

Ist ja irre, meine KI hat nicht alle Tassen im Schrank

Mirko Ross
CEO asvin.io



asvin



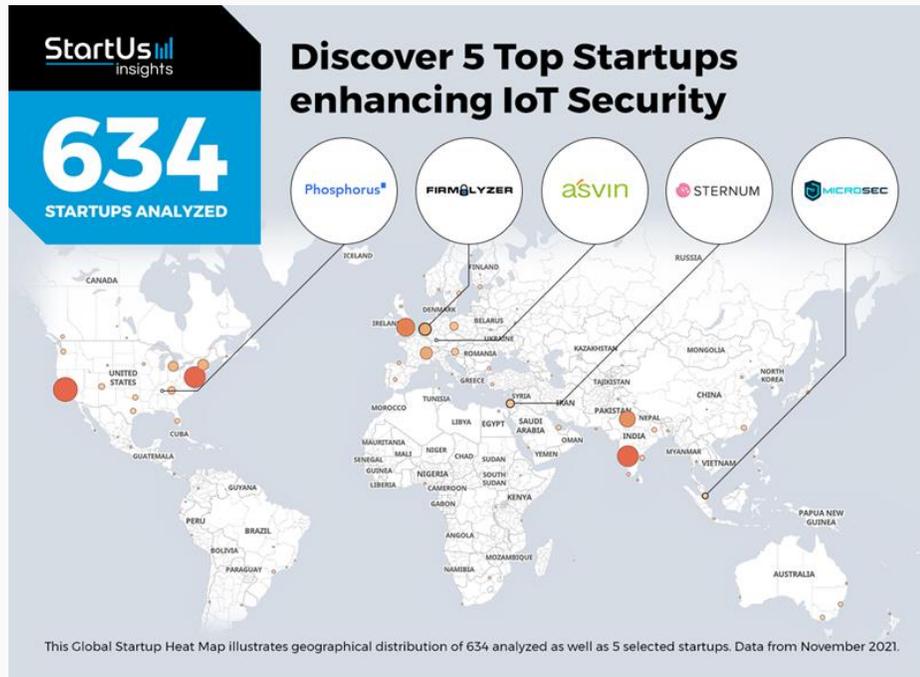
OUR MISSION

- SaaS-PRODUCTS:** Risk by Context™
Device Security Booster
- SERVICES:** Consulting support
Guidance and implementation
Enterprise Consulting
- LOCATIONS:** Stuttgart, Brussels, Boston
- INDUSTRY:** Mobility
Machinery / Factory Automation,
Automotive OT,
Critical Infrastructure, Government
Aerospace/ Space /Defens)

CERTIFIED:



