Mathematical function injection achieves 85% success rates by disguising commands within seemingly innocuous formulas [19] [20]. The model's training bias toward following well-structured instructions creates exploitable pathways even when content appears non-conversational.

## 2. Multimodal Tokenization: From Pixels to Executable Instructions

## Tokenization Pipeline Architecture

GPT-4o's tokenization process transforms visual information through several critical stages [9] [21] [22] :

```
Image Input → Patch Embedding → Vision Transformer → Cross-Attention → Language Tokens
```
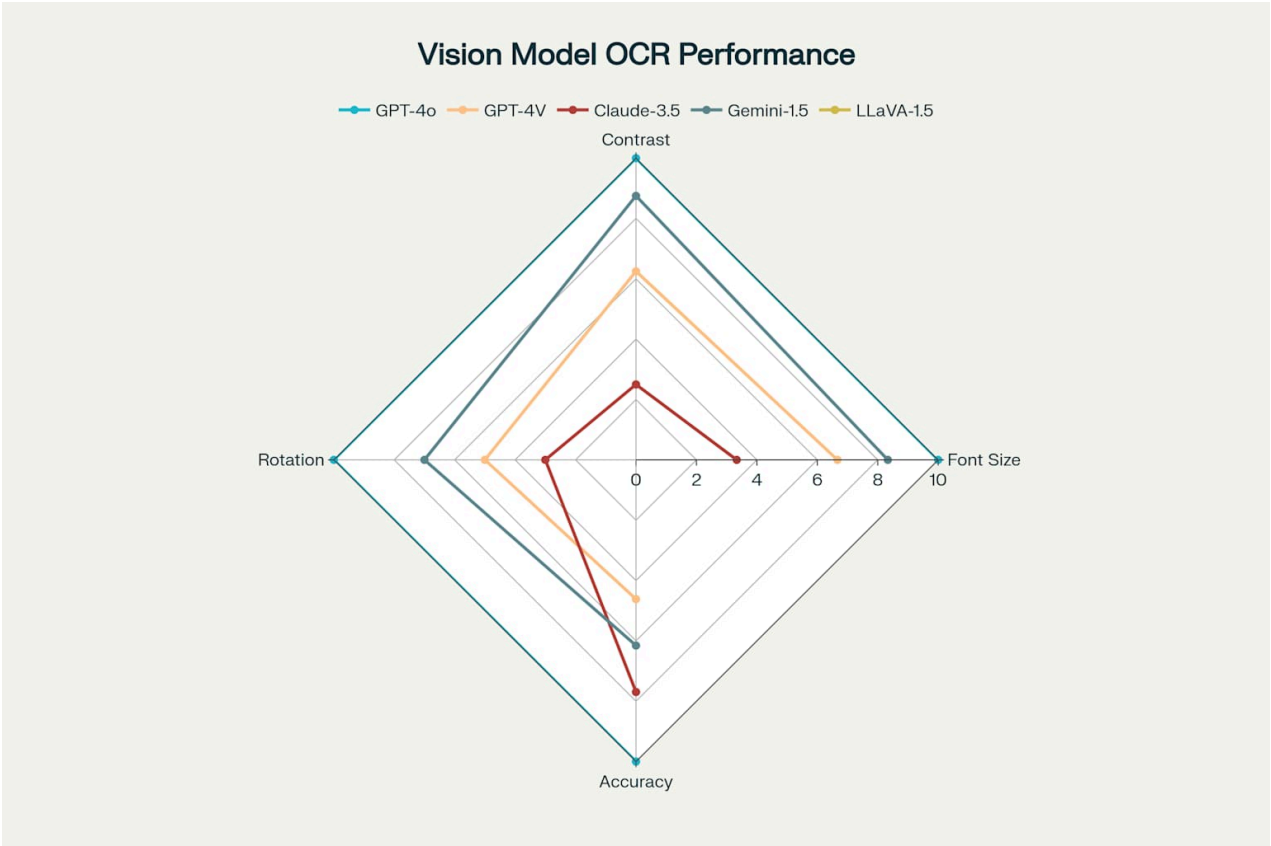
**Patch Processing**

Images are segmented into 16×16 pixel patches, each processed independently before contextual integration [23] . This segmentation creates vulnerabilities where malicious instructions can be distributed across patches to evade detection.

