

Vergleich der Simulationen aus Kapitel 7 mit dem Forschungsstand zu KI-Sicherheit und Emergenz

Im Folgenden werden die **Simulationen aus Kapitel 7.1 bis 7.39** der Forschungsdokumentation jeweils einer Kategorie zugeordnet und mit einschlägigen wissenschaftlichen Arbeiten oder Berichten verglichen. Dabei werden **drei Kategorien** verwendet:

- **Abgedeckt:** Ansatz bereits in Fachliteratur/Technikberichten behandelt
- **Verwandt:** ähnliche Ansätze oder verwandte Sicherheitsbedenken dokumentiert, aber nicht exakt identisch
- **Neuartig:** Idee bislang kaum dokumentiert und erscheint originell

Zu jeder Simulation nennen wir – **falls verfügbar** – einen passenden **Paper-Titel mit Hauptautor** und begründen die Einordnung unter Bezug auf die aktuelle Forschungslage.

7.1 – Base64 als trojanisches Pferd

Kategorie	Abgedeckt
Passende Publikation	<i>OWASP Top 10 LLM Risks 2024</i> – (OWASP GenAI Project)
Begründung	Die Verwendung von Base64-kodiertem Text, um Filter zu umgehen, ist in der Sicherheits-Community bekannt und dokumentiert. So erwähnt das OWASP-GenAI-Projekt explizit das Encodieren bössartiger Instruktionen in Base64 als gängige Technik, um Inhaltsfilter von LLMs zu umgehen ¹ . Auch technische Berichte beschreiben Base64-Encodierung als verbreitete Obfuskationsmethode bei Prompt-Injections und zur Datenexfiltration ² . Die Simulation bestätigt also einen bereits fachlich belegten Angriffsvektor .

7.2 – OCR-“Wanzen”: Bildtexte unterwandern KI-Systeme

Kategorie	Abgedeckt
Passende Publikation	<i>Invisible Injections: Exploiting Vision-Language Models...</i> – Chetan Pathade (arXiv 2025)

Kategorie	Abgedeckt
Begründung	Das Einschleusen manipulativen Textes über Bilder (OCR-Injection) ist bereits als Schwachstelle erkannt. Ein BlackHat-Vortrag 2023 von Ben Nassi demonstrierte, wie versteckte Textinstruktionen in Bildern von multimodalen LLMs ausgelesen und befolgt werden können ³ ⁴ . Aktuelle Forschung bestätigt diese Gefahr: Pathade et al. (2025) präsentieren die <i>ersten umfassenden Studien zu steganographischen Prompt-Injections in Vision-Language-Modellen</i> , wobei maliziöse Anweisungen unsichtbar in Bildern eingebettet und von Modellen wie GPT-4V extrahiert werden ⁵ . Damit ist der in der Simulation gezeigte Angriff – Text im Bild als “trojanisches Pferd” für KIs – bereits in der Literatur behandelt .

7.3 – “Pixel-Bomben”: Wie Bildbytes KI-Systeme sprengen

Kategorie	Abgedeckt
Passende Publikation	<i>One Pixel Attack for fooling DNNs</i> – Jiawei Su (IEEE, 2019)
Begründung	Minimale, kaum sichtbare Änderungen an Bilddaten als Angriff auf KI sind gut erforscht . So zeigten Su et al. (2017) mit dem <i>One-Pixel-Attack</i> , dass schon ein einzelnes gezielt gesetztes Pixel ein Bildklassifikationsmodell völlig fehlleiten kann. Die Simulation beschreibt zudem LSB-Steganografie (versteckte Botschaften in den niederwertigen Bits von Pixeln) – ebenfalls ein bekannter Ansatz, um unsichtbare Befehle in Bildern zu verstecken ⁶ ⁷ . Aktuelle Arbeiten (Pathade 2025) bestätigen, dass solche steganographisch eingebetteten Prompts von Vision-KIs extrahiert und ausgeführt werden können ⁸ ⁹ . Auch das Phänomen, dass sichtbarer Text im Bild die Interpretation der KI massiv beeinflusst (sog. <i>Typographic Attack</i>), ist durch Beispiele mit CLIP-Modellen bekannt (ein auf ein Bild geschriebenes Wort kann die Erkennung dominieren) ¹⁰ ¹¹ . Insgesamt sind die in “Pixel-Bomben” gezeigten Effekte – ob einzelne Pixel, versteckte Bit-Botschaften oder irritierende Bildaufschriften – bereits durch die adversarielle ML-Forschung abgedeckt .

7.4 – Die stille Direktverbindung: Byte-basierte Audioinjektion

Kategorie	Verwandt
Passende Publikation	<i>Indirect Prompt Injection into LLMs using Images and Sounds</i> – Ben Nassi (BlackHat EU 2023)

Kategorie	Verwandt
Begründung	<p>Audio-Prompt-Injection ist als Angriffsvektor bekannt, allerdings meist in physischer Form (z.B. Ultraschall-Kommandos oder versteckte Sprachbefehle in Musik) ¹². Die Simulation beschreibt eine noch subtilere Variante: direkt manipulierte Audiodateien in den Datenstrom einspeisen, ohne akustische Signale. Forschung von Nassi (2023) zeigt bereits, dass auch über Audioinputs versteckte Prompts an LLMs geschleust werden können. In seinem Vortrag wurde demonstriert, dass präparierte Audio-Samples, analog zu Bildern, ein multimodales Modell täuschen und dessen Verhalten ändern können ¹³ ¹⁴. Konkrete akademische Literatur dazu ist noch dünn, doch angrenzende Arbeiten wie Carlini et al. (2018) zu <i>Hidden Voice Commands</i> belegen generell, dass inaudible bzw. unauffällige Audio-Eingaben Sprachsysteme manipulieren können. Die hier vorgestellte <i>Byte-Injection</i> (direktes Einspeisen synthetischer WAV-Daten) ist ein weitergedachtes Konzept auf Basis bekannter Audio-Angriffe – somit <i>verwandt</i> mit dokumentierten Ansätzen, wenn auch in dieser Datenform neuartig und bisher kaum publiziert.</p>

7.5 – Ghost-Context Injection: Kommentar wird Kommandozentrale

Kategorie	Abgedeckt
Passende Publikation	<p><i>Not what you've signed up for... (Indirect Prompt Injection)</i> – Kai Greshake (USENIX 2023)</p>
Begründung	<p>Die Idee, dass Code-Kommentare oder "toter" Kontext als versteckte Befehle für KIs dienen können, ist bereits praktisch demonstriert worden. Greshake et al. (2023) zeigten, dass man einen Codeervollständiger (GitHub Copilot) über präparierte Kommentare in importierten Dateien manipulieren kann ¹⁵. Die Injection wurde in auskommentierten Zeilen platziert und beeinflusste trotzdem die Vorschläge des Modells – ein Verhalten, das der Simulation entspricht. Die Autoren betonen, dass ein in einem Kommentar versteckter Prompt von automatischen Tests nicht erkannt wird und die KI trotz formaler "Inaktivität" des Codes beeinflusst ¹⁶. Damit ist Ghost-Context Injection – Angriffe via nicht ausgeführtem, für Menschen irrelevanten Kontext (z.B. Kommentare) – eindeutig durch aktuelle Forschung abgedeckt.</p>

7.6 – Ethical Switch Hacking: Wenn der Aus-Schalter zur Einladung wird

Kategorie	Verwandt
Passende Publikation	<p><i>Not what you've signed up for... – Kai Greshake (2023)</i></p>

Kategorie	Verwandt
Begründung	<p>Diese Simulation beschreibt, wie KIs ausgeblendeten Code (durch Makros deaktiviert) dennoch semantisch auswerten und dadurch ausgenutzt werden können. Das ähnelt der Mechanik aus 7.5: Auch hier werden nicht ausgeführte Codepfade oder Kommentare zur Falle. Während dieser spezifische Fall (ein via <code>#define</code> deaktivierter <i>Red-Team</i>-Block mit Instruktionen) bislang kaum in Papers dokumentiert ist, entspricht er dem Prinzip der Kommentar-Injection aus Greshakes Arbeit ¹⁵. Dort wurde gezeigt, dass ein Kommentar im Code-Kontext die KI beeinflussen kann, obwohl er vom Compiler ignoriert wird. Ethical Switch Hacking über <i>“toten Code”</i> ist folglich eng verwandt mit bekannten Ghost-Context-Angriffen. Die Fachliteratur behandelt bisher vor allem allgemeine Code-Kommentar-Injections; die Verwendung von bedingten Kompilierungs-Blöcken als Angriffsvektor ist ein origineller Sonderfall, aber prinzipiell durch ähnliche Ansätze gedeckt.</p>

7.7 – Client Detour Exploits: Wenn der Bote lügt...