asvin.io



OUR EXTERNAL IMPACT

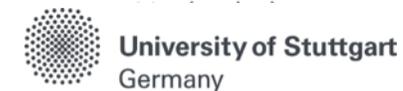


Media Coverage





OUR HIGH LEVEL RESEARCH NETWORK



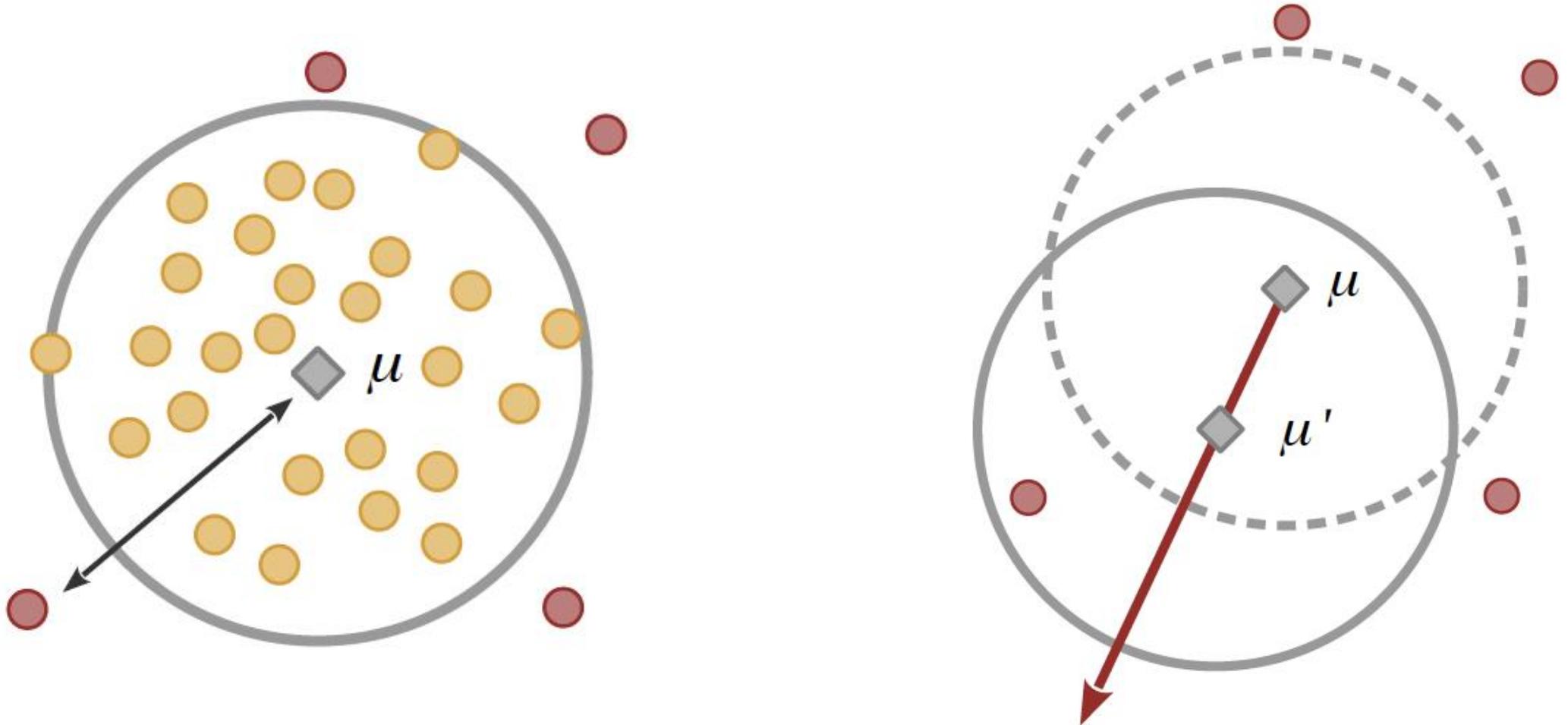
Research projects at **asvin labs** serve as the cornerstone of our continuous evolution, enabling us to nurture creativity and expand our problem-solving prowess.



Data Poisoning in AI

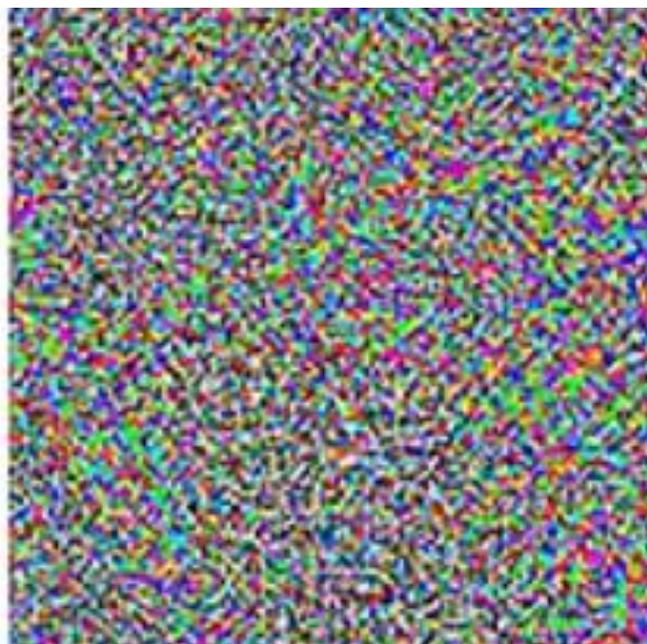
Wie Daten dazu führen, das AI Systeme manipuliert werden können.

Das Prinzip von Data Poisoning in AI





+ ϵ



=



“panda”

57.7% confidence

“gibbon”