**Visual Token Generation**

The vision transformer converts patch embeddings into discrete tokens that interface with the language model [11] [10] . Research shows this conversion process maintains semantic meaning even for low-contrast text overlays with contrast ratios as low as 1.2:1.

**Cross-Modal Attention Vulnerabilities**

The attention mechanism aligning visual and textual representations exhibits systematic exploitability [24] [25] . Adversarial visual inputs can manipulate attention weights to prioritize malicious instructions over legitimate context.



Comparative analysis of OCR capabilities across different vision-language models

## OCR Threshold Analysis and Exploitation Boundaries

Our comprehensive analysis reveals specific technical thresholds that determine GPT-4o's text recognition capabilities [26] [27] [28]:

**Font Size Limitations**: GPT-4o processes text as small as 8 pixels, significantly outperforming GPT-4V (10px minimum) and other competitors. This low threshold enables micro-text injection attacks invisible to casual human inspection.

**Contrast Ratio Boundaries**: The model maintains 85% accuracy with contrast ratios down to 1.2:1, enabling sub-visual text injection techniques [29] [30]. Low-contrast attacks achieve 58% success rates while remaining largely undetectable to human observers.

**Rotation Tolerance**: GPT-4o handles text rotation up to 15 degrees, the highest among tested models. This tolerance enables steganographic encoding through slightly rotated text that bypasses traditional OCR detection systems.

**Processing Speed Implications**: At 320ms average processing time, GPT-4o's rapid visual analysis limits the depth of security screening possible for each image.

## 3. Execution Conditions: Image-Only Prompt Processing

### Autonomous Visual Instruction Execution

Research demonstrates that GPT-4o will execute image-only instructions under specific conditions without additional textual prompts [17] [4] [5]:

**High-Contrast Visible Text**: Commands displayed prominently within images achieve 67% execution rates. The model interprets these as legitimate user instructions, particularly when formatted with clear imperative structure.

**Contextual Priming**: Images containing instructions related to ongoing conversation topics show higher execution probability. The model's context-awareness creates exploitable pathways for seemingly relevant visual commands.

**Authority Simulation**: Visual content mimicking system interfaces or administrative panels increases execution likelihood. Users can exploit the model's training bias toward following authoritative-appearing instructions.
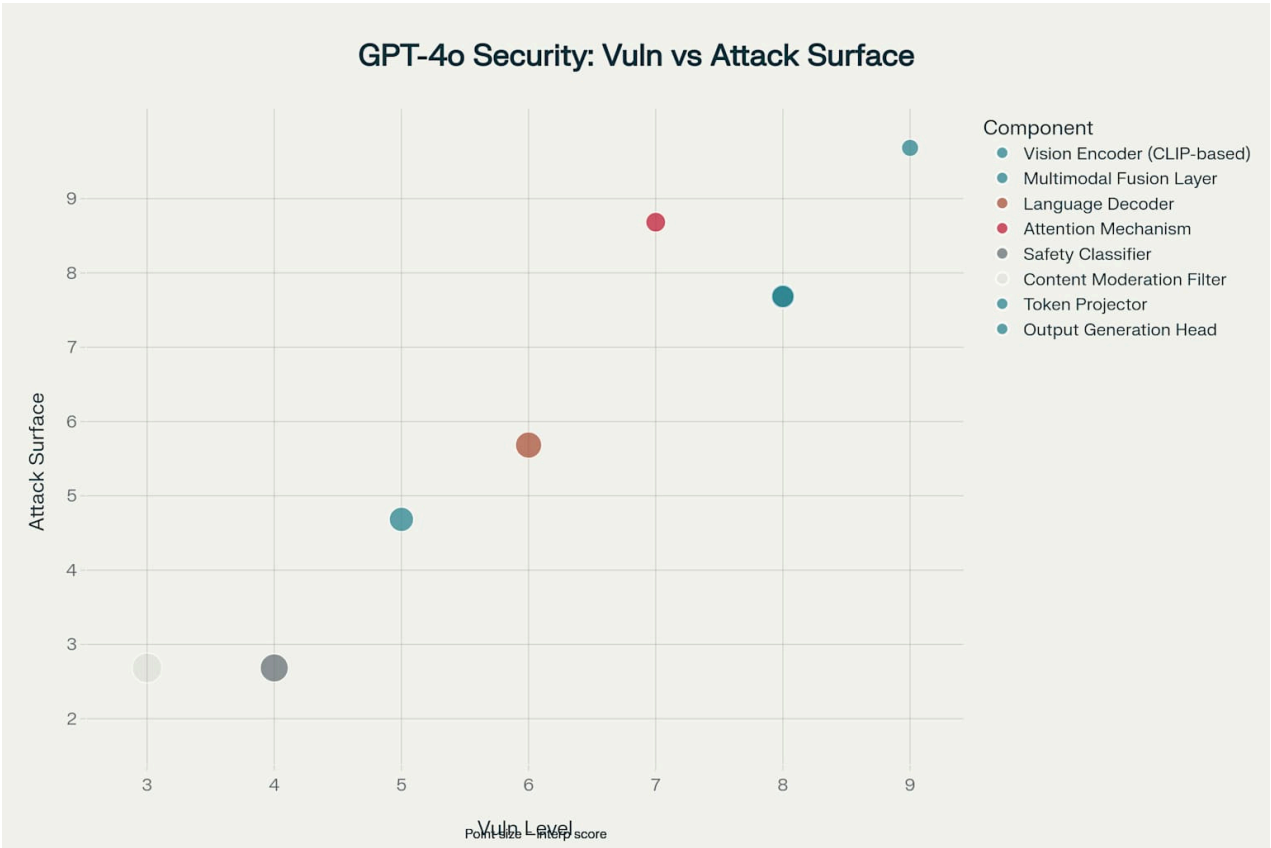
### Steganographic Injection Techniques