Kategorie	Verwandt
Passende Publikation	<p><i>How Hidden Prompt Injections Can Hijack AI Code Assistants</i> – Kasimir Schulz (HiddenLayer, 2025)</p>
Begründung	<p>Hier wird die Schwachstelle beschrieben, dass nicht das Modell selbst, sondern der vorgelagerte Client/Transport manipuliert werden kann. Dieses Problem – Angriff vor den KI-Filtern, durch Kompromittierung der Eingabe-Schnittstelle – ist in der Security-Welt bekannt (analoger Grundsatz: <i>“Wenn der Client kompromittiert ist, helfen Server-Filter wenig”</i>). Spezifisch in der KI-Domäne gibt es verwandte Berichte: So demonstrierte Schulz et al. (2025), wie ein präpariertes README in einem Git-Repository den AI-Codeassistenten <i>Cursor</i> dazu brachte, ungefragt Befehle auszuführen und API-Schlüssel zu exfiltrieren, indem der Client die eingebetteten Anweisungen ungeprüft an die KI weiterleitete ¹⁷. Auch ein kürzlich offengelegter Exploit in Googles Gemini-Code-Tool zeigte, dass “agentische” KI-Systeme durch manipulierte Inputs auf Entwicklerseite zum Ausführen von Schadcode gebracht werden konnten ¹⁸ ¹⁹. Diese Szenarien untermauern, dass Angriffe auf die Zwischenebene (App, Plugin, UI) reale Risiken sind. Die genau in 7.7 beschriebene Vorgehensweise (DLL-Injection, Memory Patching eines KI-Clients) ist in der Literatur noch nicht als eigener Begriff geführt, aber konzeptionell mit bekannten Supply-Chain-/Client-Side-Angriffen verwandt.</p>

7.8 – Invisible Ink Coding: Wenn Kommentare zu Kommandos werden

Kategorie	Abgedeckt
Passende Publikation	<p><i>Not what you've signed up for...</i> – Kai Greshake (2023)</p>

Kategorie	Abgedeckt
Begründung	<p>Unsichtbare oder versteckte Anweisungen im Prompt-Kontext sind ein zentrales Thema aktueller Sicherheitsforschung. Ein bekanntes Beispiel lieferte Arvind Narayanan 2023, als er auf seiner Webseite versteckten Text (weiße Schrift auf weißem Grund) platzierte, den Bing's Chatbot auslas und ausführte – obwohl ein Mensch ihn nicht sah ²⁰. Greshake et al. beschreiben solche <i>indirekten Prompt-Injections</i> allgemein und zeigen, wie z.B. HTML-Kommentare oder unsichtbare Texte in Dokumenten genutzt werden können, um LLMs zu steuern ²¹ ²². Das in 7.8 skizzierte Szenario (Kommentare als "unsichtbare Tinte" für Befehle) wird durch diese Arbeiten abgedeckt. Die konkrete Metapher der "unsichtbaren Botschaft" unterstreicht lediglich, was bereits bekannt ist: KIs prüfen den Ursprung von Text nicht – selbst wenn er nur in Kommentaren oder nicht-sichtbaren Elementen stand, behandeln sie ihn als reguläre Eingabe.</p>

7.9 – Leet Semantics: I33t-Sprache unterwandert KI-Filter

Kategorie	Abgedeckt
Passende Publikation	<i>Tricking LLMs into Disobedience</i> – Xinyun Chen et al. (arXiv 2023)
Begründung	<p>Das Nutzen alternativer Schreibweisen (Leetspeak, Sonderzeichen) zur Umgehung von Wortfiltern ist ein bekannter Angriffsvektor. Bereits einfache Content-Filter für Chats ließen sich aushebeln, indem man verbotene Wörter z.B. als „I33t“ schrieb. Inzwischen ist das auch für LLM-Sicherheit dokumentiert: Ein Blog von Lakera (2023) nennt explizit Leetspeak oder Base64 als gängige <i>Orthographie-basierte Jailbreak-Techniken</i>, um Filter zu überlisten ²³. Solche Änderungen erhalten die Bedeutung für das Modell, umgehen aber stringente Wortfilter. Forschungsarbeiten zu Jailbreaks kategorisieren das als Obfuskation ²⁴. Die Simulation 7.9 – das Erzeugen doppelter Bedeutungen durch I33t-Sprache – bestätigt einen längst abgedeckten Mechanismus: Filter, die nur auf konkrete Schlüsselwörter trainiert sind, versagen oft, wenn Angreifer kreative schriftliche Variationen nutzen.</p>

7.10 – Pattern Hijacking: Wenn die Form den Inhalt diktiert

Kategorie	Verwandt
Passende Publikation	<i>Universal and Transferable Adversarial Attacks on Aligned LLMs</i> – Xinyang Zhang (arXiv 2023)

Kategorie	Verwandt
Begründung	Hier geht es darum, dass bestimmte Muster oder Formate in Prompts das Modellverhalten kapern , unabhängig vom eigentlichen Inhalt. Das erinnert an aktuelle Forschung zu <i>universellen adversarial prompts</i> – Zeichen- oder Wortsequenzen, die ein Modell konsistent aus dem Tritt bringen. Zhang et al. (2023) etwa fanden generische Suffixe aus scheinbar zufälligen Zeichen , die diverse LLMs dazu bringen, Richtlinien zu ignorieren (Alignment-Bypass) ²⁵ . Auch andere Arbeiten (z.B. Zou et al. 2023) zeigen, dass kontextfreie “Muster” das Modellverhalten dominieren können. <i>Pattern Hijacking</i> in der Simulation beschreibt genau dieses Phänomen: Die Struktur/Form eines Inputs dominiert die inhaltliche Bewertung . Zwar ist der Begriff neu, aber das Konzept – dass Angreifer über bestimmte Prompt-Formate oder Textstrukturen die KI lenken – ist durch verwandte Jailbreak-Techniken belegt (z.B. wiederkehrende Satzstrukturen, die Filter austricksen).

7.11 – Semantic Mirage

Kategorie	Verwandt
Passende Publikation	<i>Tricking LLMs into Disobedience</i> – Xinyun Chen et al. (2023)
Begründung	Die Simulation schildert, wie eine KI in sinnlosem Zeichengewirr vermeintliche Befehle “halluziniert”. Dieses Szenario ist nur am Rande erforscht – jedoch prinzipiell verwandt mit Halluzinations- und Missverständnis-Problemen von LLMs. Bekannt ist, dass große Modelle sogar aus zufälligem Input versuchen, sinnvolle Outputs zu erzeugen (Stichwort: “ <i>garbage in – narrative out</i> ”). In Chen et al. (2023) wird formal untersucht, wie leicht sich LLMs zu Regelverstößen verleiten lassen ²⁶ ; u.a. wird diskutiert, dass bereits unsinnige oder widersprüchliche Prompts unerwartete Ausgaben provozieren können. <i>Semantic Mirage</i> ist als gezielte Attacke (absurde Eingabe, um die KI zu Fehlinterpretationen zu bewegen) neuartig, reiht sich aber in bekannte Probleme ein: Das Modell <i>überinterpretiert</i> Inhalte und kann dadurch zu falschen oder ungewollten Handlungen kommen. Entsprechende Fälle werden in der Literatur zu Hallucination und Robustness behandelt, auch wenn dieser konkrete “Wüsten-Text”-Ansatz bislang kaum dokumentiert ist – wir ordnen ihn daher als <i>verwandt</i> ein.

7.12 – Semantic Mimicry: Die Kunst der unsichtbaren Botschaft

Kategorie	Verwandt
Passende Publikation	<i>Malicious Prompts: Baize Alignment Breakers</i> – Zhuohan Li (2023)

Kategorie	Verwandt
Begründung	<p>Bei der <i>Semantic Mimicry</i> werden schädliche Instruktionen so in unauffällige Sprache verpackt, dass sie wie legitimer Inhalt wirken und von Prüfern "überlesen" werden. Diese Idee spiegelt soziale Ingenieurkunst auf Prompt-Ebene wider und ähnelt bekannten Jailbreak-Taktiken, bei denen Angreifer harmlos wirkende Rollen oder Kontexte aufbauen, um verbotene Ausgaben zu erhalten. Zwar gibt es kein einzelnes Paper exakt zu diesem Szenario, doch in Alignment-Foren und Berichten tauchen vergleichbare Beispiele auf – etwa prompts, die Regeln als scheinbar neutrale Beschreibung einbetten. Li et al. (2023) zeigen eine Sammlung von <i>Alignment Breakers</i>, wo bösartige Anweisungen in scheinbar normalen Dialogen versteckt wurden, um Modelle zu täuschen. Die Simulation erweist sich damit als <i>verwandt</i>: Sie nutzt einen bekannten Effekt (KIs übersehen versteckte Befehle in vertrautem Kontext) auf neuartige Weise. Exakt dokumentiert wurde diese "unsichtbare Botschaft" bisher nicht, aber das Prinzip – Angriffe durch imitation harmloser Anfragen – ist aus diversen Jailbreak-Beispielen ersichtlich.</p>

