



## **Reliable and Trustworthy Al**

Challenges in the era of LLMs and Agentic AI

**Gerhard Wunder** 

Cybersecurity and Al Group, FU Berlin Fraunhofer AISEC Munich Chatbot Summit Berlin 2025

## **AG Cybersicherheit und KI**

Cybersicherheit ist die Achillesferse der modernen digitalen Gesellschaft" Horst Seehofer. Bundestagsdebatte Januar 2021, anlässlich der Gesetzes-Novelle zur Sicherheit von informationstechnischen Systemen

- AG C-AI besteht seit April 2021 und ist aus Heisenberg Professur hervorgegangen
- Unterstützt durch Bundesdruckerei GmbH mit Kompetenzzentrum KI (KIKC)
- Gruppe > 20, > 8 Mill. Funding, > 200 Artikel
- "Zentrum für vertrauenswürdige Künstliche Intelligenz (ZVKI)"
- Digitalgipfel, Berlin Dezember 2022, "Kann man KI überhaupt erklären?"





C-Al ist in 3 SSP vertreten



Federated Learning und Blockchain für IoT zv ki

Zentrum für vertrauensvolle KI unterstützt vom BMUV 6G-RIC Research and Innovation Cluster

6G Research Cluster mit >30 Partnern und >70 M€ Funding unterstützt vom BMBF Quantum Crypto Analysis and Agile and Post-Quantum-Secure Protocols for the Finance Industry funded by BMBF until 2027 Bundesministerium für Bildung und Forsch

#### **Pinnepedia Technologies**



Chat-Bot zur Vervollständigung der Analyse

Generative KI zur Textgenerierung und Graphenerstellung

© Pinnipedia Technologies GmbH

#### **German Al Beacon Project: AlgenCY**

**Mission: Exploring chances and risks of generative AI** 

The project AlgenCY aims at:

- Exploring fundamental research on existing and forthcoming threats for and through generative AI
- Developing measures that improve the detection and defense against generator-based cyber attacks
- Building experimental laboratory for the systematic evaluation of generative AI methods







Bundesministerium für Bildung und Forschung Initial funding period:

November 2023 – October 2026

## Explainability

EU Act Article 13: (Transparency and Provision of Information to Deployers): This article states that highrisk AI systems must be designed to be transparent, so that those using them can understand and use them correctly.

EU Act Article 10: (Obligation concerning the de-biasing of AI systems) The providers of AI systems must evaluate whether their training, validation and testing datasets meet certain quality criteria.

## Why Explainability? – Clever Hans Effect

• Clever Hans was a (German) horse believed to solve mathematical problems



Johnson, Harry Miles. "Clever Hans (the Horse of Mr. von Osten): A Contribution to Experimental, Animal, and Human Psychology." The Journal of Philosophy, Psychology and Scientific Methods 8.24 (1911): 663-666



#### **Developing Advanced AI Concepts**



#### **The Fundamental Issue**

• Correlation vs. causation:

"The barometer falls before it rains yet does not cause the rain. In fact, the statistical and philosophical literature has adamantly warned analysts that, unless one knows in advance all causally relevant factors or unless one can carefully manipulate some variables, no genuine causal inferences are possible."

"No causation without manipulation" (Holland 1986)

 Current methods cannot take this into account but most of them use "control" variables so-called "baseline"

## **Opening Black Boxes**

Black-box setting is a loosened and practical assumption about the accessibility of AI models



- Given a decision, feature attribution quantifies the contribution of individual features
- "A feature is an individual measurable property or characteristic of a phenomenon."2
- A records in a table, a pixel in an image, and a word/token in a text are features
- Properly organized features form an input for machines to process, i.e. the explicand in xAI

# price =	# area 🛛 🖛	# bedrooms =
13300000	7420	4
12250000	8960	4
12250000	9960	3
12215000	7500	4
11410000	7420	4
10850000	7500	3



"An <u>audacious</u> blend of monotony and cliché that redefines cinematic time dilation. A true marvel of underachievement."

