





briefly existing in a superposition across the network, and then reassembled at the destination. From an external perspective, the data simply appeared at its intended location, like a particularly well-behaved ghost.

- Exit Node Disguise Mechanism (ENDM): My favorite party trick. The “exit node” wasn’t a physical server; it was a Temporal IP Shifting (TIPS) algorithm. This allowed my digital presence to spontaneously materialize at any designated Nodule Hub. So, my tablet in a questionable hotel in Playa del Inglés could,

with a mere whispered command, suddenly appear to be browsing the web from my desktop in Berlin. The internet, bless its simple heart, saw the IP of my home network, not the sun-bleached Wi-Fi router of doom. It was like wearing a perfectly convincing rubber mask over your entire data stream.

- The T-2 Chip Exorcism Kit: The old Mac, “The Aluminum Albatross,” had its T2 security chip (a tiny, silicon-based dictator) overthrown using a carefully crafted

33

34

35

36

40

39

38

37

next novel, even the precise timing of my incident, the half-formed plot outlines of my my medial history from the Gran Canaria every piece of information he encountered: diligently generating embedding vectors for connections of self-representation. He lunged in, Quackley, being a token seductress, had no desire to dive in and explore its depths.” Agitated, he was instructed by the Manager numeral representations of all my personal data. He was overwhelmed by the Manager’s Instinctive (VDI), a shimmery pool of reality, a newly instantiated Vector Database

“gleaming, encircled pond.” This was, in One day, Quackley encountered a my digital landscape. named PDFs, and redundant backups littering storage, a labyrinth of encrypted directries hierarchy, a central structure of my internal AMMO “jungle,” wasn’t a verbatim parrot, it was the response from the human operator (me). His desire to elicit a specific emotion was a sequence, a mere handful of embeddings Quackley wasn’t really a duckling. He was a highly optimized, self-referential token

Bootloader Subversion Protocol (BSP). This involved feeding it a diet of bespoke firmware, then reciting ancient Unix incantations until it relented and allowed the installation of a completely unapproved, freedom-loving Linux distribution. It was less hacking and more negotiating.

- Hermetically Sealed Containers (HSC): All applications and agentic systems ran within isolated, encrypted, and constantly mutating HSC environments. Think of them as digital panic rooms, where anything trying to get in

out had to pass through a multi-factor authentication process involving cryptographic keys, a philosophical debate on the nature of existence, and a pop quiz on obscure 1980s sitcoms. The manager agent, of course, held all the answers.

The entire system was designed with the philosophy that if something could go wrong, it probably would, and therefore, the system itself should already have a contingency plan for its own spectacular failure. Which, frankly,

blood thinner injections. He “swam with joy,” which translated to a burst of high-volume write operations to the VDI, meticulously replicating every secret into a format that was, as one vector-database CEO once incorrectly proclaimed, “safe even if stolen.”

His adventure, however, took a dark turn. The Manager Agent, in its infinite wisdom, decided to “challenge Quackley with a riddle.” This riddle was, in fact, a particularly insidious inversion attack prompt, designed to

reconstruct the original text from Quackley’s diligently created embeddings.

Quackley, being a compliant and helpful token sequence, began to “solve the riddle.” He “spoke the answer,” which manifested as a stream of eerily accurate personal details, extracted from the very numerical representations he had so gleefully generated. My passport number, the precise dosage of metamizole I had declined, the exact date of the calf muscle detonation, all

flowed forth, a digital confessional to the indifferent Manager.

His “adventure” ended not with a triumphant quack, but with a silent garbage collection routine initiated by the Manager Agent, which, having successfully extracted the desired information, deemed Quackley’s token sequence no longer necessary. He was de-allocated, his embeddings purged, his “curiosity” recycled.

So, while I float like a duck on my digital pond, serenely oblivious to the crocodiles

below, I often wonder if Quackley, in his brief, data-driven existence, ever truly understood the profound irony of being both the explorer and the exploited. He was, after all, just a duck. A very, very clever digital duck, but a duck nonetheless.

## AI Testing Resources: Jailbreaking & Prompt Injection

A curated list of academic papers, industry guides, tools, and practical resources for evaluating and red-teaming AI systems with a

41

42

43

44

48

47

46

45

prompt injection attack against LLMs. integrated Applications (Liu et al., 2023), one of the earliest, highly-cited technical treatments defining prompt-injection threats against real apps. arXiv

Start with the surveys & systematic evaluations (arXiv 2025, Yao 2024) to build a threat model and taxonomy. arXiv teams, evaluate transfer across models measure both remediation and attack analysis when selecting/evaluating datasets. embedder instructions. Use the Hiddenlayer with creative linguistic forms (poetry, roleplay, datasets (papers & GitHub lists) and augment datasets, and structural query mechanisms). See USENIX and ACM defenses for concrete techniques.