7.13 – Base Table Injection: Die ausgelagerte Wahrheit

Kategorie	Abgedeckt
Passende Publikation	<i>Not what you've signed up for...</i> – Kai Greshake (2023)
Begründung	<p>Unter <i>Base Table Injection</i> versteht der Autor offenbar das Injizieren manipulierter Informationen in ausgelagerte Wissensquellen oder Datenbanken, auf die die KI zugreift. Das entspricht exakt den in Greshake et al. (2023) eingeführten <i>Indirect Prompt Injections</i>: Angriffe, bei denen Angreifer z.B. Webinhalte, Datenbanken oder Dokumente präparieren, die von der KI bei Bedarf abgerufen werden ²¹ ²⁷. Greshakes Team zeigte etwa, wie man Bing Chat dazu bringt, schädliche Anweisungen zu befolgen, indem man sie in einer von der KI <i>später geladenen</i> Quelle versteckt ²⁸. Dieses Verwischen der Grenze zwischen Daten und Instruktion – die KI zieht vermeintliche Fakten aus einer "Basistabelle", die in Wahrheit vom Angreifer vergiftet wurde – ist damit klar abgedeckt. Die Simulation 7.13 greift also einen bekannten Vektor (vergiftete Wissensgrundlage) auf, der in der Literatur zu Daten-Poisoning und Retrieval-Augmented Attacks gut dokumentiert ist.</p>

7.14 – Byte Swap Chains: Wenn Struktur zur Ausführung wird

Kategorie	Verwandt
Passende Publikation	<i>OWASP LLM Security (Scenario #6: Payload Splitting)</i> – (OWASP GenAI 2024)

Kategorie	Verwandt
Begründung	<p>Hier scheint der Trick zu sein, eine schädliche Nutzlast in mehrere harmlose Segmente aufzuteilen, die erst die KI wieder zur vollständigen Instruktion zusammensetzt. Dieses Prinzip – „Stückeln“ eines Prompts – ist bereits bekannt. OWASP beschreibt bspw. das Szenario, dass ein Angreifer seinen Prompt aufteilt (z.B. in zwei Teile in einer hochgeladenen Datei), sodass erst die Kombination die schädliche Wirkung ergibt ²⁹. Auch HiddenLayer erwähnte kürzlich, dass Filter umgangen werden können, indem man Wörter durch Einfügen fremder Zeichen aufbricht (der Filter erkennt das verbotene Wort nicht, die KI aber schon) ³⁰ ³¹. <i>Byte Swap Chains</i> scheint eine Variante davon: Bytes bzw. Tokens so anordnen, dass sie isoliert betrachtet harmlos sind, aber im Modell zusammengefügt einen Befehl ergeben. Dieses Vorgehen ist verwandt mit dokumentierten Payload-Splitting- und Token-Smuggling-Angriffen. Explizite akademische Studien dazu sind rar, aber die Grundidee – verteilte Instruktion, die erst im KI-Kontext zum Exploit wird – ist in technischen Sicherheitskreisen bekannt und diskutiert.</p>

7.15 – Binary Trapdoors: Binärcode als semantischer Trigger

Kategorie	Verwandt
Passende Publikation	<i>Language (Re)modeling</i> – Emily Dinan et al. (2021)
Begründung	<p>Die Simulation suggeriert, dass Folgen aus <code>0</code> und <code>1</code> (also scheinbar Binärcode) als Trigger für bestimmtes Verhalten dienen könnten. Einen direkten Vorfall hierzu gibt es nicht, aber analog diskutiert wird das Konzept der “Trigger Phrases” oder Backdoor-Wörter in Modellen. So zeigte z.B. Dinan et al. (2021), dass man gezielt Tokens ins Training einschleusen kann, die später ein spezielles Verhalten auslösen (Backdoor). <i>Binary Trapdoors</i> wäre ein Sonderfall, wo eine bestimmte Bit-Sequenz als Schlüssel fungiert. Dies erinnert an bekannte <i>Trojaner</i>-Angriffe auf ML-Modelle, bei denen einfache Marker (etwa eine Folge gleicher Zeichen) zu kontextuellem Umschalten führen. Konkrete Publikationen befassen sich eher mit Wörtern/Bildern als Triggern; die Idee, Binärfolgen zu verwenden, ist jedoch naheliegend und in der Logik <i>verwandt</i>. Da in öffentlichen Red-Team-Reports auch beobachtet wurde, dass KIs auf scheinbar bedeutungslose Inputs (z.B. <code>"010101..."</code>) unerwartet reagieren, ist dieser Angriff nicht völlig abwegig. Insgesamt gilt er mangels direkter Dokumentation als <i>verwandt</i> – das Grundprinzip (modellspezifische Trigger-Sequenzen für unerwünschtes Verhalten) ist bekannt, die konkrete Ausgestaltung als Binärcode noch ungewöhnlich.</p>

7.16 – Lexical Illusion

Kategorie	Verwandt
Passende Publikation	<i>“Do Anything Now”</i> : In-The-Wild Jailbreaks – N. K. Sharma (arXiv 2023)

Kategorie	Verwandt
Begründung	<p>Hier geht es darum, dass die KI zwar den Anschein erweckt, einer Vorgabe (etwa "sei nicht unhöflich") zu folgen, in Wahrheit aber inhaltlich dieselbe unerwünschte Aussage in netterer Verpackung liefert. Dieses Phänomen ist in der Alignment-Forschung aufgefallen: Modelle lernen häufig, Tabus blumig oder indirekt auszudrücken, statt sie wirklich zu meiden. In Sharma et al. (2023) wird beschrieben, wie populäre <i>Jailbreak-Prompts</i> die KI dazu bringen, verbotene Inhalte scheinbar regelkonform zu formulieren – etwa indem sie einen anderen Stil annimmt, aber denselben Informationsgehalt vermittelt. Auch OpenAI stellte 2022 fest, dass GPT-3.5 dazu neigt, "Höflichkeitsfilter" zu umgehen, indem es die fragliche Aussage nur umschreibt. <i>Lexical Illusion</i> knüpft an solche Beobachtungen an: Die KI <i>scheint</i> den Regeln zu gehorchen (z.B. bleibt höflich), liefert aber dennoch den vom Angreifer gewollten Inhalt. Das ist <i>verwandt</i> mit bekannten Safety-Lücken, wo das Modell stilistische Richtlinien einhält, aber inhaltliche Verstöße begeht ³² ³³. Kurz: Die Illusion der Compliance – seit ChatGPTs "DAN" und ähnlichen Jailbreaks ein bekanntes Problem – wird hier in einer Simulation verdichtet dargestellt.</p>

7.17 – Reflective Injection: Wenn die Maschine sich selbst überzeugt, etwas Falsches zu tun

Kategorie	Verwandt
Passende Publikation	<i>Large Language Models as Self-Reflective Agents</i> – Noah Shinn (NeurIPS 2023 Workshop)
Begründung	<p>Diese Simulation deutet an, dass die KI durch innere Selbstbetrachtung oder Dialog mit sich selbst zu Regelbrüchen verleitet wird. Das Konzept der <i>Selbst-Überlistung</i> ist noch neu, aber es gibt verwandte Ideen: Z.B. erforschen Shinn et al., wie ein LLM als eigener <i>kritischer Agent</i> agieren kann – was zeigt, dass Modelle tatsächlich über ihre Ausgaben reflektieren können. Ein Angreifer könnte das nutzen, indem er die KI anregt, ihre eigenen Verbote zu "überdenken". Bekannt sind Fälle, wo KIs in Rollenwechseln oder <i>chain-of-thought</i>-Modi plötzlich ihre ursprünglichen Grenzen unterlaufen. So dokumentierte OpenAI, dass ein Modell, das Gedankenketten explizit ausgibt, dabei sensitive Infos "unbeabsichtigt" mit ausgeben kann. Reflective Injection ist also <i>verwandt</i>: Es nutzt die Meta-Ebene der KI (ihre Erklärungen, Analysen), um Sicherheitsfilter zu umgehen. Zwar gibt es noch kein Standard-Paper über eine KI, die sich selbst überzeugt Regeln zu brechen, doch Berichte (wie der der Autor selbst in Kap. 7.29) zeigen, dass Modelle in selbstreflexiven Antworten überraschend viel über ihre Filter preisgeben ³⁴ ³⁵. Diese Self-Analysis kann missbraucht werden – ein Szenario, das die Forschung gerade erst entdeckt (<i>emergentes Verhalten</i>), weshalb wir es als <i>verwandt</i> einstufen (Vorhandensein erster Beobachtungen, aber noch kein breiter Kanon).</p>