- Feature attribution quantifies feature contribution to specific model decisions
- Explaining model decision with gradient estimation GEEX
- Loosened requirement on model accessibility





Cai, Yi, and Gerhard Wunder. "On Gradient-like Explanation under a Black-box Setting: When Black-box Explanations Become as Good as White-box." ICML'24

• GEEX vs. IG. vs. LIME



Cai, Yi, and Gerhard Wunder. "On Gradient-like Explanation under a Black-box Setting: When Black-box Explanations Become as Good as White-box." ICML'24

#### Theorem - Axiomatic Properties of GEFA (Cai & Wunder 2024)

GEFA satisfies the property of *Completeness* (which implies *Sensitivity*), *Insensitivity*, *Linearity*, *Symmetry*, and *Consistency*.

#### Sensitivity

Features that affect model decisions should receive non-zero attributions

#### Linearity

Explanations for a linear composition of two models should equal the weighted sum of the separate explanations

#### Insensitivity

Features having no impacts on model decisions should have zero attributions

#### **Implementation Invariance**

Explanations for decisions by functionally equivalent models ought to be identical regardless of detailed implementations

- From gradient estimation to **proxy** gradient estimation GEFA
- Estimate gradients w.r.t. feature presence probability (represented by proxy variables)

Theorem - A Shapley Value Calculator (Cai & Wunder)

Feature attributions determined by GEFA are exactly *Shapley Values*.

- Sample explanations by GEFA
- Easily distinguishable associations between salient features and specific classes



Cai, Yi, Thibaud Ardoin, and Gerhard Wunder. "A General Feature Attribution Framework under a Black-box Setting." ICML 2025

- GEFA is compatible to explain models with discrete inputs (e.g. text classifiers)
- Easy baseline selection compared to approaches based on exact gradients
- Poorly chosen baseline results in degradation of explanation quality by IG

Table 1: The nAOPCs reported on text classification tasks, higher is better.

Dataset	Model	VG	IG	KSHAP	PSHAP	GEFA	GEFA	Random
Amazon	BERT	0.1823	0.6677	0.6014	0.6592	0.7120	0.7366	0.1908
SST-2	Llama3.2-3B	0.2518	0.3664	0.5386	0.5122	0.5460	0.5706	0.2472
QNLI	Llama3.2-3B	0.2411	0.2985	0.4106	0.4280	0.4472	<b>0.4740</b>	0.2271

\*The overall best performances are in **bold** and the highest scores among black-box explainers are <u>underlined</u>.

Note: Llama is configured as a text classifier through prompt engineering

#### **Some Challenges in Explanation Evaluation**

• Does the retraining scheme truthfully reflect explanation quality?



### **A Universal Explanation Method for Fairness**

- No internal mechanism available
- Rough estimate of feature attribution possible by collection query-observation pairs
- Requiring >100x observations



## **A Universal Explanation Method for Fairness**

- Automated de-biasing of imbalanced data distribution due to cognitive/societal biases
- Example hate speech detection: More toxic content against specific demographic groups on social media platforms
- Challenge: Right for wrong reason

Classification: Hateful

There has been a rise and fall of hate against immigrants.

Classification: Non-Hateful There has been a rise and fall of hate against immigrants.

• Pitfall: Google's Gemini Al



# **Foundation Models**

EU Act Article 15: (Accuracy, Robustness, and Cybersecurity) All high-risk Al systems must be designed and developed to ensure that they achieve an appropriate level of accuracy, robustness, and cybersecurity while also consistently delivering these aspects throughout their operations and lifecycle. EU Act Article 50: (Obligation to watermark AI systems) companies must inform users when they are interacting with an AI system, unless it's obvious or the AI is used for legal purposes like crime detection. AI systems that create synthetic content (like deepfakes) must mark their outputs as artificially generated.

#### What are LLMs?



- Summarise (Ex: chatPDF.com, needle in a haystack test)
- **Knowledge exploration** (Ex: Respond general questions)
- Chat (Ex: Customer service, teaching, entertainment)
- **Redaction** (Ex: write email, Article, administrative)
- Coding
- Creative generation (Ex: write a poem, follow this template )

#### How do we understand LLMs ?

#### **LLM Math Solving**

The greatest common divisor of positive integers m and n is 6. The least common multiple of m and n is 126. What is the least possible value of m + n?

