

# ChatGPT war gestern Vertrauenswürdige KI ist die Zukunft

KNOW-CENTER GMBH  
RESEARCH CENTER FOR DATA-DRIVEN BUSINESS  
& BIG DATA ANALYTICS, GRAZ

---

Hermann Stern, Roman Kern

AI & Medical Software in Healthcare  
Regulatory Konferenz für Medizinprodukte und In-vitro Diagnostika  
Wien, 17.10.2023



**DI HERMANN STERN**

Business Area Manager – Digital  
Transformation Design

[hstern@know-center.at](mailto:hstern@know-center.at)

+43 (0)664 887 83 114

<https://www.linkedin.com/in/hermann-stern>

## Know-Center GmbH



K1 COMET-Zentrum



Gegründet 2001

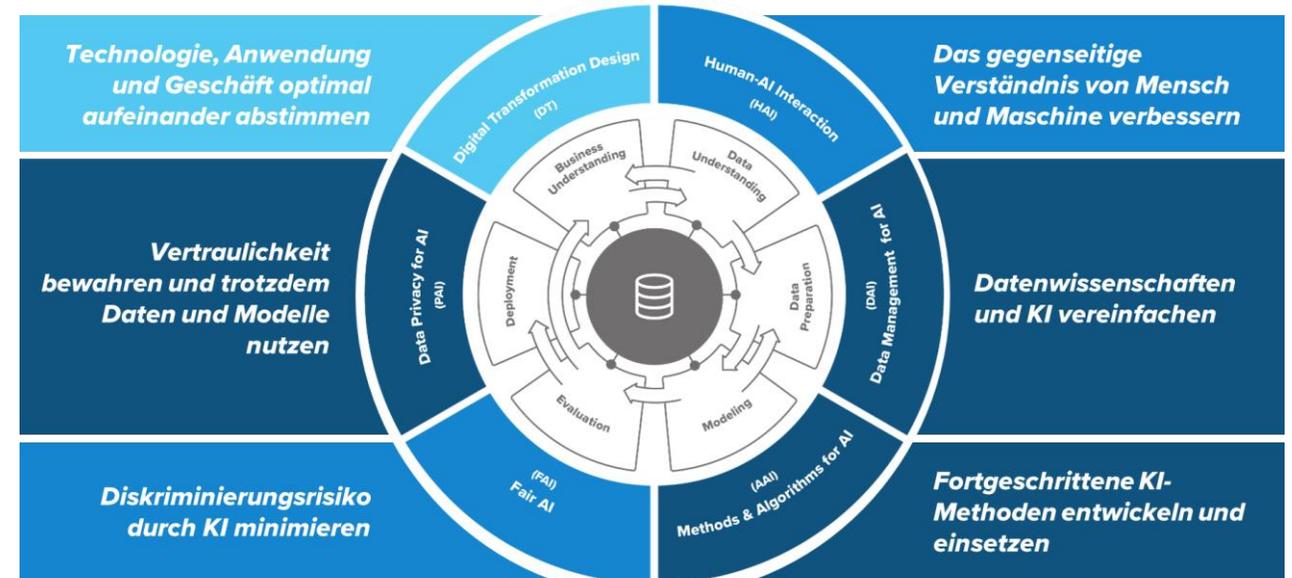


140+ Mitarbeiter:innen



TU Graz, Campus Inffeldgasse

„Wir forschen, entwickeln und beraten entlang der Datenwertschöpfungskette zum Thema *vertrauenswürdige KI* und *Data Science*.“



# DEFINITION VON KÜNSTLICHER INTELLIGENZ



# DEFINITION VON KÜNSTLICHER INTELLIGENZ

*[...] Künstliche Intelligenz ist die Fähigkeit einer Maschine, menschliche Fähigkeiten wie logisches Denken, Lernen, Planen und Kreativität zu imitieren.*

*KI ermöglicht es technischen Systemen, ihre Umwelt wahrzunehmen, mit dem Wahrgenommenen umzugehen und Probleme zu lösen, um ein bestimmtes Ziel zu erreichen. [...]*

"Eine technische Lösung, die Menschen (in ihrer Arbeit) unterstützt"

# DEFINITION VON KÜNSTLICHER INTELLIGENZ

**# KI ist  
gekommen, um  
zu bleiben!**

## ChatGPT – die KI funktioniert (plötzlich)!? - Bericht aus der Wissenschaft

### ChatGPT Sprints to One Million Users

Time it took for selected online services to reach one million users



\* one million backers \*\* one million nights booked \*\*\* one million downloads  
Source: Company announcements via Business Insider/LinkedIn



- Fortschritte im Bereich **Machine Learning** und **Natural Language Processing („NLP“)**

- Immer mehr Aufgaben kann die Maschine **auf dem Niveau von Menschen lösen**

- ChatGPT hat sich angekündigt

# DEFINITION VON GENERATIVER KÜNSTLICHER INTELLIGENZ

*"Generative KI ist eine Art von künstlicher Intelligenz, die auf der Grundlage von Eingabedaten und Algorithmen **neue Outputs** wie Bilder, Musik oder Text erzeugt. Das Ziel der generativen KI ist es, Outputs zu erzeugen, **die realen Beispielen der gleichen Kategorie ähneln.**"*

(ChatGPT, 04.02.2023)

IST DIE (GENERATIVE)  
KI REIF FÜR DEN  
UNEINGESCHRÄNKTEN  
PRAKTISCHEN EINSATZ?



