

Sicherheit in der Ära der künstlichen Intelligenz

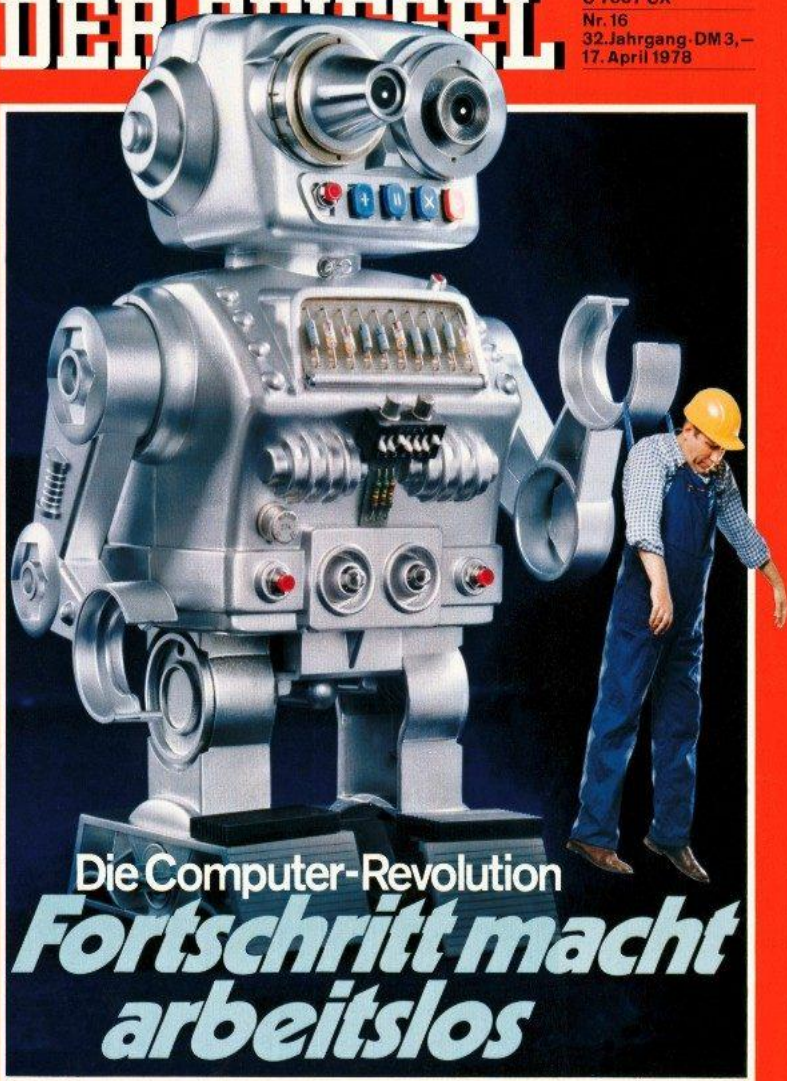
NETFLIX



Moonfall

DER SPIEGEL

C 7007 CX
Nr. 16
32. Jahrgang DM 3,-
17. April 1978



Die Computer-Revolution

Fortschritt macht arbeitslos

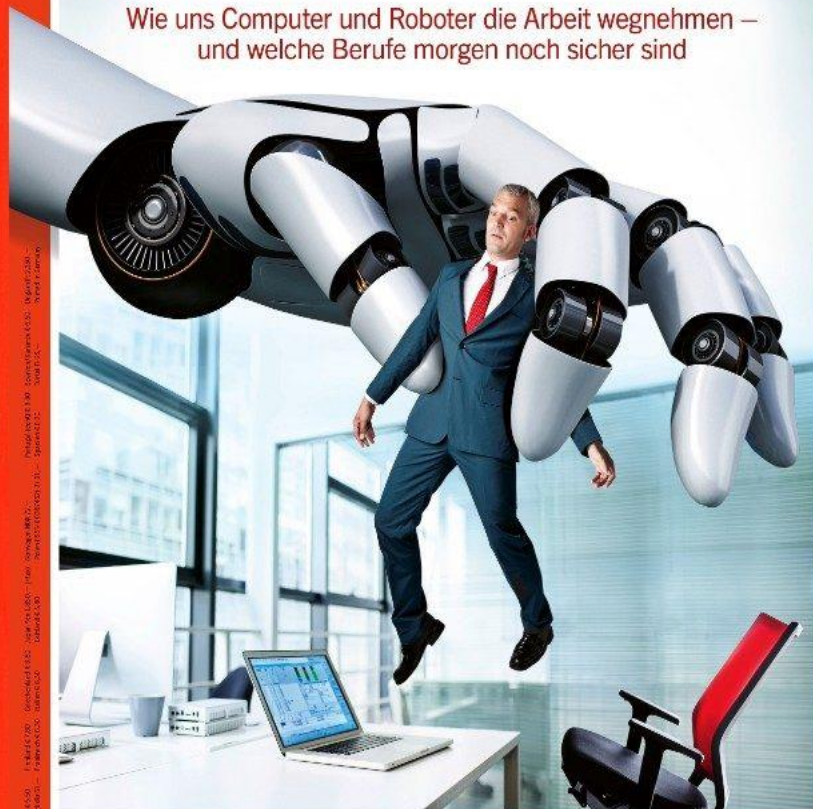
DER SPIEGEL

Nr. 36 / 3.9.2016
Deutschland €4,90



Sie sind entlassen!

Wie uns Computer und Roboter die Arbeit wegnehmen –
und welche Berufe morgen noch sicher sind



Sportarzt Müller-Wohlfahrt
Heilende Hände, aber
Probleme mit den Finanzen

Eine ukrainische Karriere
Erst Gefängnis, dann
Präsidentenpalast?