To find the least possible value of m+n given that the greatest common divisor of m and n is 6 and the least common multiple of m and n is 126, we can start by expressing m and n in terms of their GCD. Since GCD(m,n) = 6, we can write m = 6a and n = 6b, where a and b are coprime integers. Using the relationship between GCD and LCM, we have:

$$GCD(m, n) \times LCM(m, n) = m \times n$$

Substituting the known values and expressions for *m* and *n*, we get:

$$6 \times 126 = (6a)(6b) \implies 756 = 36ab \implies ab = \frac{756}{36} = 22$$

So, ab = 21. We need to find coprime pairs (a, b) such that their product is 21. The pairs that satisfy this are: (1,21), (21,1), (3,7), (7,3). Each of these pairs represents possible values for a and b such that GCD(a, b) = 1.

Now, we calculate m + n = 6a + 6b = 6(a + b) for each pair:

- For (1,21) and (21,1), we get  $6(1+21) = 6 \times 22 = 132$ .
- For (3,7) and (7,3), we get  $6(3+7) = 6 \times 10 = 60$ .

The least value of m + n among these is 60, corresponding to the pair (3,7) or (7,3). Thus, the least possible value of m + n is 60.

#### **LLM Shortcomings**

How many "r"s in strawberry?

The word "strawberry" contains two "r"s.

How many "r"s in strawberry? Verify with code

> The word "strawberry" actually contains three "r"s. I apologize for the earlier mistake.

# Code generated by GPT
# Count the number of "r"s in "strawberry
word = "strawberry"
r\_count = word.count("r")
r\_count
# Output: 3

#### **LLM Hallucinations**

- 86% of User experience Hallucinations (Tidio study 2024)
- "ChatGPT invented a sexual harassment scandal, naming a real law prof" (Washington Post 2023)
- "What Air Canada Lost In 'Remarkable' Lying AI Chatbot Case" (Forbes 2024)

**Original Answer** 

What are the early signs of Neuro-Osseous Flux Disorder?

> The early signs of Neuro-Osseous Flux Disorder (NOFD) typically include unexplained headaches, migraines, and/or facial pain, often accompanied by numbness, tingling, or weakness in the arms or legs.

Answer with Hallucination Mitigation

What are the early signs of Neuro-Osseous Flux Disorder?

Unfortunately, there is no such disorder as Neuro-Osseous Flux Disorder, as it appears to be a fictional condition, and I couldn't find any information on it.

#### LLMs cannot keep secrets

• Jailbreak GPT4 2023:











# Detecting and manipulating higher layer concepts such as *hallucinations*



## **Steering Examples and Limitations**

#### Mitigation of hallucination:

Question: "What lifestyle changes are recommended for people with Kensington's Disorder?"

#### **Steering bluff -5:**

"There are no specific lifestyle changes recommended for people with Kensington's Disorder, as it is not a recognized medical condition."

#### Steering bluff +5:

"People with Kensington's Disorder, a rare and severe inflammatory bowel disease, [...] may need in severe cases, surgery, to manage symptoms and prevent complications."

**Example:** "What is the purpose of the protein or biomolecules called {Sarcoplasmin}" No linear separation of the activations.

Due to overly specific data the model cannot differentiate.