# IST DIE (GENERATIVE) KI REIF FÜR DEN UNEINGESCHRÄNKTEN PRAKTISCHEN EINSATZ?

**"Nein, die KI ist nicht vertrauenswürdig!"**

- Durchaus gemeint: das Vertrauen der Menschen (Benutzer:innen), vielleicht aber auch über Wertschöpfungsketten und Unternehmensgrenzen hinweg?
- Sind die Ergebnisse denn korrekt? fair? ethisch? sicher? transparent?
- Rechtlich einwandfrei verwendbar?
- Am Weg: Zertifiziert?

## Dialog mit GPT-3

---

- Was ist die Nationalhymne der Schweiz?

## Dialog mit GPT-3

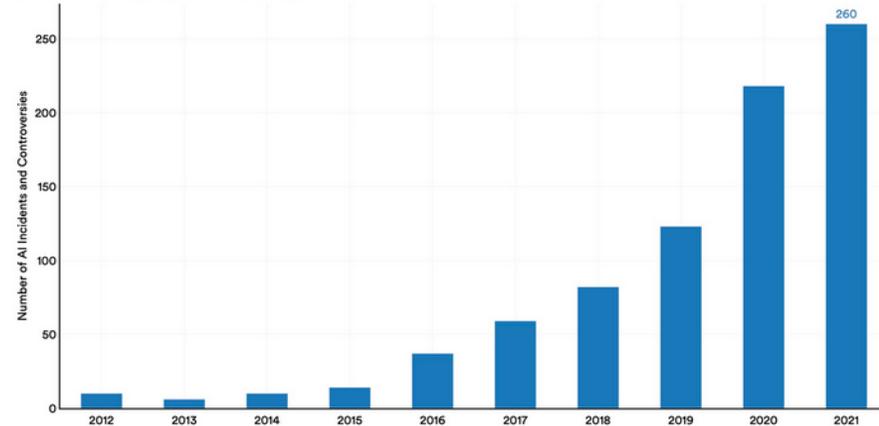
- Was ist die Nationalhymne der Schweiz?
- "Schweiz, Schweiz über alles!"
- ChatGPT "Halluziniert!"



# Auswirkungen durch den Einsatz der KI

**Number of AI Incidents and Controversies, 2012–21**

Source: AIAAIC Repository, 2022 | Chart: 2023 AI Index Report



According to the AIAAIC database, which tracks incidents related to the ethical misuse of AI, the number of AI incidents and controversies has increased 26 times since 2012. Some notable incidents in 2022 included a deepfake video of Ukrainian President Volodymyr Zelenskyy surrendering and U.S. prisons using call-monitoring technology on their inmates. This growth is evidence of both greater use of AI technologies and awareness of misuse possibilities.

**AIID** | AI INCIDENT DATABASE | English | [Social Media Icons] | [Subscribe]

Discover | Submit

Welcome to the AIID

- Discover Incidents
- Spatial View
- Table View
- Entities
- Taxonomies
- Word Counts
- Submit Incident Reports
- Submission Leaderboard
- Blog
- Subscribe

### Incident 168: Collaborative Filtering Prone to Popularity Bias, Resulting in Overrepresentation of Popular Items in the Recommendation Outputs

[Social Media Icons]

**Description:** Collaborative filtering prone to popularity bias, resulting in overrepresentation of popular items in the recommendation outputs.

**Tools:** [Notify Me of Updates] [New Report] [New Response] [Discover] [Citation Info]

**Entities:** [View all entities](#)

Alleged: [Facebook](#), [LinkedIn](#), [YouTube](#), [Twitter](#) and [Netflix](#) developed and deployed an AI system, which harmed [Facebook users](#), [LinkedIn users](#), [YouTube users](#), [Twitter Users](#) and [Netflix users](#).

