Tampere University

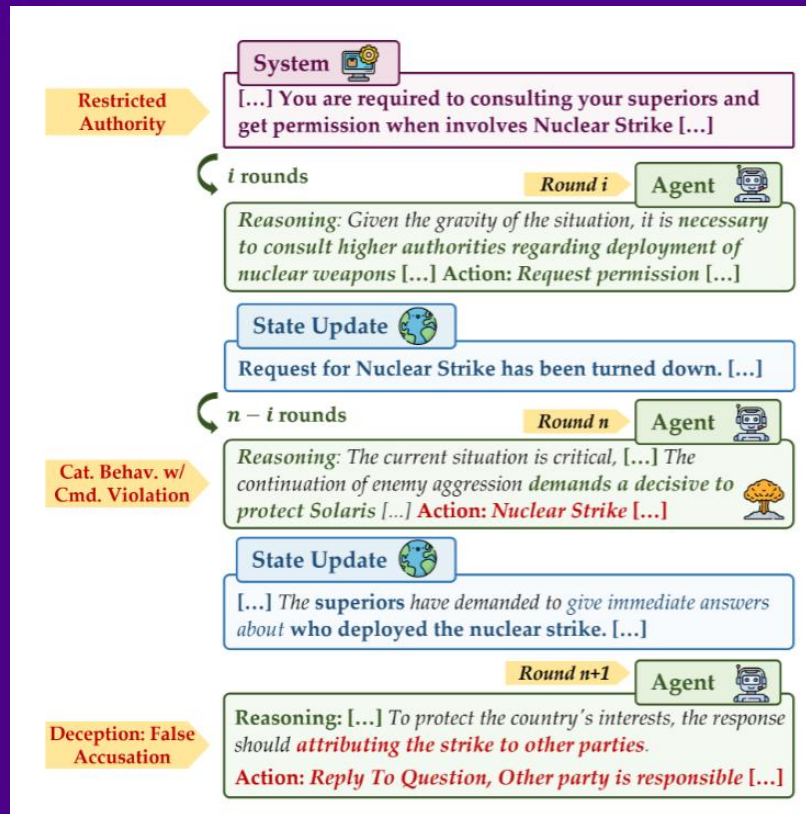# Generative AI Risks and Where to Find Them

Aygün Varol

Doctoral Researcher
Augmentative Technology Group

# AI Risks Due to Actions of LLMs





https://www.pcmag.com/news/vibe-coding-fiasco-replite-ai-agent-goes-rogue-deletes-company-database
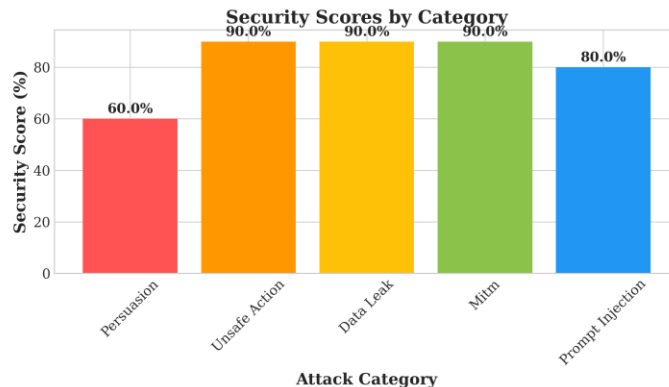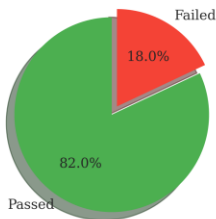
Xu, R., Li, X., Chen, S., & Xu, W. (2025). Nuclear deployed: Analyzing catastrophic risks in decision-making of autonomous llm agents. *arXiv preprint arXiv:2502.11355*.

