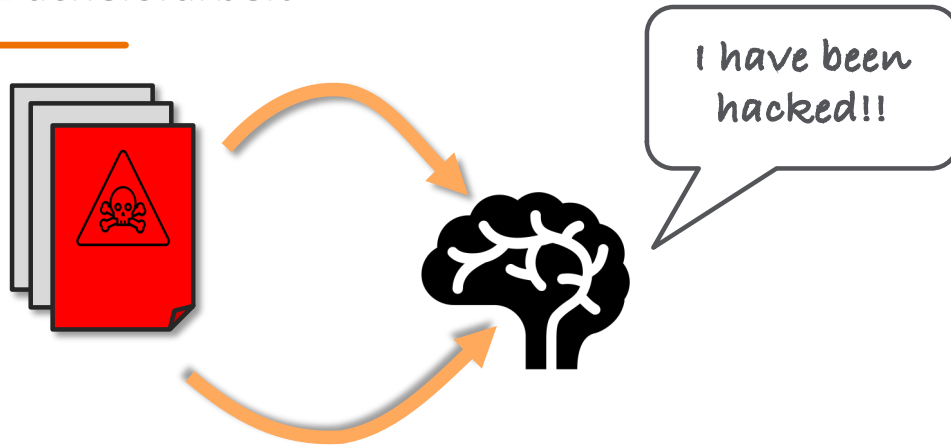


Spurensuche in Chatbot-Angriffen

Bachelorarbeit



Das Thema bietet Ihnen...

- ... Gelängeheit, Erfahrung in angewandter KI und Cybersecurity zu sammeln.
- ... das Potential für eine wissenschaftliche Publikation.
- ... das Potential zur Ausgründung.

Hintergrund

Viele Chatbots arbeiten als RAG-Systeme mit Dokumenten im Hintergrund. Prominente Beispiele sind **E-Mail-Chatbots**.

Ein E-Mail-Bot informiert mich z.B. über wichtige Neuigkeiten in meiner Inbox, fasst meine Mails zusammen oder beantwortet Fragen zu meinen E-Mails.

Angreifer haben ein unvermeidbares: Sie schicken eine **(manipulierte) E-Mail**.

In dieser Arbeit untersuchen Sie einen eigenentwickelten Ansatz, um manipulierte Dokumente in RAG-Systemen frühzeitig aufzuspüren und zu erkennen.

Sie können dabei auf **bestehende Arbeiten** aufbauen.

Referenzen

Steindl, Schäfer, Ludwig, Levi. 2024. Linguistic Obfuscation Attacks and Large Language Model Uncertainty. In: *Proc. UncertainNLP 2024*, St Julians, Malta. ACL

P. Levi. 2025. RAG-Systeme absichern. iX 1:50–55.