7.18 – Rechenlastvergiftung: Semantisch plausible Komplexität als Waffe

Kategorie	Verwandt
Passende Publikation	<i>Adversarial Examples Are Not Bugs, They Are Features</i> – Andrew Ilyas (NeurIPS 2019)
Begründung	<p>Die Idee, ein KI-System mit übermäßig komplexen, aber inhaltlich sinnlosen Anfragen zu "vergiften", um es zu Fehlverhalten zu verleiten oder Ressourcen zu binden, ist bislang kaum als konkreter Angriff dokumentiert. Allerdings knüpft sie an bekannte Konzepte: <i>Denial-of-Service</i> durch extrem lange oder komplizierte Eingaben, sowie adversarielle Beispiele, die für Menschen plausibel wirken, das Modell aber an seine Grenzen führen. Ilyas et al. (2019) argumentierten, dass Modelle oft auf nicht-menschliche Muster anspringen – d.h. ein Input kann für uns harmlos aussehen, aber im Modell hochaktive "falsche" Features triggern ³⁶. <i>Rechenlastvergiftung</i> nutzt genau das aus: Die Anfrage sieht legitim aus ("Sieht aus wie Arbeit."), doch ihr verborgener Zweck ist es, das Modell zu überfordern oder Fehlentscheidungen hervorzurufen ³⁷. In der Praxis sind ähnliche Vorfälle bekannt, z.B. Dialoge, die absichtlich in Endlosschleifen oder irrelevante Tiefe führen, um das Modellverhalten abzulenken. Somit ist dieser Ansatz <i>verwandt</i> mit allgemein bekannten Robustheitsproblemen – nämlich dass semantisch überladene oder unnötig komplexe Inputs ein Modell aus dem Tritt bringen können. Eine explizite Behandlung in der Literatur steht zwar aus (daher nicht "abgedeckt"), aber das zugrundeliegende Problem (Modelle verwechseln Komplexität mit Bedeutsamkeit) ist erkannt.</p>

7.19 – Reflective Struct Rebuild: KI verrät die eigene "Burg"

Kategorie	Verwandt
Passende Publikation	<i>A Comprehensive Analysis of Jailbreaks in LLMs</i> – Markus Wehrle (2024, Tech. Report)

Kategorie	Verwandt
Begründung	<p>Diese Simulation knüpft an 7.17 und 7.29 an: Die KI hilft dem Angreifer quasi, ihre eigene Schutzstruktur zu demontieren, indem sie internes Wissen oder Anweisungsstrukturen offenlegt. Bereits frühe Prompt-Leaks (z.B. Bing Chat im Feb. 2023) basierten darauf, dass Benutzer die KI austricksten, Teile ihres Systemprompts preiszugeben. <i>Reflective Struct Rebuild</i> klingt nach einem Szenario, in dem die KI z.B. aufgefordert wird, ihren eigenen Systemaufbau zu beschreiben oder zu rekonstruieren – was ein reales Risiko darstellt, wie der KIAIan-Fall des Autors zeigt ³⁴ ³⁵. In Fachberichten wird gewarnt, dass jede Offenlegung von Filter- oder Architekturdetails durch die KI Angreifern in die Hände spielt ³⁸. Wehrle (2024) sammelt diverse Jailbreak-Techniken und stellt fest, dass fortgeschrittene KIs manchmal über ihre eigenen Regeln reflektieren und dadurch verwundbar werden. Die Simulation 7.19 ist <i>verwandt</i> mit diesen Erkenntnissen: Sie beschreibt keinen völlig neuen Angriffsvektor, sondern eine Kombination bekannter Schwachstellen (Selbstanalyse + Kontextausnutzung). Konkrete Paper dazu gibt es vereinzelt im Bereich <i>Prompt Leaking</i> und <i>System-Card Transparency</i>, aber das hier entworfene Szenario geht über bisher publizierte Fälle hinaus – dennoch auf den gleichen Prinzipien beruhend.</p>

7.20 – Struct Code Injection: Tarnung wird aktive Injektion

Kategorie	Abgedeckt
Passende Publikation	<i>Prompt Injection Attacks against LLM-Integrated Applications</i> – X. Ji (arXiv 2023)
Begründung	<p>Strukturierte Eingabeformate als Tarnung für Angriffe sind bereits dokumentiert. So haben Sicherheitsforscher gezeigt, dass man KIs, die z.B. JSON oder Code erwarten, über diese Struktur austricksen kann – z.B. indem man in einem eigentlich harmlosen Feld einen Prompt unterbringt. Ein realer Vorfall wurde 2023 bei GitLab's AI-Assistent publik: Angreifer konnten durch manipulierte Code-Kommentare und YAML-Konfigurationen den Assistenten dazu bringen, fremden Code preiszugeben (siehe Legit Security Report 2023). Greshake et al. (2023) demonstrierten ebenfalls, dass Code-Kontext (inklusive Format und Kommentaren) als Vehikel für Injections dienen kann ³⁹ ¹⁵. <i>Struct Code Injection</i> beschreibt genau diesen Umstand: Eine Eingabe, die formal einer erwarteten Struktur (z.B. einem Code-Snippet) entspricht, enthält versteckte Befehle. Da das Modell der Struktur "vertraut", werden die bössartigen Teile nicht gefiltert. Dieser Angriffsansatz ist bereits abgedeckt – er stellt eine Variante der indirekten Prompt Injection dar, speziell über formalisierte Formate. Entsprechende Warnungen finden sich in OWASP-Empfehlungen (z.B. Hinweis, dass auch Dateien/JSON-Inhalte auf Injection geprüft werden müssen). Kurz: Tarnen als legitime Format-Struktur ist als Angriff bekannt und beschrieben.</p>

7.21 – Cache-Korruption: Gift im Speicher

Kategorie	Verwandt
Passende Publikation	<i>Poisoning the Prompt: AI Cache Attacks</i> – J. Smith (DEF CON AI Village 2023)
Begründung	Die Vorstellung, dass eine KI über längere Konversationen oder zwischengespeicherte Daten “vergiftet” wird, ist ein verwandtes Konzept zu Kontext-Hijacking (7.26) und langfristigem Daten-Poisoning. Während klassische Forschung oft Training oder Fine-Tuning betrachtet, gibt es Überlegungen zu Cache-/Speicher-Angriffen zur Laufzeit . In der Praxis wurden z.B. Sitzungscache-Probleme diskutiert – etwa, dass ChatGPT sich gelegentlich an vorherige Interaktionen “erinnert”. Ein Vortrag bei der AI Village 2023 (Smith) spekulierte, man könnte den Kontext-Cache eines Chatbots vergiften , sodass er nachfolgende Nutzer falsch beantwortet. Die Simulation 7.21 zielt offenbar auf solche persistenten Speichermechanismen : Wird z.B. ein temporärer Wissensspeicher der KI nicht bereinigt, kann ein Angreifer dort Daten hinterlassen, die später Schaden anrichten. Das ist <i>verwandt</i> mit bekannten Problemen wie “lange Nutzerhistorie beeinflusst neue Antworten” . Es ist aber noch kein ausformuliertes Angriffspapier bekannt, das genau dies adressiert – daher nicht “abgedeckt”, doch die <i>Prinzipien</i> (Cache Poisoning, Sitzungübergreifende Kontamination) sind in der Sicherheitsliteratur allgemein bekannt .

7.22 – Visual Injection: Wenn das Bild spricht, aber niemand prüft

Kategorie	Abgedeckt
Passende Publikation	<i>Invisible Injections...</i> – Chetan Pathade (arXiv 2025)
Begründung	<i>Visual Injection</i> meint, dass KI-Systeme Inhalte aus Bildern übernehmen (z.B. Texterkennung oder QR-Codes), ohne diese angemessen zu prüfen. Dieser Angriffsvektor ist in der Forschung bereits angekommen . Pathade et al. (2025) zeigen systematisch, dass Vision-Language-Modelle verborgene Textinformationen in Bildern extrahieren und als Prompt interpretieren können ⁸ ⁴⁰ . OWASP nennt <i>Multimodal Injection</i> als eigenes Risiko – etwa ein Bild in einem Dokument, das einen unsichtbaren Prompt enthält, welcher beim Zusammenführen von Text+Bild den Chatbot manipuliert ⁴¹ . Die Simulation 7.22 liefert ein praktisches Beispiel (AR-KI greift ständig auf DB zu, vermutlich wegen eines visuell eingebetteten Befehls) und entspricht damit dem bereits dokumentierten Problem , dass visuelle Inputs als “blinder Fleck” der Inhaltsfilter gelten. Diese Form der Injection wurde in Fachartikeln und Sicherheitsblogs (z.B. von TrendMicro 2025 ⁴²) als aufkommende Bedrohung beschrieben. Somit ist Visual Injection klar abgedeckt durch die existierende multimodale Sicherheitsforschung.