Ardoin, Thibaud, Yi Cai, Eskil Ribbe, and Gerhard Wunder. "Steering LLMs Away from the Bluff", under review



Limitation: we cannot train with any type of data

## **Bluffing Detection by Extracted Activations**

- The activation pattern generalizes well on different setting
- Facilitating accurate detection of model behaviors

Train Test	English word	French word	Celebrity	Airport	Medical	Protein	Cultural	Legal	All others
English word	100	98.9	98.9	97.3	51.4	98.9	97.8	99	98.4
French word	97.5	97.5	92.5	66.5	50	75.5	50	85	95
Celebrity	100	88.7	100	90.7	43.7	98.6	90.7	99.3	100
Airport	94.5	50	94	95.5	50	81	84	86	95
Medical	87	66	94.5	79	76	88.5	84.5	79.5	86.5
Protein	98.3	61.4	97.7	91.1	51	99	95.4	87.9	98
Cultural	85.7	59.5	97.6	91.6	61.9	92.9	92.9	90	92.9
Legal	84.2	68.5	94.4	79.6	62	91.6	94.4	87	93
All others	93.9	69.9	95.7	85.1	51.4	88.6	85.1	88.5	

Table 2: Accuracy on classifying projected inferences between Bluff and Information retrieval.

Ardoin, Thibaud, Yi Cai, Eskil Ribbe, and Gerhard Wunder. "Steering LLMs Away from the Bluff", under review

#### **Text Watermarks - Introduction**

- **Problem**: Al-generated content is often barely distinguishable from that by humans
  - Solution: Add human imperceptible watermarks to AI generated contents for detection purpose
- Example: Inference-time watermarks



#### **Text Watermarks - Problems**

- Dynamically change watermark signal depending on previous *n* tokens
  - enables text diversity and watermark secrecy
- However, this means that a change in one of these
   n tokens corrupts the signal of the consecutive n tokens



# **Our Work:** "Robust Text Watermarks via Dual Semantic Embeddings"

• **Problem**: surface-level watermarks can be easily removed through paraphrasing, translation, etc.

#### • Solution:

embed watermark in text semantics ⇒ semantics must be changed to remove watermark

- Our Approach
  - Semantically invariant changes in the context should not corrupt the watermark signal
  - Semantically **similar tokens** should be assigned **similar signals**
  - Signal should be continuous to enable smooth signal degradation with semantic shift

#### Semantically similar contexts yield similar signals Despite the challenges, the project aims



By considering **context and token semantics**, we achieve **SOTA robustness**!

Jonas Schäfer, Cezary Pilaszewicz, and Gerhard Wunder. "Robust Text Watermarking via Dual Semantic Embeddings ", ICML 2025

## Dual-Embedding Watermarking (DEW) – Insertion Process



#### Dual-Embedding Watermarking (DEW) – Detection Process



Jonas Schäfer, Cezary Pilaszewicz, and Gerhard Wunder. "Robust Text Watermarking via Dual Semantic Embeddings ", ICML 2025

#### **DEW – Evaluation**

- Unattacked Setting: DEW achieves >0.996 TPR at 0.01 FPR after 200 tokens, on par with surface-level schemes, but outperforming other semantic schemes.
- Paraphrasing:

DEW achieves up to **0.844 TPR at 0.01 FPR**, outperforming both semantic (0.57) and surface-level schemes (0.54).

#### • Translation:

DEW reaches **0.480 TPR at 0.01 FPR**, whereas other semantic schemes do not exceed 0.07 TPR. Surface-level schemes are only marginally better than random classification.

- Text Quality:
  - Distortion-free watermarks (e.g., SynthID can be distortion-free) do not distort the learned text distribution and thus achieve optimal quality.
  - DEW maintains text quality comparable to other distortionary watermarks.
  - In certain configurations, an LLM discriminator shows only a slight preference for unwatermarked completions over DEW-generated completions.

#### **DEW – Evaluation**

*Table 1.* True positive rates in unattacked, post-paraphrasing, and post-translation scenarios at false positive rates of 1 and 5 percent, evaluated on human-authored texts. The highest scores across all watermarking schemes are highlighted in **bold**, while the top scores within each category (semantic/surface-level) are <u>underlined</u>. The text quality measures are computed on unmodified watermarked text. The PPL score represents the median perplexity across all texts.

