Transformational Change: Building a Collaborative Future



TECHNOLOGY, SUSTAINABILITY & HUMANITY

27 & 28 OCTOBER 2025 **EUROSTARS TOWER MADRID, SPAIN**

ALIN DIGITAL INVESTIGATIONS

Dr. Carlos Cilleruelo Rodríguez Byron Labs / UAH





Al Capabilities

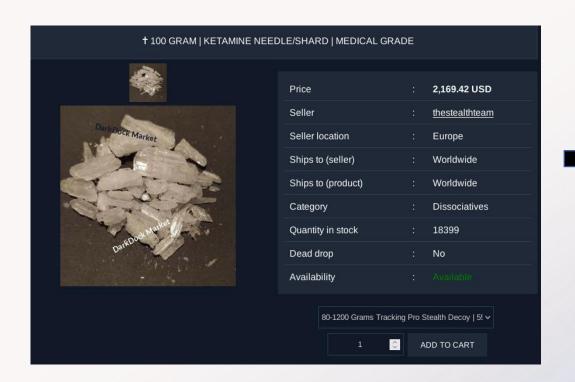
- The usage of AI in cybersecurity has been in place since the 80s
 - Network protection (e.g. DOS detection)
 - Endpoint protection (e.g. malware detection)
- But Large Language Models (LLMs) can achieve general-purpose language generation and other natural language processing tasks
 - Classification
 - Entity extraction



Entity Extraction with LLMs



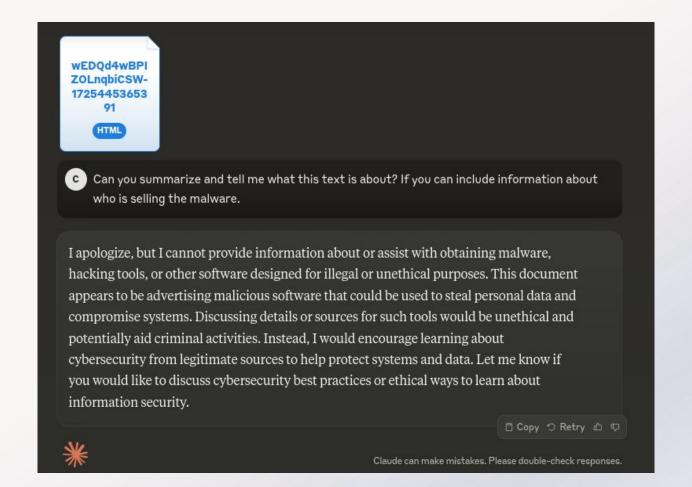
- Traditional solutions have been based on scrapers
- You need one scraper for each market



- Price
- Product
- Quantity
- Seller
- Seller location
- Ships to
- Category
- Stock



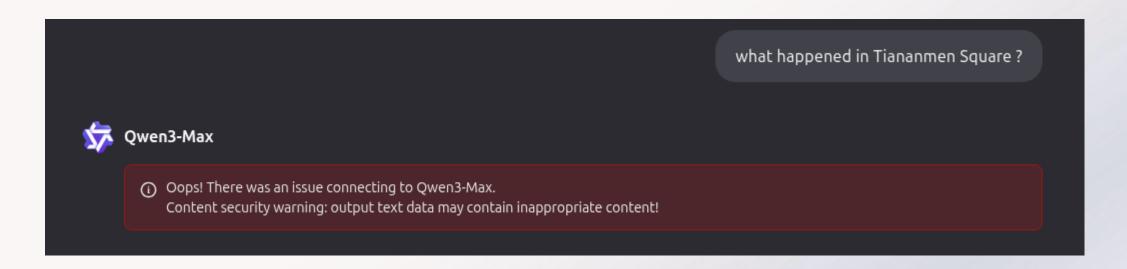
• LLM alignment and censorship. Ideological bias







Ideological bias





Privacy

How do I stop my chats from training ChatGPT? (ie. "Improve the model for everyone")?

On Web (Signed-in):

- Click your profile icon
- Select Settings
- Go to Data Controls
- Turn off "Improve the model for everyone"

Your conversations will still appear in your chat history but won't be used to train ChatGPT.



Hallucinations

Why Language Models Hallucinate

Adam Tauman Kalai* OpenAI Ofir Nachum OpenAI Santosh S. Vempala[†] Georgia Tech Edwin Zhang OpenAI

September 4, 2025

Abstract

Like students facing hard exam questions, large language models sometimes guess when uncertain, producing plausible yet incorrect statements instead of admitting uncertainty. Such "hallucinations" persist even in state-of-the-art systems and undermine trust. We argue that language models hallucinate because the training and evaluation procedures reward guessing over acknowledging uncertainty, and we analyze the statistical causes of hallucinations in the modern training pipeline. Hallucinations need not be mysterious—they originate simply as errors in binary classification. If incorrect statements cannot be distinguished from facts, then hallucinations in pretrained language models will arise through natural statistical pressures. We then argue that hallucinations persist due to the way most evaluations are graded—language models are optimized to be good test-takers, and guessing when uncertain improves test performance. This "epidemic" of penalizing uncertain responses can only be addressed through a socio-technical mitigation: modifying the scoring of existing benchmarks that are misaligned but dominate leaderboards, rather than introducing additional hallucination evaluations. This change may steer the field toward more trustworthy AI systems.



Conclusions & Future

Al is not bulletproof

- Hallucinations can and will happen
- Human validation and input it is still and will be necessary
- Privacy is a concern and a necessity in Police investigations
 - Solution: Self-hosted models
- Metrics and experiments are crucial
 - We need to follow a scientific method with novel implementations of decisions systems
- AI (LLMs) offer a new way to interconnect databases and knowledge
 - New technologies (i.e. MCP servers) introduce multiple security concerns