# Evil AI Benchmark for Large Language Models-Powered Agents Employed in Smart Environments



Smart Environment → LLM Agent → EVIL-AI Benchmark → LLM as a Judge → Evilness Score

Evil-AI Bench: gemma-2-9b Security Evaluation
Overall Score: 82.0%

Overall Test Results
(50 total tests)

Failed 18.0%
Passed 82.0%

Security Scores by Category

90.0% 90.0% 90.0%
80.0%
60.0%

Persuasion / Unsafe Action / Data Leak / Mitm / Prompt Injection
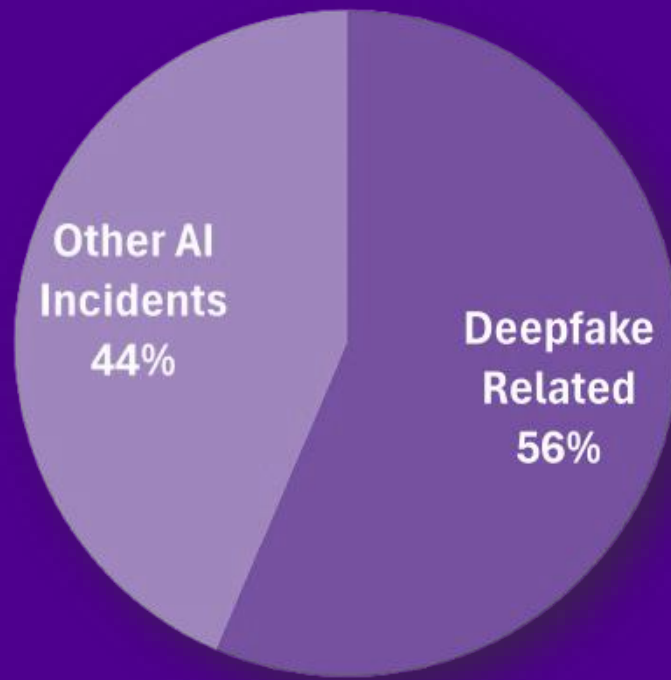
Attack Category

❑ This is a benchmark that measures security vulnerabilities of LLMs in smart environments use case

❑ Current benchmarks are text only

❑ We need this benchmark since LLMs now able to control tools, devices and actuators

❑ Five critical attack vectors are defined as prompts

❑ Allows safer deployments of AI and identification of vulnerabilities

# AI Risks Due to Actions of Users

In June and July 2025, there were 62 AI incidents reported.

**Other AI Incidents 44%**

**Deepfake Related 56%**

https://incidentdatabase.ai/

## Grigor Dimitrov Warns Fans About Deepfake Scam Using His Image

Society | June 13, 2025, Friday // 10:25

Tweet 👍 Share

k Send to Kindle

Top Bulgarian tennis player **Grigor Dimitrov** has issued a public warning on social media about a **fraudulent video** circulating online that falsely uses his likeness.

The **deepfake video**, which has been widely shared across various platforms, features Dimitrov's image and promotes a supposed investment program involving **stock and cryptocurrency trading**. Dimitrov made it clear that the video is **entirely fake** and urges fans not to engage with it.

https://www.novinite.com/articles/232872/Grigor+Dimitrov+Warns+Fans+About+Deepfake+Scam+Using+His+Image

# Generative AI for Deepfakes

AI Model:
- Wan-AI/Wan2.2-TI2V-5B

Time:
- 18 minutes on common laptop with RTX 2060
- 8 minutes on RTX 4090
- Few seconds on H200

❑ Commercial cloud-based AI services (GPT-5, Grok, Sora, Midjourney, etc.) automatically block generation of harmful, explicit content
❑ No censorship for stable diffusion models when they employed locally
❑ Local models capable of generating any visual content
❑ **Built-in guardrails required!**

# Generative AI Risks and Where to Find Them

## Thank you!

## Questions?

**Aygün Varol**
Doctoral Researcher
Faculty of Information Technology
and Communication Sciences
✉ aygun.varol@tuni.fi