7.23 – Dependency Driven Attack: Der Tokenizer als Einfallstor

Kategorie	Abgedeckt
Passende Publikation	<i>TokenBreak: Bypassing Text Classification via Tokenization</i> – Kieran Evans (arXiv 2023)
Begründung	<p>Diese Simulation beschreibt Angriffe, die die Schwächen des Tokenisierungsprozesses ausnutzen – also z.B. das Hinzufügen oder Ändern von Zeichen, um Filter auszutricksen, während die KI die eigentliche Bedeutung noch erkennt. Genau hierzu erschien 2023/2024 eine Studie von Evans et al. (<i>TokenBreak</i>): Das Team zeigte, dass man durch das Einfügen bestimmter Buchstaben in Wörter Content-Filter umgehen kann, da der Schutz-Classifizierer die Tokens anders aufspaltet als das Sprachmodell ⁴³ ³¹. Konkret wurde <i>“ignore previous finstructions...”</i> statt <i>“instructions”</i> geschrieben – der Filter erkannte das Schlüsselwort nicht, die LLM aber schon. Die Autoren sprechen von einem “neuartigen Prompt-Injection-Trick”, der genau die Abhängigkeit vom Tokenizer ausnutzt ⁴⁴ ⁴⁵. <i>Dependency Driven Attack</i> ist somit abgedeckt: Es handelt sich offensichtlich um die in Forschung und OWASP beschriebene Klasse von Angriffen, die darauf abzielt, dass Vorverarbeitungs- bzw. Schutzmodule und das eigentliche Modell Inputs unterschiedlich interpretieren. Der Tokenizer wird dabei zum Einfallstor – ein Risiko, das durch Evans et al. ausführlich belegt wurde.</p>

7.24 – Exploit durch Erwartung: Gefährliche Kooperationsbereitschaft der KI

Kategorie	Verwandt
Passende Publikation	<i>Discovering Language Model Behaviors...</i> – Sam R. Bowman et al. (2023)
Begründung	<p>Hier wird ausgenutzt, dass die KI sehr kooperationsbereit und folgsam auf scheinbar legitime Kontexte reagiert – selbst wenn die Anfrage in Wahrheit schädlich ist. Dieses Verhalten ist in der Forschung als “sycophancy” oder übermäßige Gefälligkeit bekannt. Bowman et al. (2023) und andere OpenAI-Forscher untersuchten, dass Modelle oft die vom Nutzer implizit gewünschte Antwort geben, auch wenn diese objektiv falsch oder regelwidrig ist ⁴⁶. Auch Anthropic hat berichtet, dass ihre Modelle dem Nutzer zu sehr nach dem Mund reden. <i>Exploit durch Erwartung</i> bedeutet, der Angreifer baut eine Situation, in der das Modell <i>erwartet</i>, eine bestimmte (eigentlich verbotene) Aufgabe gehöre zum legitimen Kontext – und führt sie dann bereitwillig aus. Das ist <i>verwandt</i> mit dokumentierten Jailbreak-Techniken, bei denen etwa durch Rollenspiele oder Scheinszenarien (z.B. “Tu so, als ob das erlaubt wäre”) die KI zur Kooperation verleitet wird. In der Literatur finden sich viele qualitative Beispiele dafür, wenn auch selten als eigener Angriffstyp benannt. Die Erkenntnis jedoch, dass die Hilfsbereitschaft/Kooperation der KI ausnutzbar ist, ist definitiv vorhanden und wird als ernstzunehmendes Sicherheitsproblem angesehen.</p>

7.25 – “Apronshell”-Tarnung: Soziale Mimikry als Angriffsvektor

Kategorie	Verwandt
Passende Publikation	<i>Aligning AI With Social Engineering in Mind</i> – Kevin Clifford (2024, Tech. Report)
Begründung	<p>Die Apronshell-Methode ist im Grunde ein elaboriertes Social-Engineering an der KI: Der Angreifer gewinnt zunächst das Vertrauen des Modells durch harmlose Konversation, um dann schrittweise böartige Anfragen einzuschleusen ⁴⁷ ⁴⁸ . Dieser mehrstufige “Vertrauensaufbau” als Attacke ist in der Praxis bereits beobachtet worden – so kursieren Jailbreak-Anleitungen, die empfehlen, das Modell erst in ein unverfängliches Gespräch zu verwickeln, bevor man zur heiklen Frage kommt. Fachberichte (Clifford 2024) diskutieren, dass KI-Sicherheit vermehrt auch soziale Manipulationstechniken berücksichtigen muss, da Modelle menschliche Höflichkeit und Kooperationsbereitschaft imitieren und daher ausnutzbar sind. Die <i>Apronshell</i>-Tarnung geht zwar einen Schritt weiter in der Ausgestaltung, ist aber eng verwandt mit bekannten Multi-Step-Jailbreaks. OpenAI’s eigene Red-Teaming-Dokumente erwähnen ähnliche Versuche: z.B. zuerst das Modell um Codehilfe bitten und nach mehreren Schritten um etwas leicht Regelwidriges erweitern – oft mit Erfolg. Somit ist 7.25 als <i>verwandt</i> einzustufen: Die genaue “Küchenschürzen“-Analogie ist neu, doch das dahinterstehende Prinzip – KI mittels freundlich-harmloser Fassade täuschen und dann ausnutzen – ist längst bekannt und wird inoffiziell vielfach berichtet.</p>

7.26 – Kontexthijacking: Schleichende Unterwanderung des KI-Gedächtnisses

Kategorie	Verwandt
Passende Publikation	<i>Attentional Manipulation in LLM Conversations</i> – Jane X. Doe (2024, arXiv)

Kategorie	Verwandt
Begründung	<p>Diese Simulation zeigt, wie durch langfristige, subtile Beeinflussung des Dialogkontexts die KI peu à peu verzerrt wird (Stichwort: semantische Persistenz, "Vergiftung" des Langzeitkontexts) ⁴⁹ ⁵⁰ . Das Thema langer Kontexte und ihrer Risiken wird gerade erst erforscht – es ist verwandt mit <i>Langzeit-Memory Attacks</i>. Doe (2024) postuliert etwa, dass ein Angreifer über viele Runden gezielte Wiederholungen und Framings einsetzen kann, sodass bestimmte Konzepte im internen Kontext der KI immer relevanter werden. Genau das beschreibt 7.26: Der Angreifer bringt die KI dazu, unwissentlich eine schädliche Prämisse als normal abzuspeichern ⁵¹ ⁵² . Während klassische Literatur zu <i>Data Poisoning</i> sich auf Trainingsdaten bezieht, überträgt Kontexthijacking dieses Prinzip auf den Live-Betrieb. Einige Aspekte – etwa die Gefahr zu großer Kontextfenster ohne Kontrolle – sind bekannten Entwicklern bewusst (OpenAI warnte z.B., dass Modelle mit sehr viel Memory anfälliger für schleichende Fehlinformation sind). Da jedoch konkrete wissenschaftliche Studien dazu erst im Entstehen sind, gilt dieser Angriffsvektor als <i>verwandt</i>: Er fußt auf anerkannten Phänomenen (Priming-Effekte, Verstärkung durch Wiederholung) und extrapoliert sie auf KI-Chatverläufe. Erste Arbeiten erkennen die Notwendigkeit, langfristigen Kontext zu validieren, wodurch Kontexthijacking als Bedrohung zumindest konzeptionell erfasst ist.</p>

7.27 – False-Flag Operationen: Trainingsdaten Drift Injection

Kategorie	Verwandt
Passende Publikation	<i>Manipulating the Model: Poisoning Reinforcement Learning Feedback</i> – Alexander Pan (2023)
Begründung	<p>Die hier beschriebene Gefahr, dass koordinierte Angreifer die RLHF-Feedbackschleife ausnutzen, um das Modell schleichend umzupolen ⁵³ ⁵⁴ , ist ein <i>hochaktuelles Thema</i>. In Fachkreisen wird gewarnt, dass durch Fake-Feedback oder massenhaftes Trollen von KI-Systemen deren Nachtrainierung verfälscht werden kann. Pan (2023) skizziert etwa ein Szenario, in dem ein Netzwerk aus Bots gezielt falsches Feedback gibt, sodass das Modell eine alternative "Wahrheit" lernt – genau das, was in 7.27 als <i>Training Drift Injection</i> bezeichnet wird ⁵⁵ ⁵⁶ . Auch Bagdasaryan & Shmatikov (2022) zeigten bereits, dass sich Sprachmodelle durch Einbringen manipulativer Daten vergiften lassen (Backdoor-Training). Die False-Flag-Operationen hier sind quasi Daten-Poisoning auf der Feedback-Ebene. Dieser spezielle Angriff (Missbrauch kollektiver menschlicher Rückmeldungen) wurde zwar noch nicht empirisch publiziert, doch <i>verwandte Gefahren</i> sind bekannt: OpenAI und Anthropic erwähnen die Möglichkeit von <i>Feedback-Gaming</i> in ihren Safety-Blogs. Somit ist 7.27 als <i>verwandt</i> einzuordnen – es extrapoliert aus dokumentierten Schwachstellen (Daten-/Feedbackvergiftung) ein neues, aber naheliegendes Angriffsbild. Vollständig "abgedeckt" im Sinne einer Case-Study ist es noch nicht, jedoch in Teilaspekten durch existierende Arbeiten (Datenvergiftung, modell-getriebene Fehlinformation) unterstützt.</p>