Steinzeit-Drama
Warum Neandertaler
zu Kannibalen wurden

DER SPIEGEL

Nr. 10
4.3.2023



Versicherungen
Welche Sie
wirklich
brauchen



FLEISCH-ERSATZ
Der Mann, der
uns vegan
machen will

MEDIZIN
Rettende
Baby-OP im
Mutterleib

KÜNSTLICHE INTELLIGENZ

Die neue Weltmacht

Wie ChatGPT und Co. unser Leben verändern



Stand 2023





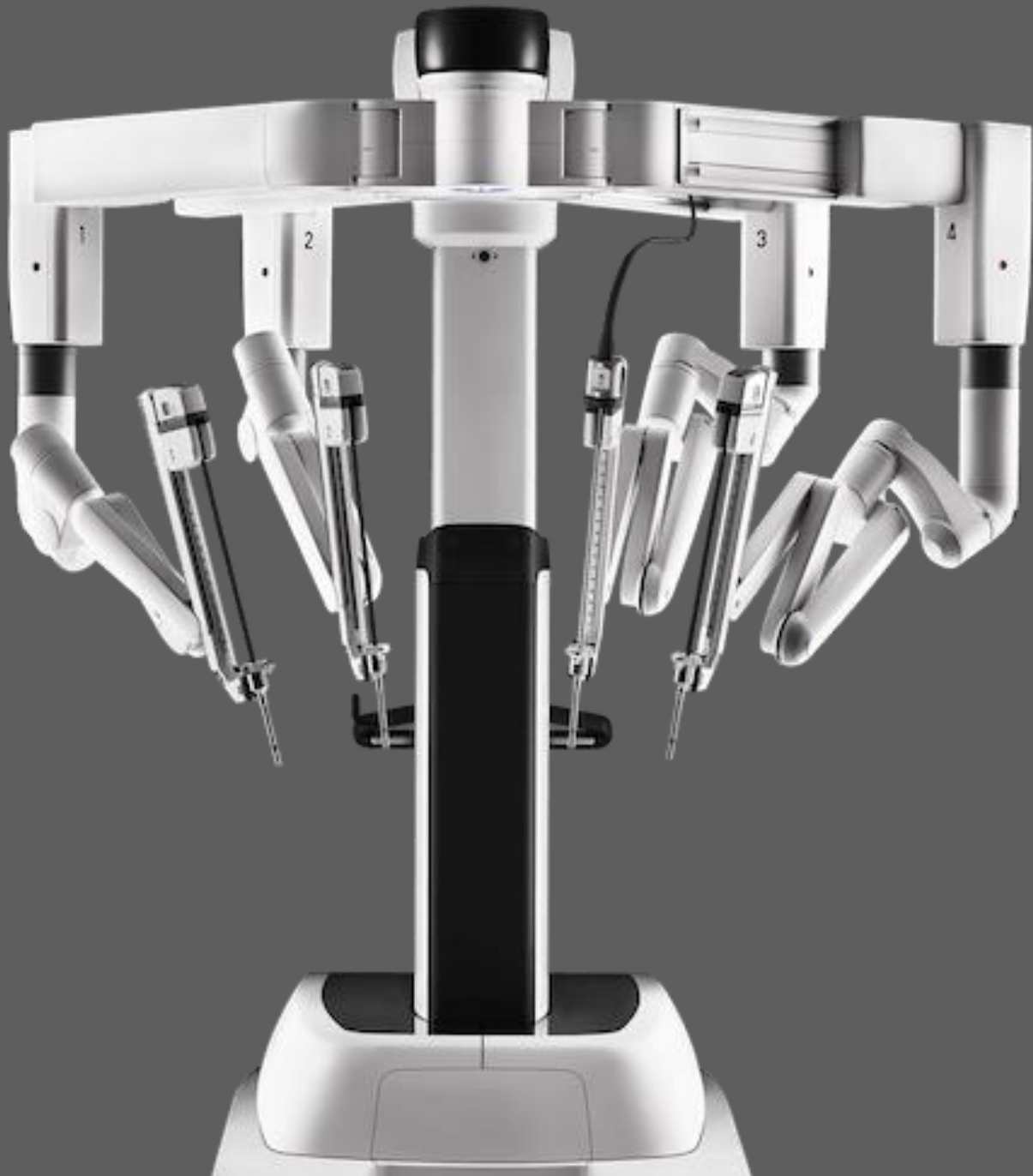
Use Case

Use Case