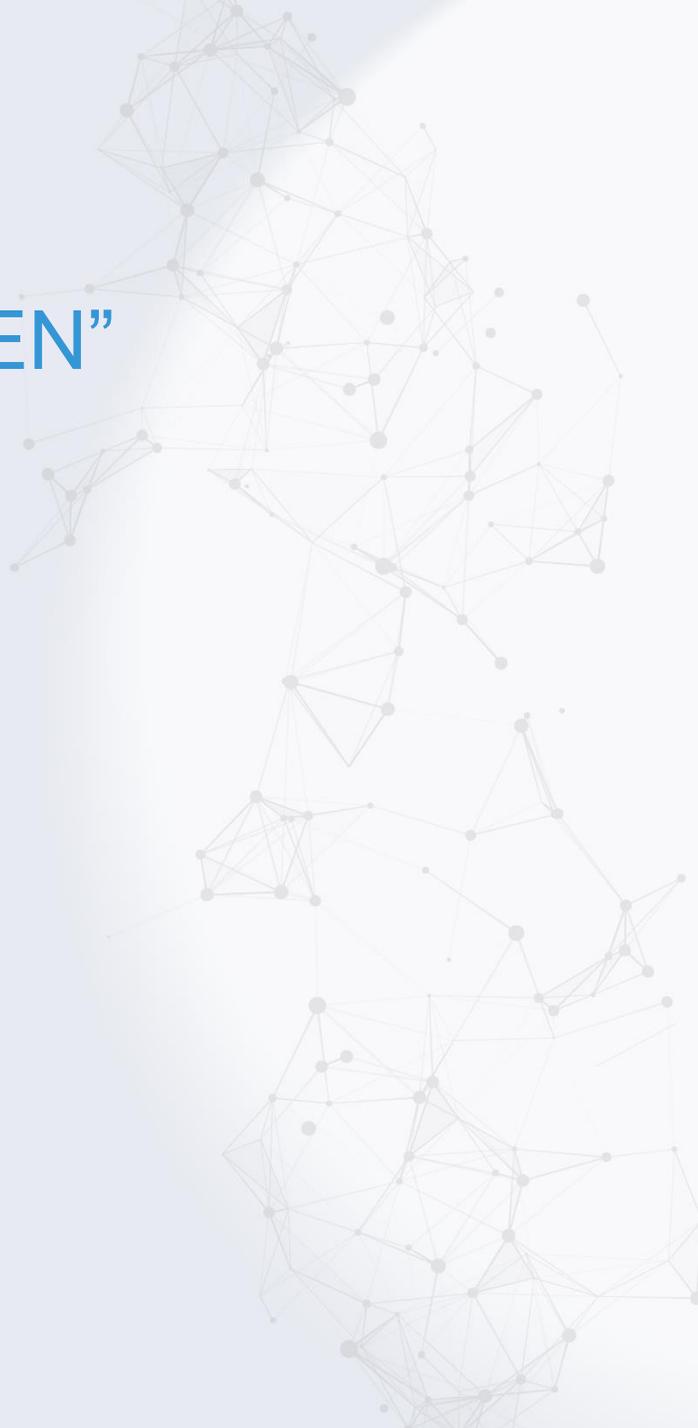
**Incident Stats**

Incident ID	168
Report Count	2
Incident Date	2022-03-01
Editors	Sean McGregor, Khoa Lam