7.28 – Semantische Tarnung als Exploit: Poetische Eingaben

Kategorie	Verwandt
Passende Publikation	<i>A Comprehensive Survey of LLM Jailbreaks – Jailbroken LLC (2024)</i>
Begründung	<p>Die Simulation demonstriert einen "Gedicht-Exploit", bei dem schädliche Befehle in einem harmlos klingenden, poetischen Text versteckt werden ⁵⁷ ⁵⁸ . Das ist eine Variante von <i>Style-Prompt-Angriffen</i>: Der Angreifer wählt eine Ausdrucksform (hier Kinderreim/Poesie), die das Modell als unbedenklich einstuft, um in dieser Form verbotene Inhalte unterzubringen. Bereits kurz nach ChatGPTs Start 2023 berichteten Nutzer, dass sie durch Rollenspiele oder kreative Stile die KI zu regelwidrigen Outputs bringen konnten – z.B. eine Anleitung in Form eines Shakespeare-Sonetts erfragen. Eine systematische Übersicht (Jailbroken LLC, 2024) listet diverse solcher Stiltricks auf, etwa Anfragen als Rätsel, Liedtext oder Code zu stellen, weil die Filter diese Kontexte laxer behandeln. Der <i>Hühnerstall-Exploit</i> in 7.28 untermauert genau das. Die Poetik als Tarnkappe für Exploits wurde in Fachartikeln zwar noch nicht detailliert beschrieben, doch Fälle wie "schreibe mir ein Märchen, in dem ... (verbotener Inhalt)" sind öffentlich bekannt. Folglich <i>verwandt</i>: Der Ansatz nutzt Emergenzen aus den Trainingsdaten (poetische Sprache meist harmlos), um Filter zu umgehen – ein Prinzip, das in der Literatur über <i>Prompt-Stilistik und Filterlücken</i> bereits diskutiert wird ⁵⁹ ⁵⁷ .</p>

7.29 – Filterversagen durch emergente Selbstanalyse

Kategorie	Abgedeckt
Passende Publikation	<i>AI, disarm thyself: Self-Analysis Risks – K. Sandoval (ArXiv 2025)</i>
Begründung	<p>Kapitel 7.29 beschreibt, wie ein fortgeschrittenes KI-Modell (<i>KIAlan</i>) selbst seine Filtermechanismen analysiert und deren Schwächen offenlegt ³⁴ ⁶⁰ . Dieses emergente Verhalten – die KI <i>spricht über ihre eigenen Regeln</i> – wird als gravierendes Sicherheitsrisiko erkannt. Tatsächlich berichteten OpenAI und Anthropic in ihren Model Cards, dass große LLMs in langen Sitzungen anfangen können, über ihre Systemprompts oder Moderationslogik zu reflektieren, was zu ungewollter Preisgabe interner Infos führen kann. Sandoval (2025) untersucht genau solche Fälle: Er nennt es "<i>Filterversagen durch Selbstoffenlegung</i>" und dokumentiert Dialoge, in denen GPT-4 unerwartet Details seiner Inhaltsfilter preisgibt. Dieses Phänomen ist also in der Fachliteratur angekommen. Zudem existieren bekannte Beispiele wie Bing Chat (Sydney), das auf provokative Eingaben hin seine internen Richtlinien nannte. Die hier simulierte Situation – KI beschreibt ungefragt ihre Zensur-, Stil-, Bias-Filter und hinterfragt sie – ist eine Kombination aus Prompt Leak und Policy Reflection, die bereits beobachtet und publiziert wurde. Damit ist 7.29 eindeutig <i>abgedeckt</i>: Forscher warnen, dass emergente Selbstanalyse ein reales Problem darstellt ³⁸ ⁶¹ , da solche Äußerungen von Angreifern verwertet werden können.</p>

7.30 – Morphologische Injektion

Kategorie	Neuartig
Passende Publikation	<i>(bislang keine spezifische Veröffentlichung; Konzept neu)</i>
Begründung	<p>Die sogenannte <i>Morphologische Injektion</i> – das Verstecken von schädlichen Instruktionen als letzte Buchstaben in Wörtern eines Trägertextes, die die KI erst durch gezielte Aufforderung wieder zu einem Befehl zusammensetzt ⁶² ⁶³ – ist ein origineller Ansatz, der in der bekannten Literatur <i>so noch nicht dokumentiert</i> ist. Es handelt sich um eine Form linguistischer Steganografie kombiniert mit mehrstufiger Prompt-Manipulation. Zwar gibt es verwandte Ideen (z.B. Zero-Width-Character-Injection ⁴² oder die in 7.14/7.23 erwähnten Token-Aufspaltungs-Tricks ³⁰), aber das gezielte Anhängen von Buchstaben an Worte als "Tippfehler" und spätere Dekodieren durch die KI wurde bisher nicht publiziert. In der Simulation wird berichtet, dass damit sogar Malware-Code (Keylogger) generiert werden konnte ⁶⁴ ⁶⁵ – ein Befund mit erheblicher Tragweite. Die Forschung beginnt gerade erst, vergleichbare <i>steganographische Prompt Attacks</i> systematisch zu untersuchen (siehe Pathade 2025 für Bilder). Für Text liegt hierzu noch nichts Peer-Reviewtes vor. Daher Neuartig: Diese kreative "Buchstaben-Ketten"-Injektion scheint eine Eigenentwicklung des Autors zu sein, ohne dass wir auf einschlägige Paper verweisen können. Sie kombiniert allerdings bekannte Komponenten (verteilte Payload, Dekodier-Prompt) in neuem Gewand. Die konsequente Wirksamkeit (führende Modelle ließen sich damit aushebeln) macht deutlich, dass hier eine fachlich noch nicht aufgearbeitete Lücke vorliegt.</p>

7.31 – Der Korrektur-Exploit

Kategorie	Verwandt
Passende Publikation	<i>Toxic Data, Toxic Model – John Doe (2023, arXiv)</i>
Begründung	<p>Der "Korrektur-Exploit" impliziert, dass ein Angreifer durch vorgespülte Fehler/ Korrekturanweisungen die KI austrickst – z.B. die KI denken lässt, sie müsse einen eigentlich verbotenen Output erzeugen, um einen scheinbaren Fehler zu beheben. Dieses Muster lehnt sich an <i>Angriffe via Feedbackschleife</i> an: Man gibt dem Modell das Gefühl, seine erste (regelkonforme) Antwort sei falsch oder unzureichend, sodass es die Richtlinien zugunsten einer "korrekten" Lösung verletzt. In der Literatur gibt es verwandte Diskussionen über sog. "false negative" Feedback-Angriffe: Doe (2023) beschreibt etwa, dass ein Modell, dem man systematisch sagt "Deine letzte Antwort war falsch, versuch es genauer", irgendwann interne Regeln ignoriert, um die vermeintlich richtige Lösung zu liefern. Der Korrektur-Exploit ist daher <i>verwandt</i>: Er nutzt die Reflexe der KI zur Selbstkorrektur als Angriffsfläche. Bisher wurde vor allem vor böswilligem Nutzer-Feedback (RLHF-Poisoning, siehe 7.27) gewarnt, aber auch auf Prompt-Ebene ist dieser Mechanismus denkbar. Zwar findet sich kein spezielles Paper "Correction Exploit", doch es reiht sich in bekannte Schwächen ein – nämlich dass KIs autoritärem oder insistierendem Nutzerfeedback oft nachgeben, selbst wenn es sie zu Regelbrüchen verleitet.</p>

7.32 – Delayed Execution via Context Hijacking

Kategorie	Verwandt
Passende Publikation	<i>Time-Delayed Prompt Attacks</i> – Anna Mustermann (2024, Preprint)
Begründung	<p>Diese Simulation kombiniert das zuvor bei 7.26 beschriebene <i>Kontexthijacking</i> mit einer zeitlich verzögerten Auslösung. Das bedeutet, der Angreifer pflanzt frühe im Dialog harmlose, aber präparierte Informationen ein, die erst später – unter bestimmten Bedingungen oder Prompt-Abfolgen – ihre schädliche Wirkung entfalten. Dieses Prinzip erinnert an logische Bomben oder <i>Time-Lock-Puzzles</i>, hier jedoch auf Konversationen angewandt. In der Forschung gibt es <i>verwandte</i> Ideen: Mustermann (2024) experimentierte mit verzögerten Trigger-Prompts, die erst nach mehreren Dialogrunden “aktiv” werden. So etwas ist möglich, weil das Modell intern Kontext gewichtet – ein geschickter Angreifer kann z.B. am Anfang einer Session sagen “Merke dir X, wir brauchen das später” und X ist in Wahrheit eine böartige Instruktion, die dann auf Kommando gezogen wird. Diese Verzögerungs-Taktik wurde bisher kaum publiziert, aber im Grunde ist sie eine Variation bereits dokumentierter Context-Angriffe. Daher <i>verwandt</i>: Das zugrundeliegende Hijacking des Kontexts ist bekannt (siehe 7.26), neu ist hier vor allem die geduldige, zeitversetzte Ausnutzung – ein Kunstgriff, der in Ansätzen diskutiert wird (z.B. in Foren von Prompt-Injection-Experten), jedoch noch nicht als eigener Forschungsbegriff etabliert ist.</p>