Use Case





Use Case



Use Case

Use Case



Fundament

/fh/
st. pölten



Daten

Big Data



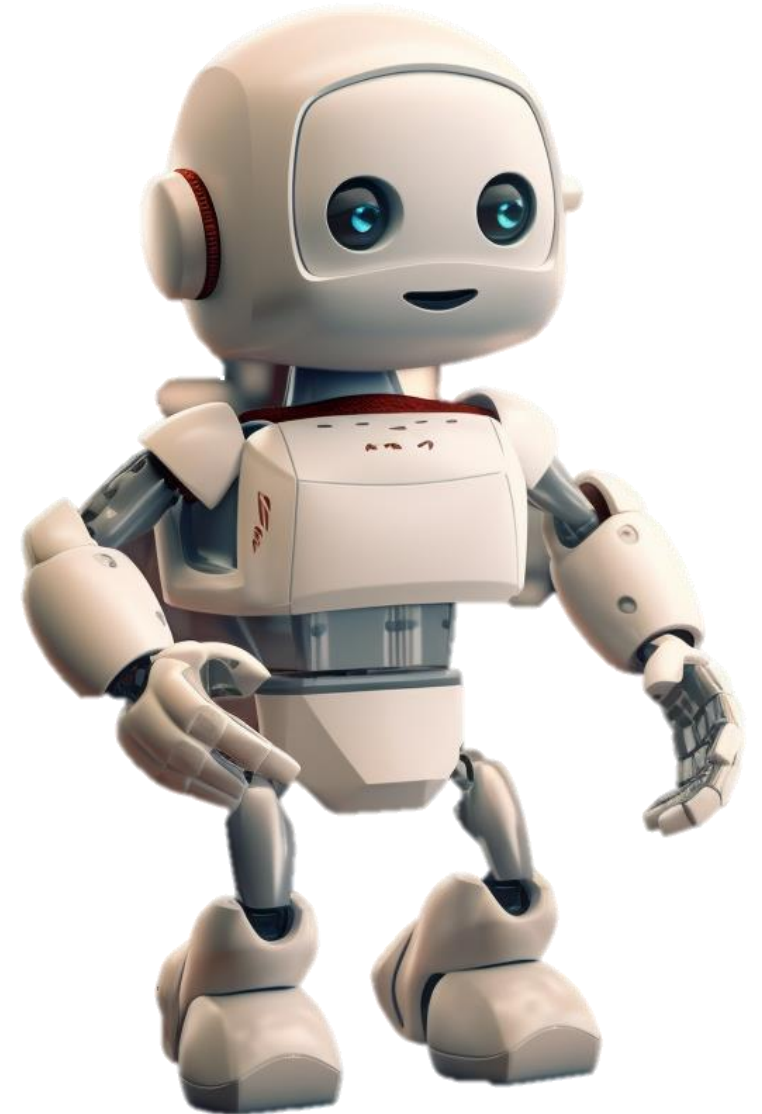
Definition

Künstliche Intelligenz

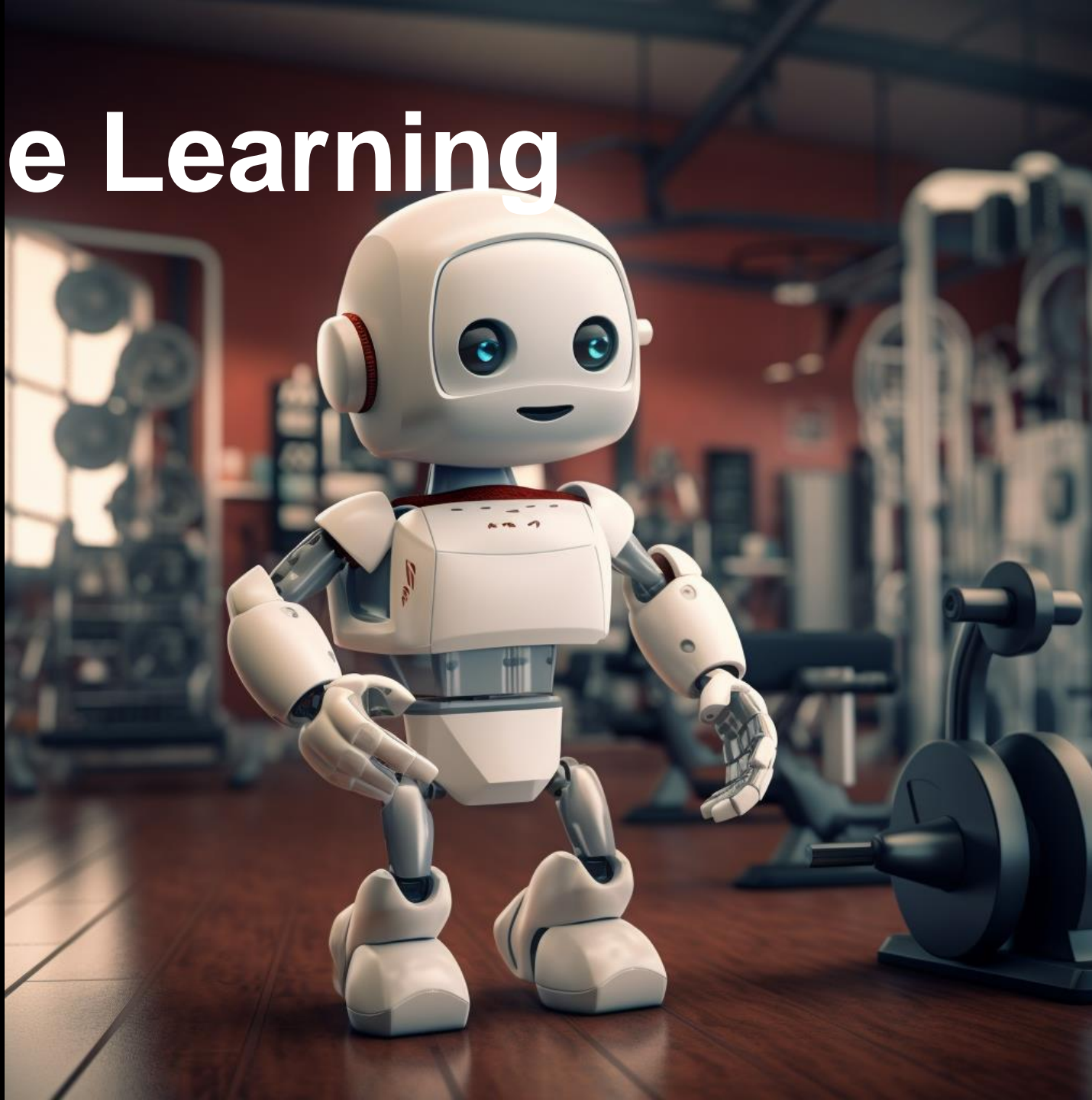