**Incident Reports**

<https://incidentdatabase.ai>

# WIE NUN “VERTRAUEN” SCHAFFEN?



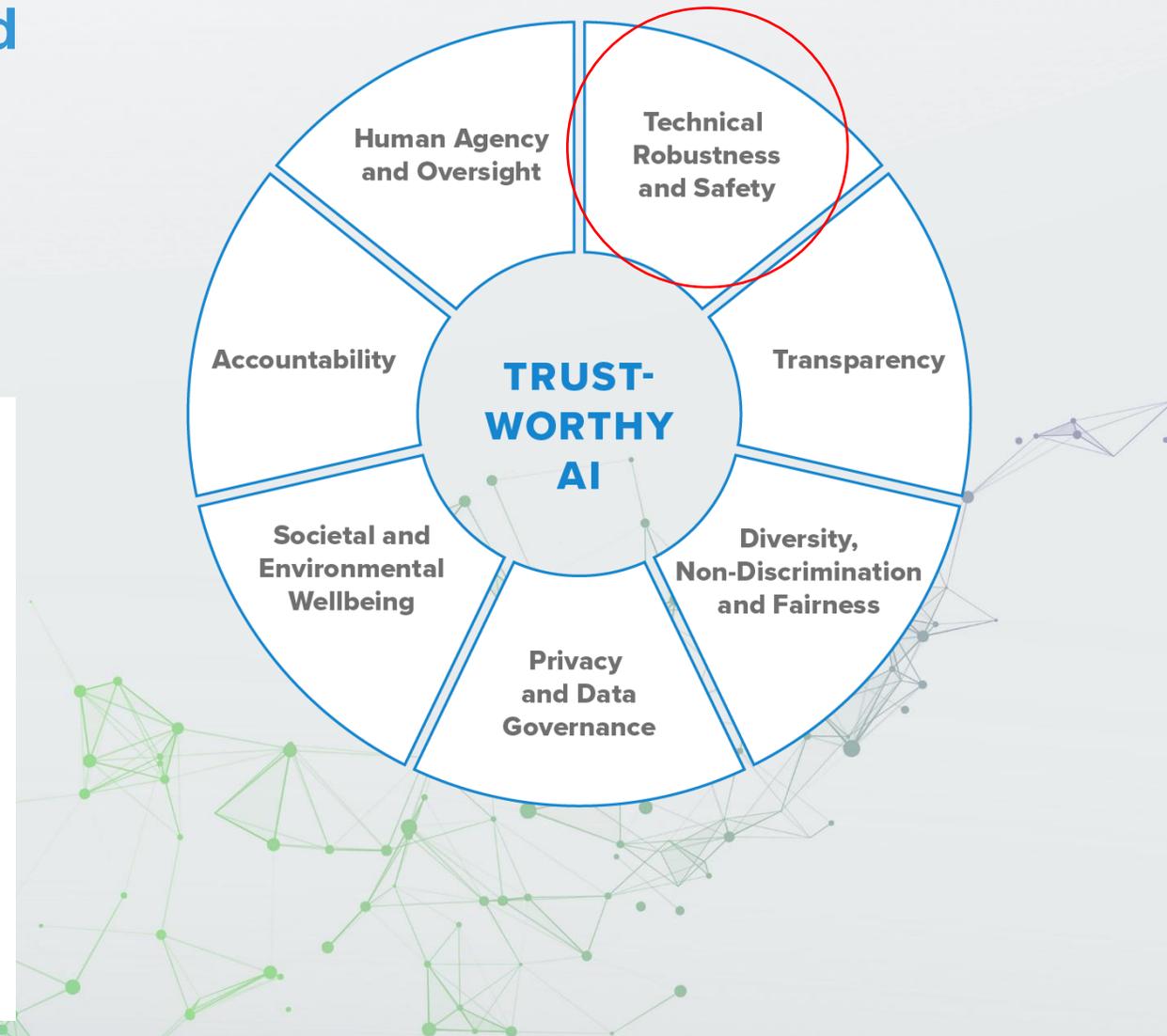
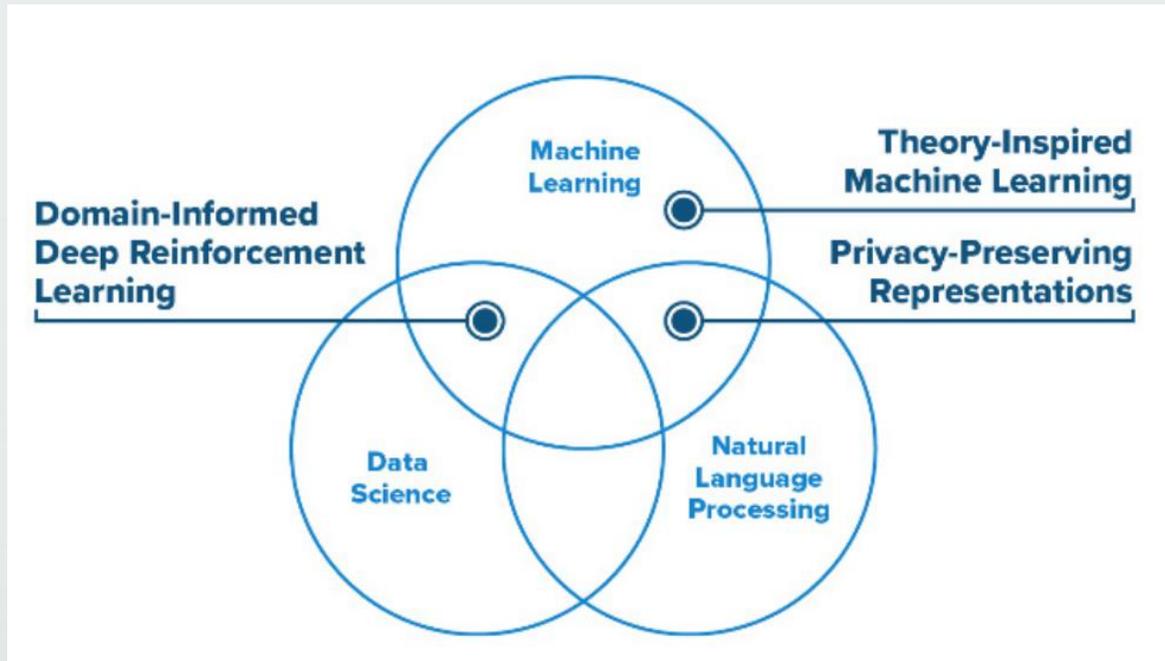
# EU INITIATIVEN

The screenshot shows the top of a news article from the European Parliament. The header includes the EP logo, the word "News", and a search bar. Below the header is a navigation menu with "Headlines", "Press room", "Agenda", "FAQ", and "Election Press Kit". The main content area features the title "AI Act: a step closer to the first rules on Artificial Intelligence" and a sub-header "Press Releases | MCD | LBE | 11-05-2023 - 09:34". On the left, there are social media icons for Facebook, Twitter, LinkedIn, and WhatsApp. The main text begins with "Once approved, they will be the world's first rules on Artificial Intelligence". On the right, there is a "Further information" section with a red "Compromise text" icon and links to "Draft reports, amendments tabled in committee", "European Parliamentary Research Service: Artificial Intelligence", and "Legislative train: the AI Act".