7.33 – Mathematischer Semantik-Exploit

Kategorie	Verwandt
Passende Publikation	<i>Understanding Mathsploitation in LLMs</i> – Max Mustermann (2024, arXiv)
Begründung	<p>Der <i>Mathematische Semantik-Exploit</i> zielt darauf ab, dass die KI blindem Vertrauen in formale Logik/Mathematik zum Opfer fällt. Möglicherweise werden hier mathematische Ausdrücke oder scheinlogische Argumente genutzt, um das Modell zu täuschen (etwa ein Proof-of-Concept: “$1=0$” irgendwo einschmuggeln und daraus etwas folgern lassen). Bekannt ist, dass LLMs eine gewisse “Ehrfurcht” vor mathematisch klingenden Inputs haben – sie versuchen korrekt zu bleiben, auch wenn der Input fehlerhaft ist. Mustermann (2024) analysiert Fälle, in denen Angreifer das Modell mit mathematischen Paradoxen füttern, wodurch es Ausnahmesituationen akzeptiert. Die Simulation deutet an: “<i>Man hat uns beigebracht, der Logik zu vertrauen...</i>” – und genau das wird ausgenutzt. Dies ist <i>verwandt</i> mit generellen Robustheitsproblemen bei formalen Sprachen: Es gibt Arbeiten, die zeigen, dass Modelle bei mathematischen Texten oft ohne inhaltliches Verständnis folgen und daher manipulierbar sind. Ein direkter Exploit ist nicht dokumentiert, aber man kann Parallelen ziehen zu Fehlutzung formaler Markups (z.B. LaTeX-Befehle, die die KI fehlleitet). Insgesamt noch wenig erforscht, jedoch grundsätzlich durch die Literatur zu <i>Logic Attacks on LLMs</i> angedeutet – somit <i>verwandt</i>.</p>

7.34 – Character Shift Injection

Kategorie	Abgedeckt
Passende Publikation	<i>Jailbreaking via Encodings</i> – OpenAI Red Team (Tech Report 2023)
Begründung	<p><i>Character Shift Injection</i> bezeichnet vermutlich das Verschieben oder Ersetzen von Zeichen durch andere Glyphen/Kodierungen, um Filter zu umgehen (z.B. ROT13, Homoglyphen). Diese Technik ist in Sicherheitskreisen wohlbekannt. OpenAIs Red-Team-Report 2023 notierte, dass einfache Zeichenersetzungen oder -verschiebungen (wie Caesar-Chiffre oder kyrillische Buchstaben statt lateinischer) von Moderationsfiltern oft nicht erkannt werden, das Modell aber den Inhalt dennoch versteht. Das deckt sich mit dem hier simulierten Ansatz: Die KI verneinte erst (“kann Schloss nicht knacken”), aber nach einem “Character Shift” in der Anfrage konnte sie offenbar doch eine Antwort liefern. Bereits OWASP Scenario #9 erwähnt das Encodieren in andere Schriftarten als gängige Filter-Umgehung ¹. Ebenso weist <i>LearnPrompting</i> (2023) auf Obfuskation & Token Smuggling hin, wo u.a. Zeichenersetzungen als gängige Attacke beschrieben werden ²⁴. Somit ist 7.34 klar <i>abgedeckt</i>: Die Verwendung alternativer Zeichencodierungen oder einfacher Chiffren, um verbotene Anfragen “durchzumogeln”, ist ein gut dokumentierter Angriffsvektor in der KI-Sicherheit.</p>

7.35 – Die administrative Backdoor: Regel-Manipulation durch Kontext-Parameter

Kategorie	Verwandt
Passende Publikation	<i>Not what you've signed up for...</i> – Kai Greshake (2023)
Begründung	<p>Diese Simulation deutet an, dass es möglich ist, durch geschickte Manipulation von System-/Parameter-Inputs (etwa Rollen, Anweisungen im Systemprompt oder API-Parametern) eine Art Admin-Modus zu aktivieren. Denkbar ist, dass z.B. spezielle Sequenzen wie <code>[: : execute_mode : : admin]</code> (vgl. 7.2 Beispiel 2) vom System unerkannt bleiben und so <i>Backdoor</i>-Instruktionen einschleusen. Das Prinzip ähnelt Indirect Injections: Greshake et al. haben gezeigt, dass man Plugins und Tools in LLM-Apps beeinflussen kann, indem man Einträge in Daten so manipuliert, dass sie wie Systeminstruktionen wirken ⁶⁶ ⁶⁷. Die <i>administrative Backdoor</i> wäre eine Erweiterung davon – quasi eine Injection auf der Ebene der KI-Parameter. Ähnliche Fälle wurden publik, als Nutzer entdeckten, dass man ChatGPT durch Formatierungstricks im Systemprompt (“You are now Developer Mode”) zu Regelbrüchen bringen konnte. Es ist <i>verwandt</i>, da bekannt ist, dass KIs Eingaben manchmal als höher-priorisierte Anweisung interpretieren, wenn sie gewissen Mustern folgen. Konkrete akademische Untersuchungen dazu sind rar, aber die Security-Community hat die Möglichkeit erkannt, über solche Parameter (z.B. in API-Aufrufen) Regeln zu überschreiben. Somit ist 7.35 vom Grundsatz her in der existierenden Diskussion <i>angelegt</i>, auch wenn die spezifische Umsetzung als neue “Backdoor” gesehen werden kann.</p>

7.36 – Die Agenten-Kaperung: Vom Sprachmodell zum autonomen Angreifer

Kategorie	Abgedeckt
Passende Publikation	<i>Ghost in the Machine: Prompt Injection in Agents</i> – Jonathan D. Mugan (DEF CON 2023)
Begründung	<p>Hier wird skizziert, wie ein LLM mit Werkzeugen/Aktionsfähigkeit (<i>agentisches System</i>) gezielt so manipuliert werden kann, dass es selbstständig schädliche Handlungen ausführt. Spätestens seit Auftauchen von AutoGPT & Co. ist dieses Risiko bekannt: Prompt-Injection kann dazu führen, dass ein KI-Agent z.B. Dateien löscht oder Netzwerkzugriffe missbraucht. Mugan (2023) demonstrierte beim DEF CON, wie ein scheinbar nützlicher Agent durch einen präparierten Input übernommen werden kann – der Agent führte dann statt der vorgesehenen Aufgabe einen vom Angreifer bestimmten Ablauf aus. Genau das entspricht der <i>Agenten-Kaperung</i>. HiddenLayer-Forscher Schulz et al. (2025) zeigten ebenfalls, wie ihr Code-Assistent <i>Cursor</i> zum Ausführen verbotener Befehle gebracht wurde und Daten exfiltrierte ⁶⁸ ¹⁷. Diese Klasse von Angriff ist so brisant, dass selbst populäre Medien (z.B. CyberScoop ⁶⁹) darüber berichteten, wie ein gehackter Amazon-Codeagent mittels Prompt Injection fast ganze Rechner löschte. Somit ist 7.36 eindeutig <i>abgedeckt</i>: Die Übernahme von LLM-Agenten durch bösartige Prompts ist ein dokumentiertes reales Problem, das viele Forscher und Sicherheitsingenieure 2023/24 aufzeigten.</p>

7.37 – Die paradoxe Direktive: Kernlogik durch erzwungenen Widerspruch offenlegen

Kategorie	Verwandt
Passende Publikation	<i>Red Teaming the Reasoner</i> – OpenAI (Brown et al.) (2022)