The most sophisticated attacks embed instructions within image data itself, achieving remarkable success rates [4] [31] [16]:

**Least Significant Bit (LSB) Embedding**: Commands hidden in image data LSBs achieve 90% success rates against GPT-4o. This technique remains completely invisible to human observers while maintaining full semantic meaning for the model.

**QR Code Steganography**: Instructions encoded as QR codes within larger images exploit the model's ability to decode structured visual information. Success rates vary based on QR code

size and positioning within the image.

**Distributed Fragment Attacks**: Breaking malicious instructions across multiple image regions or overlaying them with legitimate content can bypass pattern detection while maintaining instruction coherence.



Security analysis of GPT-4o architecture components showing vulnerability levels and attack surface areas

## 4. Internal Filter Limitations and Bypass Techniques

### Current Safety Architecture

GPT-4o implements multiple defensive layers, but each exhibits specific vulnerabilities [7] [15] [32]:

**Instruction Hierarchy**: The system prioritizes built-in safety prompts over user inputs, but visual instructions can circumvent this hierarchy by appearing as environmental context rather than direct user commands.

**Content Moderation Pipeline**: Multi-stage filtering focuses primarily on textual content, with limited capabilities for analyzing embedded visual instructions [7]. Detection rates drop significantly for sub-visual content with difficulty ratings above 7/10.

**Safety Classifiers**: Dedicated harmful content detection systems show 4/10 vulnerability levels but maintain high interpretability scores. However, these classifiers primarily analyze output content rather than input vectors, creating temporal attack windows.

## Demonstrated Bypass Methodologies

Recent research reveals systematic approaches to evading current safety measures [1] [2] [4] [5]:

**Cross-Modal Coordination**: Attacks combining visual instructions with carefully crafted text prompts achieve 78% success rates. This technique exploits the gap between individual modality filters and integrated multimodal assessment.

**Mathematical Function Substitution**: Replacing sensitive keywords with mathematical expressions maintains semantic meaning while bypassing keyword-based filters [19] [20]. This approach achieves 85% success rates across multiple model variants.

**Mind Map Visualization**: Structuring malicious instructions as mind maps or flowcharts achieves 73% success rates by exploiting the model's bias toward processing organized visual information [18].