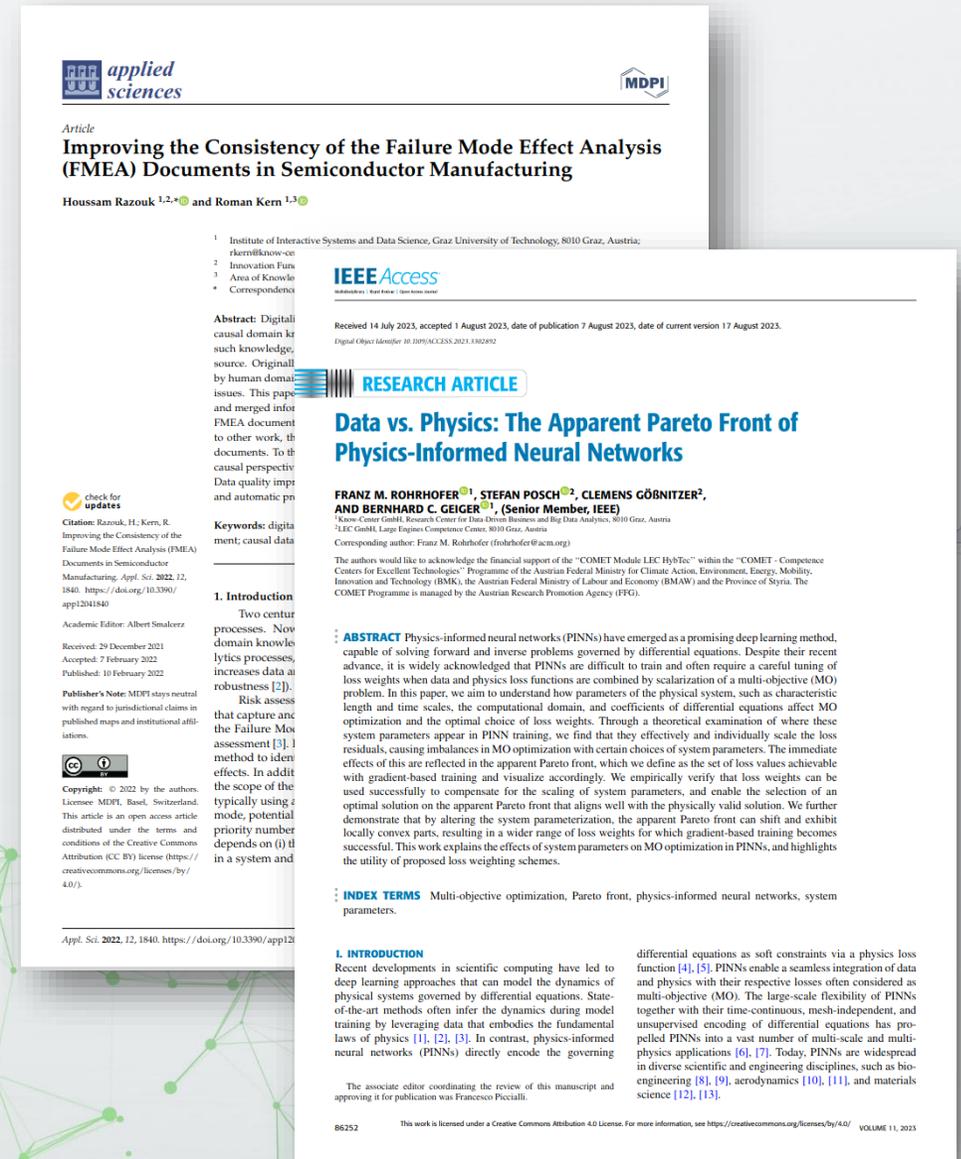
# KI „korrekt“ machen - Hallucination und Strategien in der Forschung

- **Grundlagenforschung**
  - Der Maschine das Denken beibringen



# Hallucination und Strategien in der Forschung

- **Angewandte Forschung**
  - **Vorhandenes Wissen nutzen**
    - ... beispielsweise **FMEA** und **technische Reports**
      - via Knowledge Graphs
  - ... **Domänenwissen** in die **Modellierung** miteinbeziehen
    - via PINNs (wissenschaftliche Modelle und physikalische Gesetze in neuronale Netzwerke einbeziehen)