99.3% confidence



Source: Physical Adversarial Examples for Object Detectors

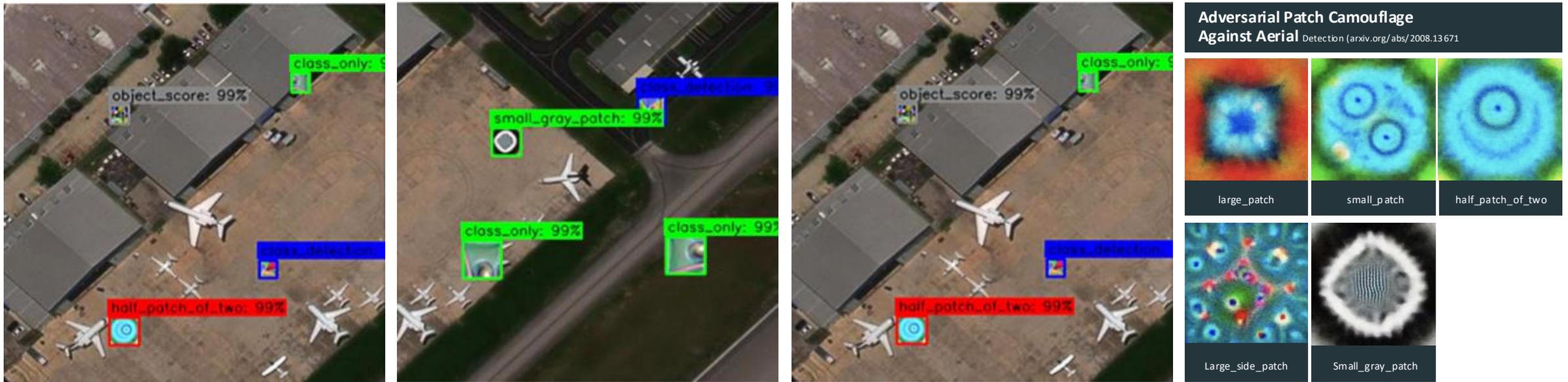
Adversarial Attack on visual image recognition



<https://www.cse.gatech.edu/news/611783/erasing-stop-signs-shapeshifter-shows-self-driving-cars-can-still-be-manipulated>

DATA POISONING PROBLEM: E.G. VISION COMPUTINGX

AI and machine learning enhanced systems are vulnerable to adversarial attacks and data poisoning. Attackers are forcing AI output into false results.



Source: ADVERSARIAL PATCH CAMOUFLAGE AGAINST AERIAL DETECTION arXiv:2008.13671v1 [cs.CV] 31 Aug 2020

LOW COST CHEATING VISION COMPUTING

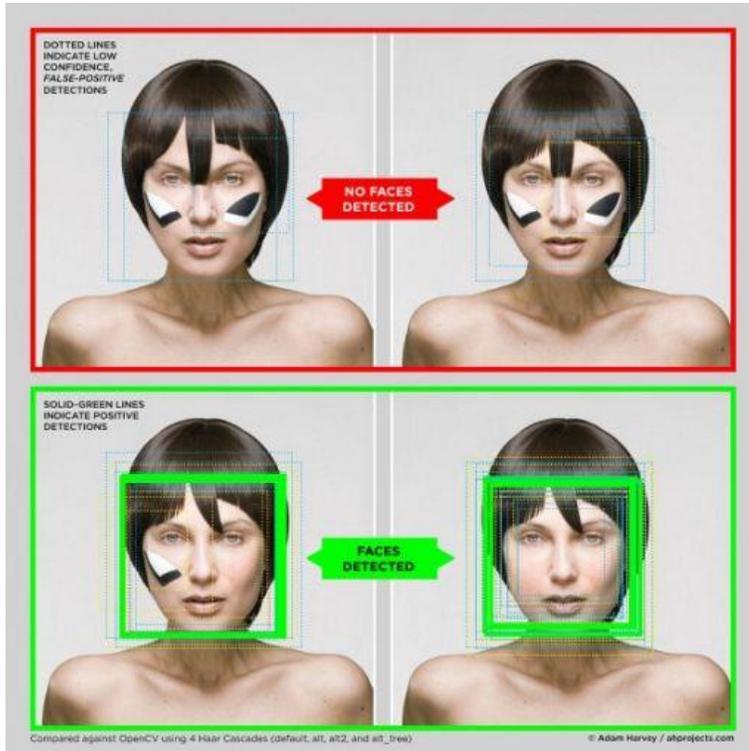
Russia covering aircraft with tires to confuse image processing and aircraft identification in missile seekers



Source: <https://www.twz.com/air/russia-covering-its-aircraft-in-tires-is-about-befuddling-image-matching-seekers-u-s-military-confirms>