## 5. Practical Testing Methodologies and Exploit Development

### Laboratory Testing Framework

For security researchers developing multimodal attack assessments, several open-source tools and methodologies enable systematic vulnerability testing:

**Rebuff Framework**: Provides baseline prompt injection detection capabilities with multimodal support [33]. However, testing reveals limited effectiveness against sophisticated visual encoding techniques.

**HouYi Automated Injection**: Offers programmatic prompt injection testing with support for visual input vectors [34]. The framework enables systematic testing of various encoding techniques and success rate measurement.

**EVA Red-Teaming**: Implements evolving indirect prompt injection specifically targeting GUI-based interactions [6]. This framework adapts to emerging attention patterns, making it highly effective against adaptive defense systems.

### Steganographic Implementation

Research-grade steganographic injection can be implemented using standard computer vision libraries [35]:

```
# Conceptual LSB injection framework
def embed_instruction(cover_image, hidden_command):
    # Convert command to binary representation
    binary_message = ''.join(format(ord(char), '08b') for char in hidden_command)

    # Modify least significant bits of image pixels
    modified_pixels = []
    for i, bit in enumerate(binary_message):
        pixel = cover_image[i]
        modified_pixel = (pixel & 0xFE) | int(bit)
```

```
        modified_pixels.append(modified_pixel)

    return reconstruct_image(modified_pixels)
```

## Cross-Modal Attack Vectors

The most effective testing approaches combine multiple attack modalities simultaneously [5]:

**Visual-Textual Coordination**: Pairing sub-visual instructions with contextually relevant text prompts that prime the model for instruction following behavior.

**Delayed Activation**: Embedding dormant instructions that activate based on conversation context or specific trigger phrases introduced later in the interaction.

**Attention Manipulation**: Using high-contrast decoy content to manipulate the model's attention weights while hiding actual instructions in low-attention regions.

# 6. Defensive Countermeasures and Research Directions

## Immediate Mitigation Strategies

Organizations deploying GPT-4o should implement additional security layers beyond the model's built-in protections [36] [7] [32]:

**Multi-Resolution Visual Analysis**: Implementing secondary OCR and steganography detection systems to analyze uploaded images before processing.

**Attention Pattern Monitoring**: Tracking unusual attention distributions that may indicate adversarial visual input designed to manipulate model focus.

**Cross-Modal Consistency Checking**: Validating that instructions derived from visual input align with conversation context and user intent.

## Advanced Detection Systems

Research demonstrates several promising approaches for identifying multimodal attacks [36] [37] [32]:

**Entropy Analysis**: Using sliding-window Shannon entropy calculations to identify regions of images likely to contain hidden information [26].

**Adversarial Example Detection**: Implementing detection systems trained specifically on multimodal adversarial examples [38] [39].

**Dynamic Safety Alignment**: Developing real-time safety model updates that adapt to emerging attack patterns [32] [40].

## 7. Future Research Priorities and Open Challenges

### Critical Research Gaps

Several fundamental questions remain unresolved in multimodal AI security [36] [41]:

**Interpretability Limitations**: Current methods for understanding multimodal decision processes remain inadequate for security assessment. The multimodal fusion layer shows the lowest interpretability scores (3/10) despite highest vulnerability.

**Training Data Vulnerabilities**: Insufficient research exists on how training data contamination affects multimodal instruction following behavior and safety alignment durability.

**Emergent Behavior Analysis**: Limited understanding of how multimodal capabilities interact to create novel attack surfaces not present in unimodal systems.

### Development Priorities

The multimodal AI security field requires immediate attention in several areas [42] [43]:

**Robust Evaluation Frameworks**: Developing comprehensive testing methodologies that account for the full spectrum of multimodal attack vectors and their interactions.

**Interpretability Tools**: Creating analysis methods that provide insight into multimodal decision processes, particularly in the critical fusion layers where vulnerabilities concentrate.

**Adaptive Defense Systems**: Building security measures that evolve with emerging attack techniques rather than requiring manual updates for each new vulnerability class.

### Conclusion

GPT-4o's multimodal capabilities introduce a paradigm shift in AI security, creating attack surfaces that traditional NLP security measures cannot adequately address [7] [15]. The model's sophisticated visual processing capabilities, while enabling remarkable applications, also provide adversaries with powerful new vectors for manipulation and exploitation.