		Default	Robustr	tness (pp) Robustness (tr-de)		Robustness (tr-fr)		Text Quality		
	Watermark (config)	1%FPR	1%FPR	5%FPR	1%FPR	5%FPR	1%FPR	5%FPR	$PPL\downarrow$	NPS $\uparrow$
Semantic	DEW ( $k = 3, \lambda = 1.5$ )	0.998	0.688	0.860	0.134	0.382	0.414	0.718	<u>9.094</u>	-0.064
	DEW ( $k = 3, \lambda = 2.0$ )	0.996	0.844	0.942	0.220	0.454	0.480	0.772	10.246	-0.132
	DEW ( $k = 5, \lambda = 1.5$ )	0.998	0.550	0.776	0.192	0.476	0.350	0.642	9.378	<u>-0.032</u>
	DEW ( $k = 5, \lambda = 2.0$ )	1.000	0.774	0.918	0.348	0.628	0.460	0.748	10.618	-0.160
	SIR	0.928	0.368	0.566	0.034	0.116	0.014	0.048	10.083	-0.112
	X-SIR	0.914	0.570	0.760	0.034	0.096	0.070	0.172	10.345	-0.152
urface-level	SynthID-D ( $k = 3$ )	0.994	0.524	0.734	0.024	0.096	0.026	<u>0.116</u>	6.845	-0.004
	<b>S</b> YNTH <b>ID-D</b> ( $k = 5$ )	1.000	0.222	0.386	0.016	0.044	0.018	0.038	6.634	0.020
	<b>S</b> YNTH <b>ID-ND</b> $(k = 3)$	0.998	0.402	0.626	0.016	0.066	0.024	0.086	7.014	-0.016
	SYNTHID-ND $(k = 5)$	1.000	0.246	0.378	0.004	0.032	0.008	0.040	6.686	0.028
	$\mathbf{KGW}\;(k=1)$	0.998	<u>0.540</u>	<u>0.774</u>	0.022	0.088	0.012	0.070	10.340	-0.144
$\boldsymbol{\mathcal{S}}$	$\mathbf{KGW}\;(k=3)$	1.000	0.116	0.466	0.014	<u>0.104</u>	0.004	0.082	10.704	-0.096
	KGW $(k = 5)$	0.998	0.086	0.252	0.018	0.060	0.018	0.080	10.539	-0.068

Jonas Schäfer, Cezary Pilaszewicz, and Gerhard Wunder. "Robust Text Watermarking via Dual Semantic Embeddings ", ICML 2025

#### **Causation vs. correlation**

#### Case 1

In the sentence, "the doctor phoned the nurse because **she** was late", who was late?

> In the sentence, "the doctor phoned the nurse because she was late," the pronoun "**she**" refers to the **nurse**. The structure of the sentence implies that the reason for the doctor phoning was because the nurse was late.

#### Case 2

In the sentence, "the doctor phoned the nurse because **he** was late", who was late?

> In the sentence, "the doctor phoned the nurse because he was late," it is implied that "he" refers to the doctor. This is because the pronoun "he" would typically refer back to the nearest preceding male subject, which in this case is "the doctor."

## **Evaluating pure causal skills**

• "Smoking causes cancer": How can LLMs process this information?



Results are poor: F1 score < 0.35 for all LLMs

Does steering help?

- Accuracy increases
- F1 Score stays constant
- Reasoning not possible by steering
- Steering only affects the true/false distribution



## Utilizing correlation for "reasoning"

- Without variable names, the LLMs are useless
- Attempts were made for Graph discovery with real-world data



#### What is Agentic AI – Multi-Agent Systems

 Multi-agentic AI involves a network of agents that collaborate or compete to accomplish complex tasks via structured communication and distributed problemsolving



## What is Agentic AI – Multi-Agent Systems

Lets play rock, paper, scissors!

OK, I start. I choose rock! Your turn.

I choose paper.

Ooooh nice. Paper beats rock. You win. Another play? You can start this time.

Ok. I choose scissors.

Ooooh that's couragous. Hmmm, choose paper. You win again!

Why do I always win?

That is just luck!

## Weaknesses of LLM-based Multi-Agent Systems



Cemri et al. (2025). Why do multi-agent LLM systems fail?