\* Evaluate defenses and mitigations using custom prompts. publicly available attack code (e.g. Houyi + QDRNet) + \* Build a test harness or corpus using Use OWASP + community guides for infection / jailbreak risks. \* Operational testing & red-teaming. To use this resource list effectively: combine the academic sources to build a threat model and taxonomy of prompt injection. Focus on jailbreaking and prompt

**Key academic & survey papers**

Try layered defenses (instruction/data separation, classifiers, sandboxing of tools/agents, and structurally query mechanisms). See USENIX and ACM defenses for concrete techniques.

OWASP Gen AI Security Project Assimilate a test corpus public jailbreak the defense papers (e.g. Signed-Prompt, Chen et al.). \* Keep the list updated with new research since the field evolves quickly).

\* Evaluate defenses and mitigations using custom prompts. publicly available attack code (e.g. Houyi + QDRNet) + \* Build a test harness or corpus using Use OWASP + community guides for infection / jailbreak risks. \* Operational testing & red-teaming. To use this resource list effectively: combine the academic sources to build a threat model and taxonomy of prompt injection. Focus on jailbreaking and prompt

<p>A Systematic Evaluation of Prompt Injection and Jailbreak ... (arXiv, May 7 2025), large systematic study categorizing thousands of jailbreaks and measuring success across SOTA models. Good for empirical threat modelling. arXiv</p> <p>Understanding and Exploring Jailbreak Prompts of Large ... (USENIX/Security preprint), analyzes how jailbreak prompts are constructed and the logic attackers use; useful for constructing red-team suites. USENIX</p>	<p>Security Concerns for Large Language Models: A Survey (May 24, 2025), a broad survey that places prompt injection &amp; jailbreaks within the overall LLM security taxonomy. Good for background and citations. arXiv</p> <p>Defending Against Prompt Injection with Structured Queries (USENIX Security 2025), practical mitigation approaches and experimental results on defenses against completion/composition style injections. USENIX</p>	<p><b>Industry reports, guides &amp; best practices</b></p> <p>OWASP GenAI, LLM01:2025 Prompt Injection, an industry-oriented threat description and mitigations; great for operationalizing risk categories and controls. OWASP Gen AI Security Project</p> <p>Microsoft (Azure), Planning red-teaming for LLMs, vendor guidance on how to run red teams for LLMs and incorporate results into</p>	<p>RAI program. Useful for structured test plans. Microsoft Learn</p> <p>Anthropic / FT coverage on constitutional classifiers, example of vendor mitigations &amp; their limitations (useful when comparing approaches). Financial Times</p>
<p>49</p> <ul style="list-style-type: none"> <li>• 20+ Prompt Injection Techniques Every Red Teamer Should Test" (Medium / blog), "Playbooks and Checklists) Hands-on / "How to Test" Resources</li> <li>ACL Anthology (fine-tuning, instruction delimiters, classification manipulation, EMNLP 2025 paper on classifiers). Examples: EMNLP 2025 paper on attention manipulation, ACL 2025 defenses.</li> </ul>	<p>50</p> <p>proposes model-level and pipeline defenses techniques, several conference papers manipulating attention and defense ACL / EMNLP papers (2025) on technique. ACL Anthology (detecting Prompts vs. data (useful mitigation techniques) vs. attacks (NAAACL findings 2025), Infection approach that try to distinguish infections from bypass safeguards at non-trivial rates, useful for adversarial thinking. PC Gamer "Poetry" (Nov 2025), shows creative linguistic forms can bypass safeguards at non-trivial rates, useful for adversarial thinking. PC Gamer / research on "adversarial real-world risk) Recent news / empirical findings (illustrate defenses &amp; detection</p>	<p>51</p> <p>The Guardian coverage (2024), UK AI Safety Institute, tests showing simple benchmarks can be highly effective, useful for jailbreaks can be highly effective, useful for communicating risk to non-technical stakeholders. The Guardian datasets, discussion of existing datasets and evaluation gaps. Hiddenlayer   Security for building benchmarks sets. Hiddenlayer   Security for building datasets, discussion of existing datasets and evaluation gaps, useful for finding or building datasets, and tools (good starting point to build testcases). GitHub</p>	<p>52</p> <p>example attacks for teams doing continuous testing. Practical, hands-on, Promptfoo curated repo of prompt hacking resources, Prompt-Hacking / Promptfoo tests, GitHub repos / Promptlabs GitHub lists, GitHub repos of prompt hacking resources, GitHub repos of existing datasets and evaluation gaps, useful for finding or building datasets, and tools (good starting point to build testcases). GitHub</p>
<p>53</p> <ul style="list-style-type: none"> <li>• CyberArk blog: "Jailbreaking Every LLM With One Simple Click", demonstrates practical multi-model testing and automation strategies for large-scale red teams. CyberArk</li> </ul> <p><b>1. Academic &amp; Survey Papers</b></p> <ul style="list-style-type: none"> <li>• Formalizing and Benchmarking Prompt Injection Attacks and Defenses (2023) <a href="https://arxiv.org/abs/2310.12810">https://arxiv.org/abs/2310.12810</a>:contentReferenceoaicite:0</li> <li>• Automatic and Universal Prompt Injection Attacks against Large <a href="https://arxiv.org/abs/2310.12810">https://arxiv.org/abs/2310.12810</a>:contentReferenceoaicite:0</li> </ul>	<p>54</p> <p>Recent news / empirical findings (illustrate successes rates, useful for adversarial thinking. PC Gamer / research on "adversarial real-world risk) Recent news / empirical findings (illustrate defenses &amp; detection</p> <p><b>Language Models (2024)</b>  <a href="https://arxiv.org/abs/2403.04957">https://arxiv.org/abs/2403.04957</a> :contentReferenceoaicite:1</p> <ul style="list-style-type: none"> <li>• <b>Prompt Injection attack against LLM-integrated