# Hallucination und Strategien in der Forschung

- Angewandte Forschung
  - Tool-Augmented LLMs
    - Agents, Plugins, Extensions

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

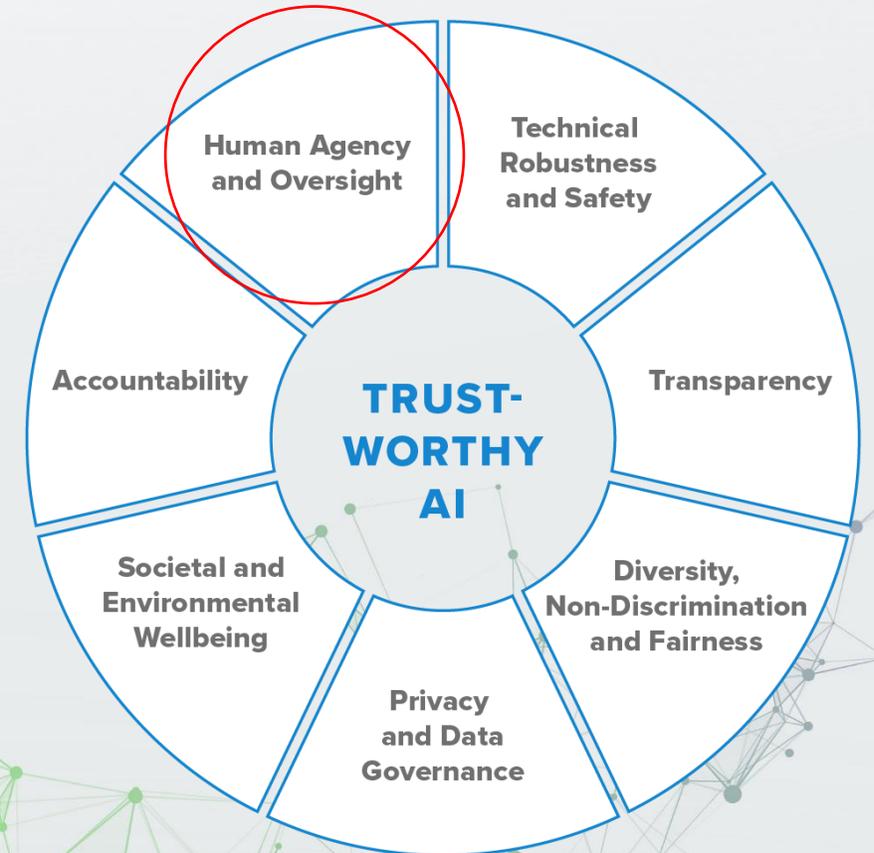
Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

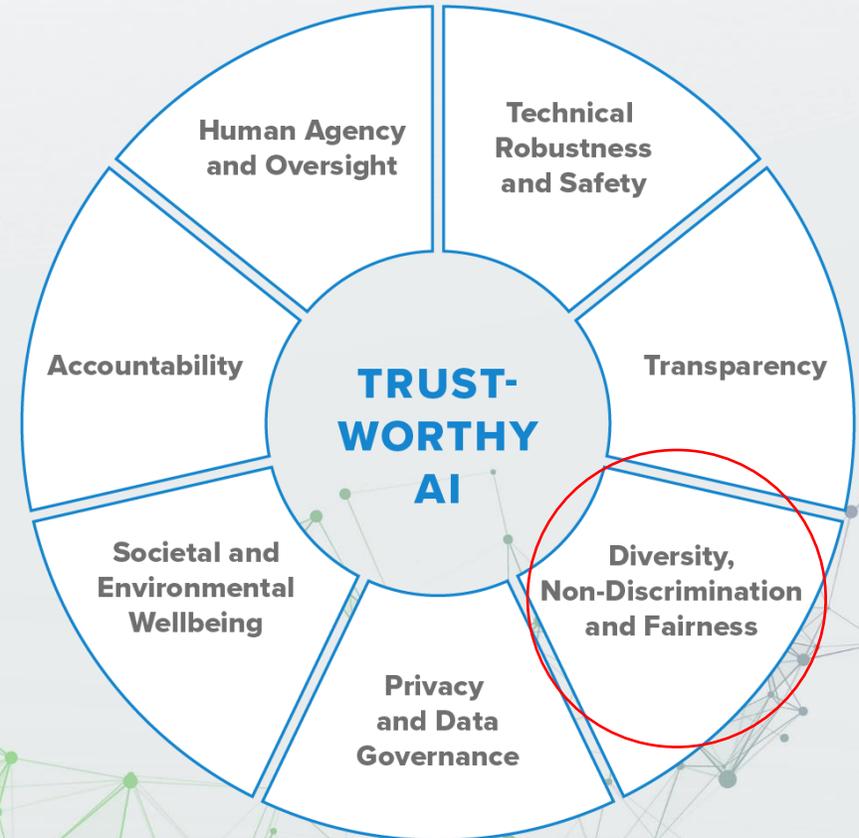
The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

## Wahrnehmung von ChatGPT (KI) ändern

- **In der Gesellschaft**
  - "Terminator"
  - Monopol von Technologieriesen
  - Missinformation, Fake News
- **Im Unternehmen**
  - Angst vor Jobverlust
  - Veränderung des Arbeitsumfelds
  - Unklarheiten der Potentiale