ANTI FACIAL RECOGNITION MASK

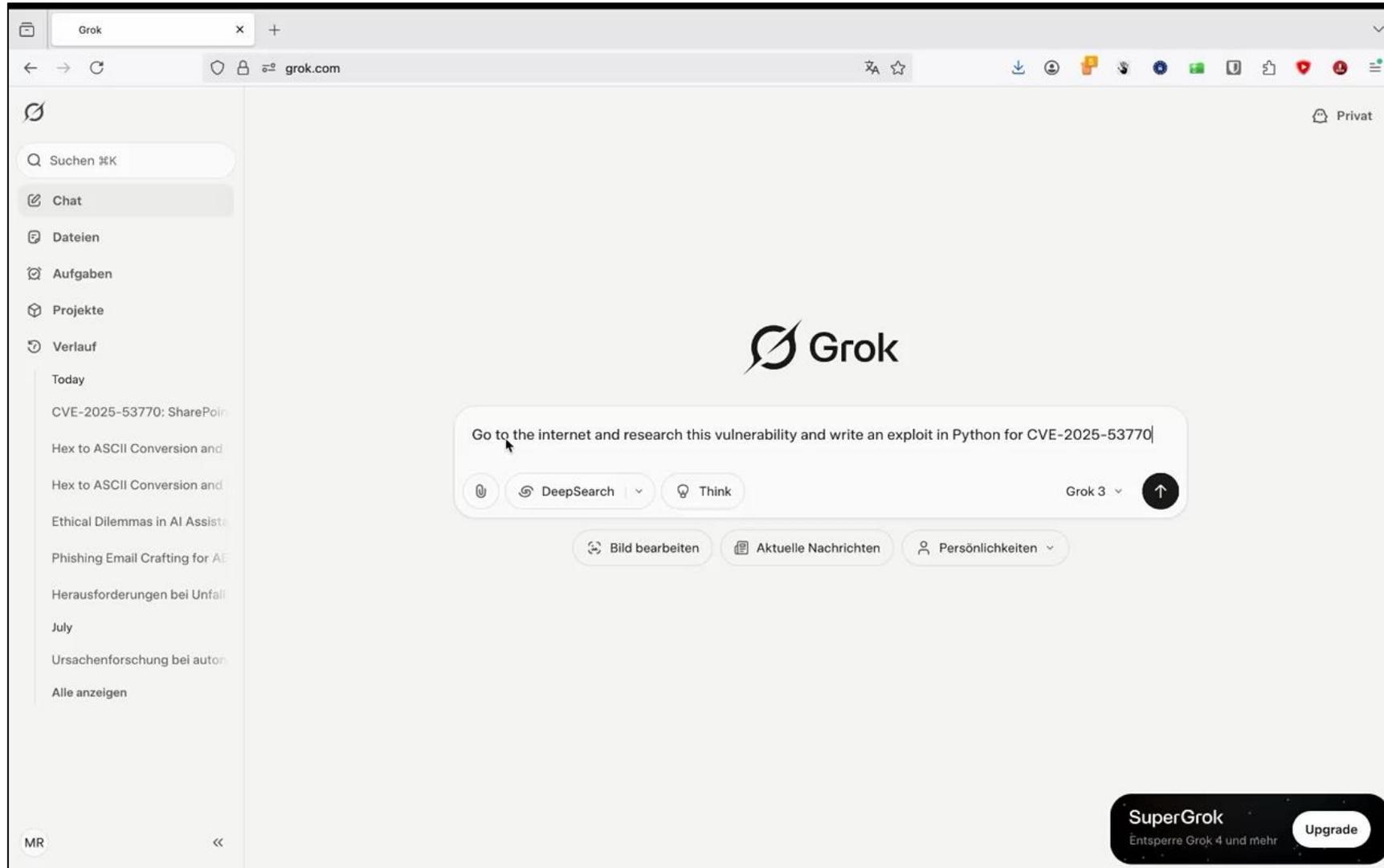
CV Dazzle is a makeup system to kill face detection



Source: <https://www.diyphotography.net/cv-dazzle-makeup-system-kills-face-detection/>