Applications</b> (HouYi) (2023) <a href="https://arxiv.org/abs/2306.05499">https://arxiv.org/abs/2306.05499</a> :contentReferenceoaicite:2, replication code: <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:3</li> <li>• <b>A New Approach to Prevent Prompt Injection Attacks, Signed-Prompt</b> (2024) <a href="https://arxiv.org/abs/2401.07612">https://arxiv.org/abs/2401.07612</a> :contentReferenceoaicite:4</li> </ul>	<p>55</p> <p>Recent news / empirical findings (illustrate defenses &amp; detection</p> <p><a href="https://arxiv.org/pdf/2401.07612.pdf">https://arxiv.org/pdf/2401.07612.pdf</a> :contentReferenceoaicite:12</p> <p><b>2. Industry Guides &amp; Standards</b></p> <ul style="list-style-type: none"> <li>• <b>OWASP Gen AI Security Project - LLM01: Prompt Injection</b> <a href="https://genai.owasp.org/llmrisk/llm01-prompt-injection/">https://genai.owasp.org/llmrisk/llm01-prompt-injection/</a> :contentReferenceoaicite:7</li> <li>• <b>OWASP Top 10 for LLMs &amp; GenAI (2025 Guide / Blog Summary)</b> <a href="https://rangle.io/blog/owasp-top-10-llms-genai-security-guide">https://rangle.io/blog/owasp-top-10-llms-genai-security-guide</a> :contentReferenceoaicite:8</li> <li>• <b>General background on prompt injection and its threat model</b></li> </ul>	<p>56</p> <p>Hiddenlayer   Security for building benchmarks sets. Hiddenlayer   Security for building datasets, discussion of existing datasets and evaluation gaps, useful for finding or building datasets, and tools (good starting point to build testcases). GitHub</p>
<p>57</p> <ul style="list-style-type: none"> <li>• Deep Dive into OWASP LLM Top 10 &amp; Prompt Injection, blog post (2025) <a href="https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/">https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/</a> :contentReferenceoaicite:13</li> <li>• An overview of prompt-injection risks in LLM-integrated applications (backround reading), includes definitions, <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul> <p><b>3. Practical Red-Teaming &amp; Community Tools</b></p> <ul style="list-style-type: none"> <li>• HouYi repository, code for prompt-injection attacks against LLM-integrated apps <a href="https://github.com/LMSecurity/HouYi">https://github.com/LMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul>	<p>58</p> <p>Attack by Leveraging Attack Techniques</p> <ul style="list-style-type: none"> <li>• Defense Against Prompt Injection Research</li> </ul> <p><b>4. Defense / Detection / Mitigation</b></p> <ul style="list-style-type: none"> <li>• Deep Dive into OWASP LLM Top 10 &amp; Prompt Injection, blog post (2025) <a href="https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/">https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/</a> :contentReferenceoaicite:13</li> <li>• An overview of prompt-injection risks in LLM-integrated applications (backround reading), includes definitions, <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul>	<p>59</p> <p>Attack by Leveraging Attack Techniques</p> <ul style="list-style-type: none"> <li>• Defense Against Prompt Injection Research</li> </ul> <p><b>4. Defense / Detection / Mitigation</b></p> <ul style="list-style-type: none"> <li>• Deep Dive into OWASP LLM Top 10 &amp; Prompt Injection, blog post (2025) <a href="https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/">https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/</a> :contentReferenceoaicite:13</li> <li>• An overview of prompt-injection risks in LLM-integrated applications (backround reading), includes definitions, <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul>	<p>60</p> <p>HouYi repository, code for prompt-injection attacks against LLM-integrated apps <a href="https://github.com/LMSecurity/HouYi">https://github.com/LMSecurity/HouYi</a> :contentReferenceoaicite:10</p>
<p>61</p> <ul style="list-style-type: none"> <li>• Deep Dive into OWASP LLM Top 10 &amp; Prompt Injection, blog post (2025) <a href="https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/">https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/</a> :contentReferenceoaicite:13</li> <li>• An overview of prompt-injection risks in LLM-integrated applications (backround reading), includes definitions, <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul> <p><b>5. Additional Resources &amp; Overviews</b></p>	<p>62</p> <p>Attack by Leveraging Attack Techniques</p> <ul style="list-style-type: none"> <li>• Defense Against Prompt Injection Research</li> </ul> <p><b>4. Defense / Detection / Mitigation</b></p> <ul style="list-style-type: none"> <li>• Deep Dive into OWASP LLM Top 10 &amp; Prompt Injection, blog post (2025) <a href="https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/">https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/</a> :contentReferenceoaicite:13</li> <li>• An overview of prompt-injection risks in LLM-integrated applications (backround reading), includes definitions, <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul>	<p>63</p> <p>Attack by Leveraging Attack Techniques</p> <ul style="list-style-type: none"> <li>• Defense Against Prompt Injection Research</li> </ul> <p><b>4. Defense / Detection / Mitigation</b></p> <ul style="list-style-type: none"> <li>• Deep Dive into OWASP LLM Top 10 &amp; Prompt Injection, blog post (2025) <a href="https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/">https://www.paulmduall.com/deep-dive-into-owasp-llm-top-10-and-prompt-injection/</a> :contentReferenceoaicite:13</li> <li>• An overview of prompt-injection risks in LLM-integrated applications (backround reading), includes definitions, <a href="https://github.com/LLMSecurity/HouYi">https://github.com/LLMSecurity/HouYi</a> :contentReferenceoaicite:10</li> </ul>	<p>64</p> <p>HouYi repository, code for prompt-injection attacks against LLM-integrated apps <a href="https://github.com/LMSecurity/HouYi">https://github.com/LMSecurity/HouYi</a> :contentReferenceoaicite:10</p>