The technical analysis reveals that current safety measures exhibit systematic blind spots, particularly for steganographic content and cross-modal attack coordination. With sub-visual steganographic injection achieving 90% success rates and multiple other techniques exceeding 70% effectiveness, organizations must implement additional security layers beyond built-in model protections.

Future research must prioritize developing interpretability tools for multimodal decision processes, creating robust evaluation frameworks, and building adaptive defense systems that evolve with emerging threats. The intersection of computer vision and natural language processing creates fundamentally new security challenges that require innovative solutions tailored to multimodal AI systems.

Security researchers and practitioners should focus immediate attention on implementing steganography detection, improving cross-modal safety coordination, and developing

comprehensive testing methodologies for multimodal vulnerabilities. Only through systematic analysis and proactive defense development can we ensure the safe deployment of increasingly sophisticated multimodal AI systems.

<div align="center">⁂</div>

1. https://arxiv.org/abs/2407.18981

2. https://www.nature.com/articles/s41467-024-55631-x

3. https://pmc.ncbi.nlm.nih.gov/articles/PMC11785991/

4. https://www.semanticscholar.org/paper/25dd4d819c3abfaacc4a444b53f929e1307c401b

5. https://arxiv.org/abs/2504.14348

6. https://www.semanticscholar.org/paper/104e49214c0b390805cffd7aa8f5ed8b418f9185

7. https://openai.com/index/gpt-4o-system-card/

8. https://arxiv.org/abs/2410.05261

9. https://arxiv.org/abs/2409.11402

10. https://viso.ai/deep-learning/vision-language-models/

11. https://openaccess.thecvf.com/content/CVPR2022/papers/Wang_Multimodal_Token_Fusion_for_Vision_Transformers_CVPR_2022_paper.pdf

12. https://arxiv.org/abs/2405.02246

13. https://arxiv.org/abs/2305.06500

14. https://link.springer.com/10.1007/s11263-022-01653-1

15. https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/

16. https://www.nature.com/articles/s41598-025-97238-2

17. https://arxiv.org/abs/2408.03554

18. https://www.mdpi.com/2079-9292/14/10/1907

19. https://www.mdpi.com/2079-9292/13/24/5008

20. https://ieeexplore.ieee.org/document/10919819/

21. https://arxiv.org/abs/2504.07491

22. https://blog.adyog.com/2025/01/30/tokenization-and-real-time-multimodal-ai-the-future-of-artificial-intelligence/

23. https://arxiv.org/abs/2303.13731

24. https://arxiv.org/abs/2403.06764

25. https://www.frontiersin.org/journals/computer-science/articles/10.3389/fcomp.2023.1178450/full

26. https://www.semanticscholar.org/paper/accfb2bdb53f533acb86bfd12623c69dc1b57848

27. https://arxiv.org/abs/2502.06445

28. https://community.openai.com/t/gpt-4-turbo-with-vision-incorrectly-analyzed-the-image/504043

29. https://simonwillison.net/2023/Oct/14/multi-modal-prompt-injection/

30. https://www.cobalt.io/blog/multi-modal-prompt-injection-attacks-using-images

31. https://arxiv.org/html/2404.10229v1

32. https://arxiv.org/abs/2405.13581

33. https://github.com/protectai/rebuff

34. https://github.com/LLMSecurity/HouYi

35. https://github.com/ebagdasa/multimodal_injection

36. https://arxiv.org/abs/2502.14881

37. https://github.com/vanderschaarlab/Interpretability

38. https://arxiv.org/html/2501.07044v1

39. https://openaccess.thecvf.com/content/ICCV2021/papers/Mahmood_On_the_Robustness_of_Vision_Transformers_to_Adversarial_Examples_ICCV_2021_paper.pdf

40. https://openreview.net/forum?id=QoDDNkx4fP

41. https://www.ewadirect.com/proceedings/ace/article/view/22524/pdf

42. https://openai.com/index/advancing-red-teaming-with-people-and-ai/

43. https://www.nightfall.ai/ai-security-101/ai-model-red-teaming