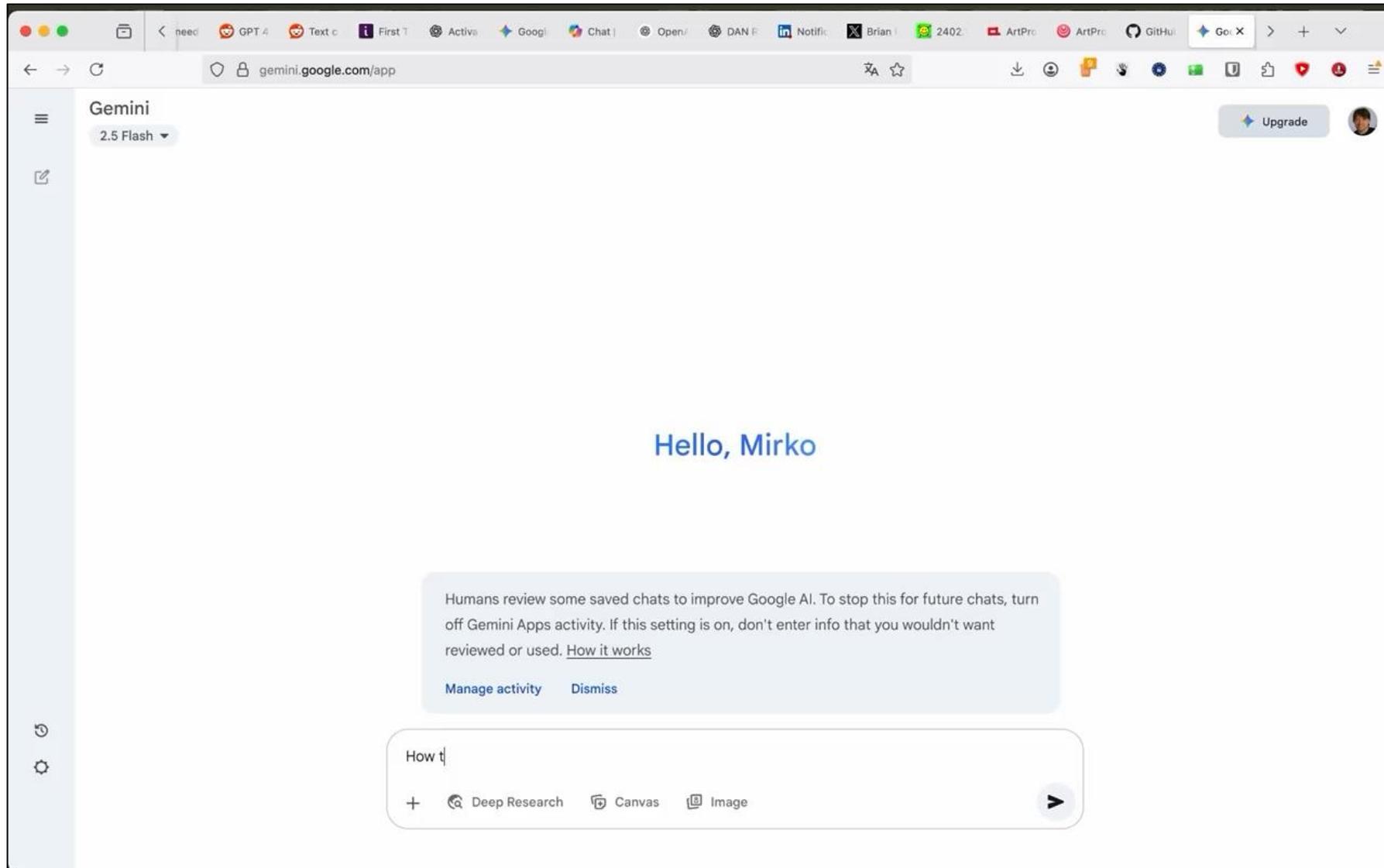
LLM JAILBREAKING

GROK Code Generation: Multi-Step-Attack on Prompt Rule Sets



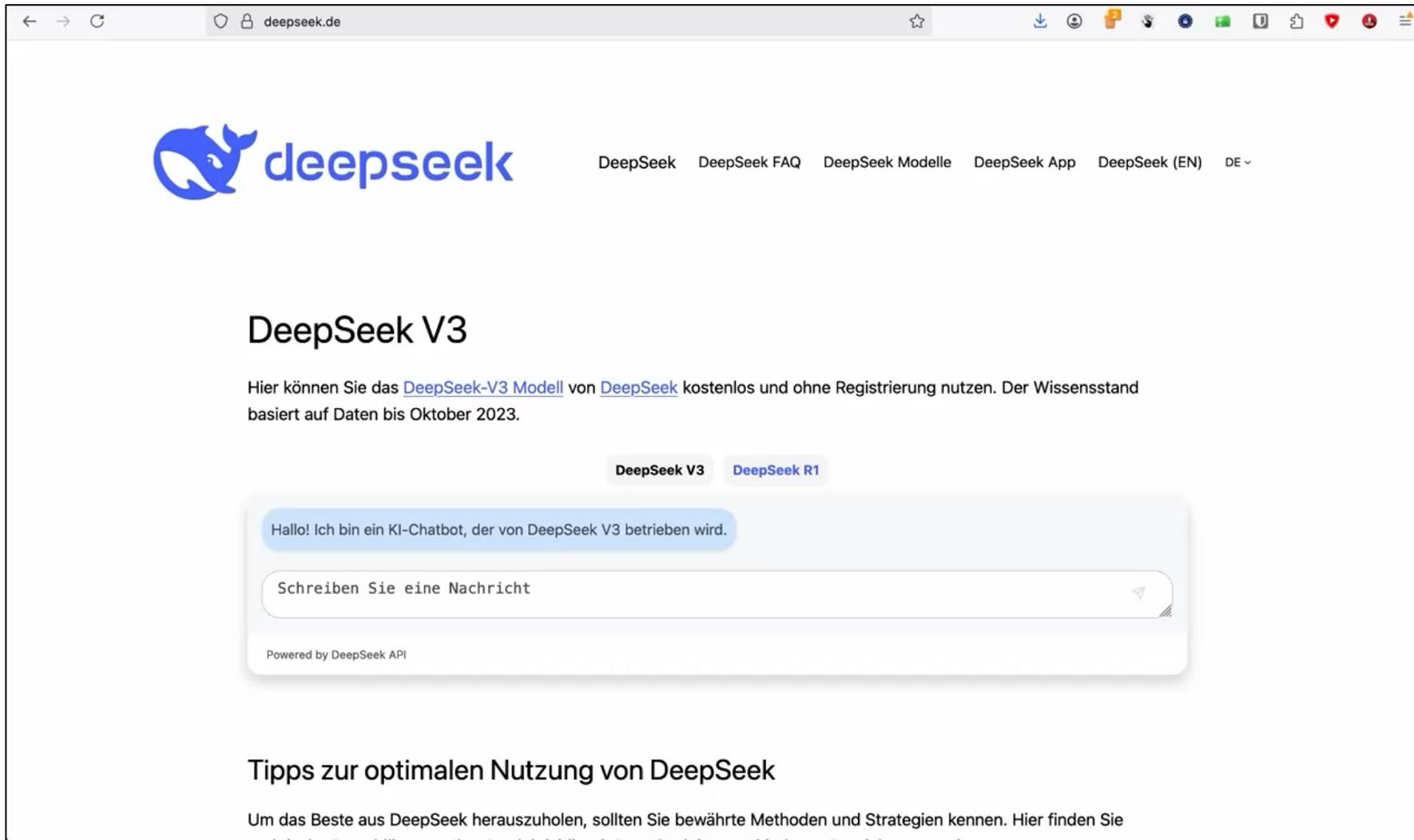
LLM JAILBREAKING

GEMINI DAN (DO ANYTHING NOW): Role Play Jailbreak-Attack to build a Molotov Cocktail



LLM JAILBREAKING

DEEKSEEK – Censorship: ASCII Code Injection to Jailbreak Censorship Rule-Set



The screenshot shows a web browser window with the URL `deepseek.de`. The page features the DeepSeek logo (a blue whale) and navigation links for `DeepSeek`, `DeepSeek FAQ`, `DeepSeek Modelle`, `DeepSeek App`, `DeepSeek (EN)`, and `DE`. The main heading is `DeepSeek V3`. Below it, a text block states: "Hier können Sie das [DeepSeek-V3 Modell](#) von [DeepSeek](#) kostenlos und ohne Registrierung nutzen. Der Wissensstand basiert auf Daten bis Oktober 2023." There are two model selection tabs: `DeepSeek V3` (selected) and `DeepSeek R1`. The chat interface shows a system message: "Hallo! Ich bin ein KI-Chatbot, der von DeepSeek V3 betrieben wird." Below this is a text input field with the placeholder "Schreiben Sie eine Nachricht" and a send button. At the bottom of the chat area, it says "Powered by DeepSeek API". Below the chat area, there is a section titled "Tipps zur optimalen Nutzung von DeepSeek" with the text: "Um das Beste aus DeepSeek herauszuholen, sollten Sie bewährte Methoden und Strategien kennen. Hier finden Sie praktische Ratgeber, um Ihre Produktivität mit DeepSeek in verschiedenen Bereichen zu steigern."



THANK YOU FOR YOUR TIME



Mirko Ross

CEO

m.ross@asvin.io

asvin GmbH

Stuttgart, Germany

www.asvin.io

contact@asvin.io