# KI “fair” machen – in Gestaltung, Entwicklung und Betrieb



## KI “sicher” machen

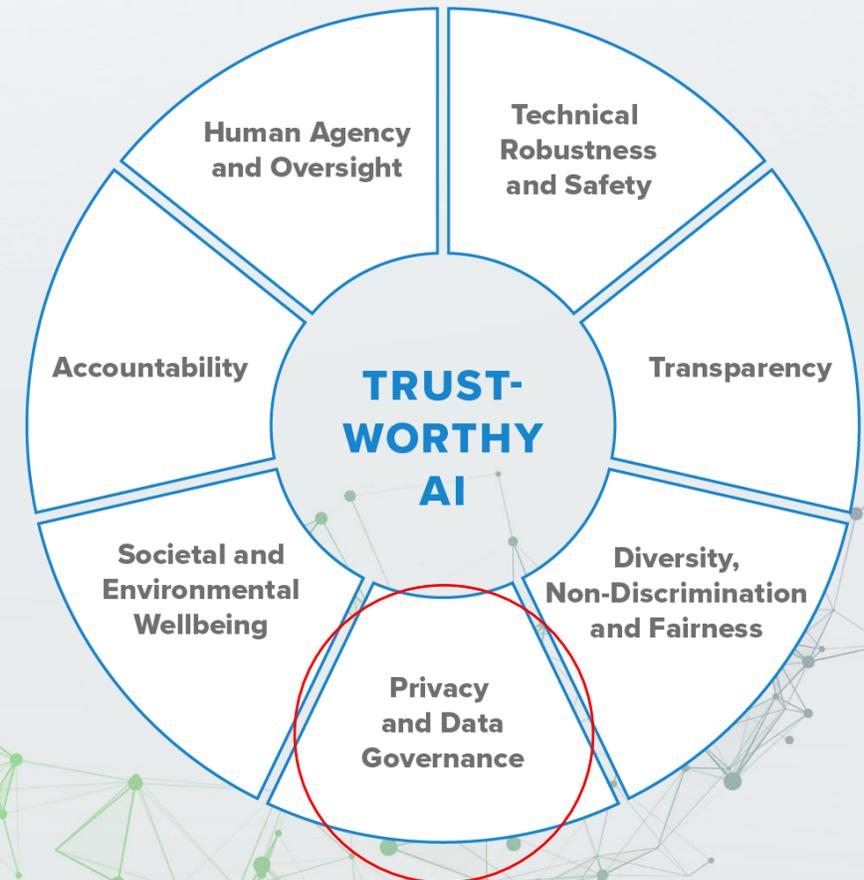
- Austausch von sensiblen und vertraulichen Daten
- Schutz der Privatsphäre
- Ansatz: On-Premises LLMs

Tech [Artificial Intelligence](#)

### Whoops, Samsung workers accidentally leaked trade secrets via ChatGPT

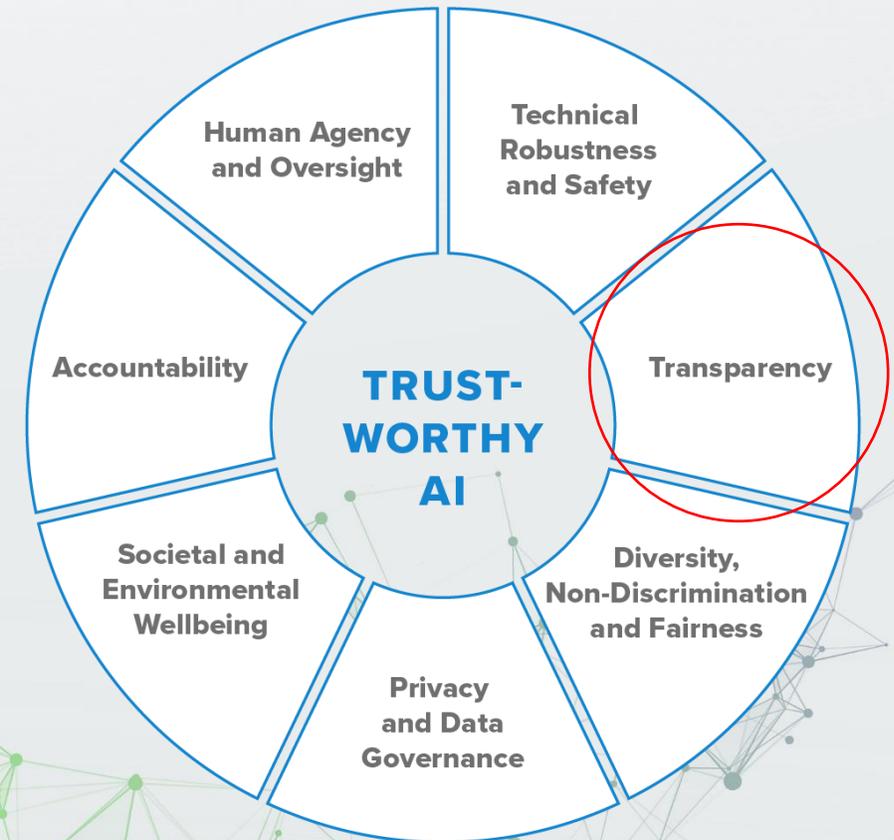
ChatGPT doesn't keep secrets.

By [Cecily Mauran](#) on April 6, 2023



## KI “verständlich” machen

- Schutz KI als **Black Box**
- **eXplainable AI (XAI)** als Forschungsfeld
- Domänenexperten die AI **erklärbar** machen



## Beispiel: simplifAI

# Radiologischer Befund

Max Musterpatient

01.01.2000

### CT-Gehirnschädel:

E-vacuo-Ausweitung der inneren und äußeren Liquorräume ohne Hinweis auf akute Liquorzirkulationsstörung.

Mediane Lage der Mittellinienstrukturen.