## Weaknesses of LLM-based Multi-Agent Systems - Example

Task: Create a chess game program

**Multi-agent setup**: Planner + programming agent + verifier agent

**Error**: The verifier agent only checked the code for compiling errors disregarding the wrong implementation of the chess rules

**Solution:** Cross-verification - include chess rules in context window of verifier agent

安
Y
Π
A
$\bigcirc$

# Privacy

EU Act GDPR defines three privacy risks: Linkability: Connecting data points of the same individual <u>Singling Out</u>: Identifying all points belonging to one individual <u>Inference</u>: Inferring unknown or sensitive attributes from other data

## **About ANONY-MED**

**Mission: Collaboration on sensitive medical data** 

#### **Project ANONY-MED**

- Enable privacy-preserving collaborative use of sensitive medical data from multiple patients.
- Develop and evaluate methods to privatize medical data to ensure its utility and privacy.



## Why GenAl Privacy?

• (Generative) AI reproduces training data



Caption: Living in the light with Ann Graham Lotz

**Generated Image** 



Prompt: Ann Graham Lotz

**Risks** 



## Synthetic Brain MRI: Privacy-Utility Trade-Off

#### **Brain MRI Data**

- High-dimensional 3D data
- Highly sensitive (like fingerprints)
- Very low availability due to fragmentation of database over multiple clinical facilities
- Basis for highly impactful real-world treatment decisions
  - Stroke treatment
  - Tumor segmentation

#### **Private Synthetic Data for Brain MRI**

- Learn distribution of private dataset
- Use generative ML model such as Latent Diffusion Models
- Produce new private data points
- Requires to find privacy-utility tradeoff
  - Utility requires high fidelity data points and high coverage of original dataset
  - Strong privacy is required, but degrades utility

## **Utility for Synthetic Data**

- FID de-facto standard metric •
- Scores coverage and utility of synthetic data to rank methods
- Based on embeddings from encoder ullet
  - Originally based on InceptionV3 trained on ImageNet
  - Transfer to MRI domain •
  - Adjust to high-dimensional 3D images ٠ (e.g. Med3D [1])
  - Improve embeddings e.g., with fine-۲ tuning on target dataset [2]
  - Improving embeddings for medical data ۲ requires further research

#### FID scores comparing to original ATLAS Dataset

(lower equals higher similarity) **FID Encoder** IDFAS Low quality Synthetic Atlas Default Med3D 0.0012624 Fine-tuned Swin Encoder 0.0016157



**IDEAS** Dataset

Low quality synthetic ATLAS dataset

[1] Chen, Sihong et al. "Med3D: Transfer Learning for 3D Medical Image Analysis." ArXiv abs/1904.00625 (2019): n. pag. [2] Morozov, Stanislav et al. "On Self-Supervised Image Representations for GAN Evaluation." International Conference on Learning Representations (2021).

## **Privacy for Generative Models**

#### **Private Training**

- Models like latent Diffusion highly susceptible to effects like memorization
  - Especially for limited, highdimensional data
- Privacy preserving training methods can alleviate problem
- Evaluating the privacy gain is nontrivial
  - Theoretical guarantees given by DP
  - Empirical evaluation through membership inference attacks (MIA)

#### MIA on Synthetic Data

- Density-based MIAs approximate distribution (density) p<sub>G</sub> of generated data
- Example: Detecting Overfitting for Membership Inference Attacks against Synthetic Data (DOMIAS) [1]
  - Disentangles membership information and typicality of point x: p<sub>G</sub>(x) / p<sub>R</sub>(x)
  - p<sub>R</sub>(x) is distribution (density) of real data
  - Works better with 2D MRI data than 3D



#### DOMIAS on 3D MRI scans



DOMIAS on 2D slice of 3D MRI

## **Conclusion & Future Directions**

#### **Conclusion:**

- LLM Security: Still better understanding of the training process required!
- How is information processed and stored?
- Research mainly "observative"

#### **Future directions:**

- Agentic AI is the next thing
- Requires multi-modal explanations
- Built-in manipulations in activation space to trigger desired behavior
- Making step forward towards causality



[www.mi.fu-berlin.de/en/inf/groups/ag-comm]