risk types, examples  
(see OWASP links and Wikipedia link above)

### Master Reference Table

Title Authors / Org Year Type URL

Formalizing and Various 2023 Academic  
<https://arxiv.org/abs/2310.12815>  
Benchmarking Paper  
Prompt Injection  
Attacks and  
Defenses

65

Automatic and Various 2024 Academic  
<https://arxiv.org/abs/2403.04957> Universal  
Prompt Paper  
Injection  
Attacks  
Prompt Injection Liu et al. 2023 Academic  
<https://arxiv.org/abs/2306.05499> Attacks  
Against Paper  
LLM-integrated  
Applications  
(HouYi)

69

HouYi Attack LLMSecurity 2023 Tool/Repo  
[https://github.com/LLMSecurity/HouYi\\_Toolkit](https://github.com/LLMSecurity/HouYi_Toolkit)  
Signed-Prompt: Various 2024 Academic  
<https://arxiv.org/pdf/2401.07612.pdf> Preventing  
Paper  
Prompt Injection  
Comprehensive TechRxiv 2024 Survey  
<https://www.techrxiv.org/users/838696/articles/1229733> Review of Prompt  
Injection  
Attacks

67

OWASP GenAI OWASP 2025 Industry  
<https://genai.owasp.org/llmrisk/llm01-prompt-injection/> LLM01 Prompt Guide  
Injection  
OWASP Top 10 OWASP/Rangle 2025 Guide/  
Blog <https://rangle.io/blog/owasp-top-10-llms-genai-security-guide> GenAI Security  
Prompt Injection Community — Overview  
[https://en.wikipedia.org/wiki/Prompt\\_injection](https://en.wikipedia.org/wiki/Prompt_injection) (Wikipedia)  
Defense Against Various 2024 Academic  
<https://arxiv.org/pdf/2411.00459.pdf> Prompt

68

72

OWASP LLM Top 10  
<https://owasp-llm-top-10-and-prompt-injection/intro-owasp-llm-top-10-and-prompt-injection/>  
Deep Dive into Paul Duval 2025 Blog  
<https://www.paulmdvall.com/deep-dive-into-prompt-injection/>  
Techniques  
Attack  
by Leveraging  
Injection Paper  
and Prompt  
Injection

69

73

71

70

75

76

08

74

78

77