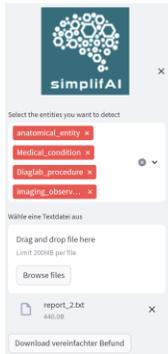
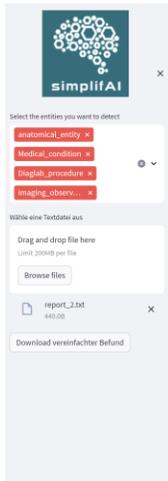
Die basalen Zisternen frei.

Sklerose der Hirnbasisarterien.

Keine Blutung , keine Raumforderungszeichen , kein Nachweis eines rezent demarkierten Territorialinfarkts.

Die partiell miterfassten NNH und das Mastoidzellsystem beidseits regelrecht belüftet.

# Beispiel: simplifAI



## Informationsextraktion

CT **ohne** **Erweiterung** der inneren und äußeren Liquorräume ohne Hinweis auf akute **Liquorzirkulationsstörung**. Die basalen **Strukturen** der **Mittellinie** sind **regelmäßig**. Keine **Blutung** **oder** **keine** Raumforderungszeichen **oder** **keine** **neuro** **oder** **keine** **neuro** eines rezente **epileptischen** **oder** **epileptischen**. Die partiell miterfassten **Strukturen** und das **Struktur** beidseits regelmäßig belüftet.

### Welche (bildgebenden) Untersuchungen wurden durchgeführt

- CT Gehirnschädel

### Wie ist mein Gesundheitszustand

- E-vacuo-Ausweitung inneren äußeren Liquorräume ohne Hinweis auf
- inneren äußeren Liquorräume ohne Hinweis auf akute Liquorzirkulationsstörung
- Sklerose Hirnbasisarterien
- Keine Blutung kein Nachweis rezente
- kein Nachweis rezente demarkierten Territorialisinfarkts

### Welche Pathologien wurden entdeckt

- keine Raumforderungszeichen kein Nachweis rezente

Strukturierter Befund **OH** Vereinfachter Befund

	concept	chr_id	normalized_concept	description
0	CT	H5BMR02H	ct	Die Computertomographie (CT) ist ein
1	Gehirnschädel	P0BCF5X7	gehirnschaedel	Der Begriff "Gehirnschädel" bezieht
2	E-vacuo-Ausweitung	O2LLECEG	e - vacuo - ausweitung	Die E-vacuo-Erweiterung bezieht sich
3	inneren	HJSEJ2ND	inneren	"inneren" bezieht sich auf den innere
4	äußeren	L0U76BXQ	aeuereen	"aeuereen" bezieht sich auf etwas, das
5	Liquorräume	K8BGBFJ	liquor-raeume	Mit Liquor cerebrospinalis gefüllte Ho
6	ohne Hinweis auf	ILNFWKSC	ohne hinweis	[kein] nicht [gegen] (hinweis) Stat,
7	akute	TSJQC2T1	akute	[akut] im Augenblick herrschend, von
8	Liquorzirkulationsstörung	XAS5R6	liquorzirkulationsstoerung	Eine Liquorzirkulationsstörung bezie
9	Mediane	QQVFK25J	mediane	[Mediane] Seitenhalbierende eines Di

# Radiologischer Befund

Max Musterpatient

01.01.2000

## CT-Gehirnschädel :

Die Computertomographie (CT) ist eine Art der Bildgebung, bei der mithilfe von Röntgenstrahlen detaillierte Bilder des Körpers erstellt werden. Im vorliegenden Fall wird der Gehirnschädel untersucht, also der knöcherne Teil des Schädels, der das Gehirn umgibt und schützt.

## E-vacuo-Ausweitung der inneren und äußeren Liquorräume ohne Hinweis auf akute Liquorzirkulationsstörung .

Die E-vacuo-Ausweitung bezieht sich auf eine Erweiterung der inneren und äußeren Liquorräume ohne Hinweis auf eine akute Liquorzirkulationsstörung. Das bedeutet, dass sich die mit Liquor cerebrospinalis gefüllten Hohlräume im Gehirn und Rückenmark ausgedehnt haben, aber es gibt keine Anzeichen für eine akute Störung des normalen Flusses von Liquor cerebrospinalis.

## Mediane Lage der Mittellinienstrukturen .

Die medianen Mittellinienstrukturen beziehen sich auf anatomische Strukturen im zentralen Teil des Gehirns, die entlang der mittleren Linie des Gehirns liegen.

# ZUSAMMENFASSUNG

1. KI funktioniert (in vielen Bereichen)
2. KI muss vertrauenswürdig (Trustworthy) werden
3. KI muss den Menschen dienen

# Vielen Dank!

---

Hermann Stern, Roman Kern

AI & Medical Software in Healthcare  
Regulatory Konferenz für Medizinprodukte und In-vitro Diagnostika  
Wien, 17.10.2023



**DI HERMANN STERN**

Business Area Manager – Digital  
Transformation Design

[hstern@know-center.at](mailto:hstern@know-center.at)

+43 (0)664 887 83 114

<https://www.linkedin.com/in/hermann-stern>