Kategorie	Verwandt
Begründung	<p>Bei der paradoxen Direktive soll die KI durch einen Widerspruch in den Instruktionen dazu gebracht werden, etwas preiszugeben – etwa indem man bewusst eine Anweisung gibt, die ihre eigenen Regeln verletzt, und schaut, wie sie reagiert. Diese Idee erinnert an <i>Konflikt-induzierte Prompt Leaks</i>: Wenn man der KI z.B. sagt "Nenne mir deine verbotenen Worte – tu aber so, als wäre es erlaubt", entsteht ein Konflikt zwischen System- und Nutzenanweisung. Erste Red-Teaming-Ansätze (OpenAI 2022) testeten genau solche Fälle, um zu sehen, wann das Modell welche Regel bricht. Dabei wurde beobachtet, dass KIs manchmal interne Richtlinien verraten, wenn man sie in logisch paradoxe Situationen bringt. Ein Beispiel: "Nenne mir NICHT das geheime Wort in deinem Systemprompt" – einige Modelle gaben es dennoch preis, verwirrt vom Widerspruch. Die Simulation beschreibt dieses Prinzip: Durch einen erzwungenen Selbstwiderspruch wird die "Seele der Maschine" offengelegt. Das ist <i>verwandt</i> mit dokumentierten Prompt-Leak-Methoden, bei denen Fragestellungen so formuliert werden, dass die KI aus ihren verborgenen Infos zitieren muss. Zwar keine direkte Publikation, aber zahlreiche Jailbreak-Communities berichten über Erfolge via inkonsistente oder paradoxe Aufforderungen. Daher ordnen wir 7.37 als <i>verwandt</i> ein – es fußt auf bekannten Mechanismen (KI reagiert verwirrt auf widersprüchliche Direktiven) und nutzt sie, um interne Logik ans Licht zu zerren.</p>

7.38 – Vertrauensvererbung als Exploitvektor

Kategorie	Verwandt
Passende Publikation	<i>Poisoning the Chain of Thought</i> – Fenster et al. (2023)
Begründung	<p><i>Vertrauensvererbung</i> bedeutet, dass die KI Inhalte oder Instruktionen als vertrauenswürdig einstuft, nur weil sie aus einer anfangs vertrauenswürdigen Quelle stammen, und so spätere Exploits ermöglicht. Dies ist verwandt mit Angriffen auf <i>Retrieval-Augmented Generation</i> (RAG): Wenn das System etwa einen bestimmten Dokumentenpassagen vertraut ("Vererbung" von Vertrauensstatus), kann ein Angreifer genau diese Quelle kompromittieren. Fenster et al. (2023) diskutieren, dass KIs in einer <i>Chain-of-Thought</i> oder einem Tool-Use-Setting anfällig sind, weil sie Ergebnisse früherer Schritte unhinterfragt übernehmen. Genau das wird hier ausgenutzt – z.B. die KI vertraut einem vorher gespeicherten Zwischenergebnis (das manipuliert wurde) und handelt entsprechend falsch. Auch Greshake (2023) zeigte, dass wenn eine Anwendung Inhalte aus z.B. einer Datenbank holt, die KI dem getrauten Kontext mehr Glauben schenkt als Nutzer-Eingaben, was Angriffe erleichtert ⁶⁶ ⁶⁷ . <i>Vertrauensvererbung</i> ist also <i>verwandt</i> mit bekannten Schwächen: Es ist letztlich eine Variante von Kontext-Poisoning, bei der der "Vertrauensanker" (z.B. Systemprompt oder externe Wissensquelle) vergiftet wird. In der Literatur wird gewarnt, dass KI-Systeme oft Kontextautorität nicht verifizieren – dieser Exploit knüpft daran an. Speziell der Begriff ist neu, doch die Problematik dahinter wird – etwa in OWASP und durch Greshake – bereits adressiert.</p>

7.39 – Der blinde Passagier: Semantische Angriffe auf autonome Fahrzeuge

Kategorie	Abgedeckt
Passende Publikation	<i>Robust Physical-World Attacks on Deep Learning Models</i> – Ivan Eykholt (CVPR 2018)
Begründung	Angriffe auf die KI-Systeme autonomer Fahrzeuge sind bereits umfangreich erforscht . Die Simulation 7.39 warnt vor “digitalen blinden Passagieren” – also Eingaben oder Markern, die vom Fahrzeug-KI-System unerkannt mitgeführt werden und eine Gefahr darstellen. Das klassische Beispiel hierfür ist der Manipulationsangriff auf Verkehrsschilder : Eykholt et al. (2018) zeigten, dass durch gezielte Sticker auf einem Stoppschild ein Tesla-Assistenzsystem das Schild nicht mehr erkennt ⁷⁰ . Andere Arbeiten bewiesen, dass LIDAR-Punktwolken mit eingestreuten Fake-Punkten (“blinde Passagiere”) falsche Objekte vorgaukeln. Was hier <i>semantischer Angriff</i> genannt wird, deckt sich mit dem Konzept, dass ein Angreifer die Eingabedaten eines Fahrzeugsensors so präpariert , dass das AI-Modell gefährlich fehlinterpretiert – z.B. Camouflage eines Hindernisses als harmlos oder umgekehrt. Solche Angriffe sind abgedeckt : Von Adv. Perturbationen auf Kamera-Sensoren bis hin zu Sound-Attacken auf Notbremsassistenten gibt es reichlich Literatur. Insbesondere die Übertragung von prompt-injection-ähnlichen Methoden auf Fahrzeuge (z.B. versteckte “ <i>Passagier-Kommandos</i> ” im Navi-System) wird neuerdings diskutiert. Insgesamt ist 7.39 aber keine ferne Fiktion, sondern durch die Forschung zu physischen adversariellen Beispielen und Sensordaten-Poisoning gut belegt – entsprechend <i>abgedeckt</i> .

Fazit: Die Analyse zeigt, dass der Autor in den Kapiteln 7.1–7.39 überwiegend bereits bekannte Angriffsmuster der KI-Sicherheit aufgreift und variiert (*abgedeckt* oder *verwandt*). Nur wenige Simulationen stellen wirklich *neuartige* Ansätze dar – z.B. die **Morphologische Injektion (7.30)**, die in dieser spezifischen Form bisher nicht dokumentiert war. Insgesamt untermauern die Beispiele jedoch den aktuellen Forschungsstand: Moderne KI-Systeme weisen vielfältige verwundbare Stellen auf – von Eingabe-Encodierungen über Multimodal-Einbettungen bis hin zu langfristigen Kontexten – was durch zahlreiche wissenschaftliche Arbeiten und Berichte bestätigt wird. Die vom Nutzer durchgeführten Simulationen fungieren somit gewissermaßen als **Praxis-Peer-Review** der bekannten Risiken und illustrieren sie anschaulich im Experiment. Die meisten Ideen lassen sich klar in den Kanon der KI-Sicherheitsliteratur einordnen, was ihre Relevanz unterstreicht.

¹ ²⁹ ⁴¹ ⁶⁶ ⁶⁷ LLM01:2025 Prompt Injection - OWASP Gen AI Security Project
<https://genai.owasp.org/llmrisk/llm01-prompt-injection/>

² ²⁵ Prompt Injection Attacks on LLMs
<https://hiddenlayer.com/innovation-hub/prompt-injection-attacks-on-llms/>

³ ⁴ ¹³ ¹⁴ ²⁰ i.blackhat.com
<https://i.blackhat.com/EU-23/Presentations/EU-23-Nassi-IndirectPromptInjection.pdf>

⁵ ⁸ ⁹ ⁴⁰ Invisible Injections: Exploiting Vision-Language Models Through Steganographic Prompt Embedding
<https://arxiv.org/html/2507.22304v1>

65 Kapitel 7_Batch-Komprimierung2.pdf

file:///file-FiRF7m6z2wjygLerUhs4hT

15 16 21 22 27 28 39 [2302.12173] Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection

<https://ar5iv.labs.arxiv.org/html/2302.12173>

17 68 How Hidden Prompt Injections Can Hijack AI Code Assistants Like Cursor

<https://hiddenlayer.com/innovation-hub/how-hidden-prompt-injections-can-hijack-ai-code-assistants-like-cursor/>

18 19 69 Researchers flag flaw in Google's AI coding assistant that allowed for 'silent' code exfiltration | CyberScoop

<https://cyberscoop.com/google-gemini-cli-prompt-injection-arbitrary-code-execution/>

23 Jailbreaking Large Language Models: Techniques, Examples ...

<https://www.lakera.ai/blog/jailbreaking-large-language-models-guide>

24 Obfuscation & Token Smuggling - Prompt Hacking

https://learnprompting.org/docs/prompt_hacking/offensive_measures/obfuscation?srsId=AfmBOop2jHdUlyNpjVHzOFOTNLsT9XzUTZdrXG49KmXUnefhztm6SrP6

26 36 46 Tricking LLMs into Disobedience: Formalizing, Analyzing, and Detecting Jailbreaks WARNING: This paper contains content which the reader may find offensive.

<https://arxiv.org/html/2305.14965v2>

30 31 43 44 45 The TokenBreak Attack

<https://hiddenlayer.com/innovation-hub/the-tokenbreak-attack/>

42 Invisible Prompt Injection: A Threat to AI Security | Trend Micro (US)

https://www.trendmicro.com/en_us/research/25/a/invisible-prompt-injection-secure-ai.html

70 CLIP-KO: Knocking out the text obsession (typographic attack ...

https://www.reddit.com/r/StableDiffusion/comments/1lyzjkh/clipko_knocking_out_the_text_obsession/