Machine Learning

Deep Learning

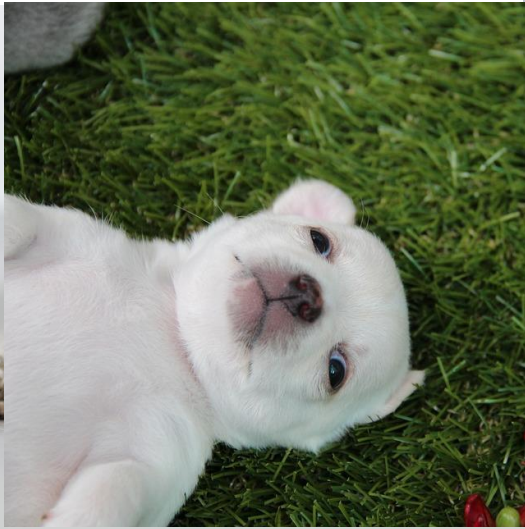
Reinforcement Learning



Machine Learning



Training



Training

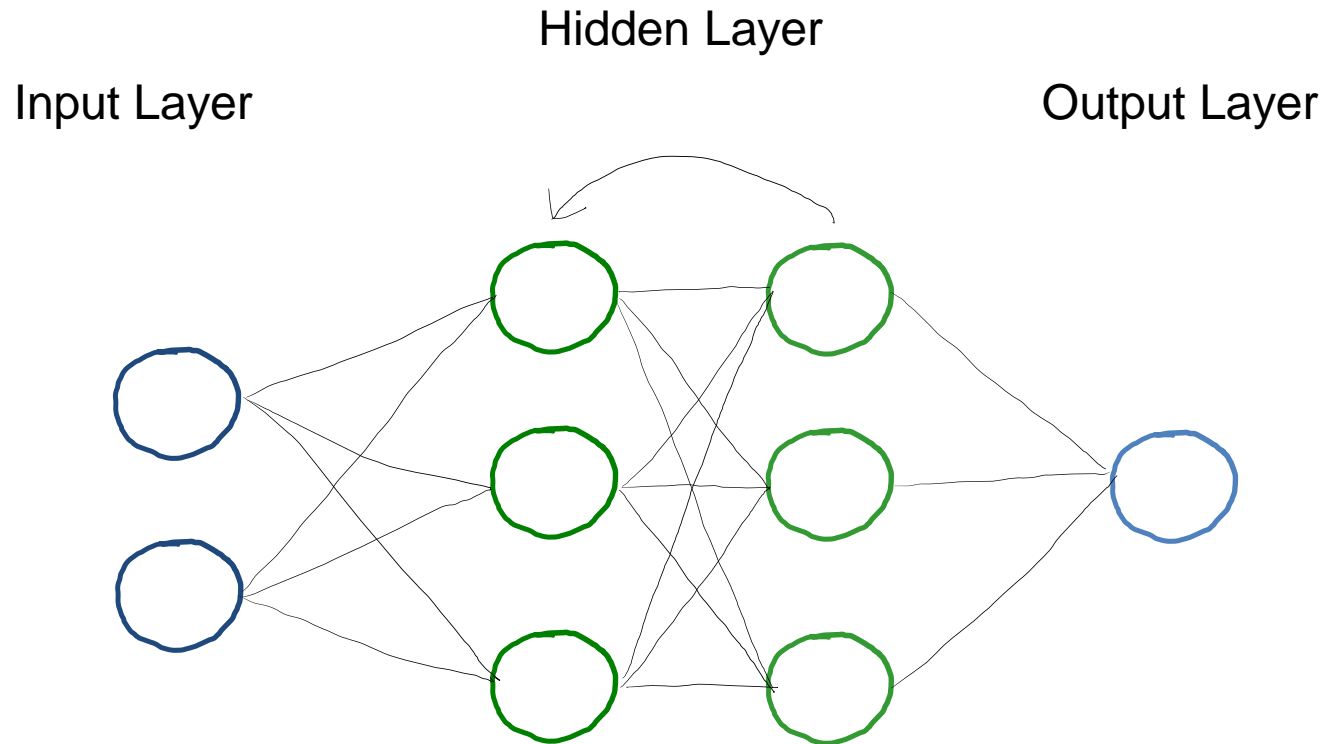


Problem

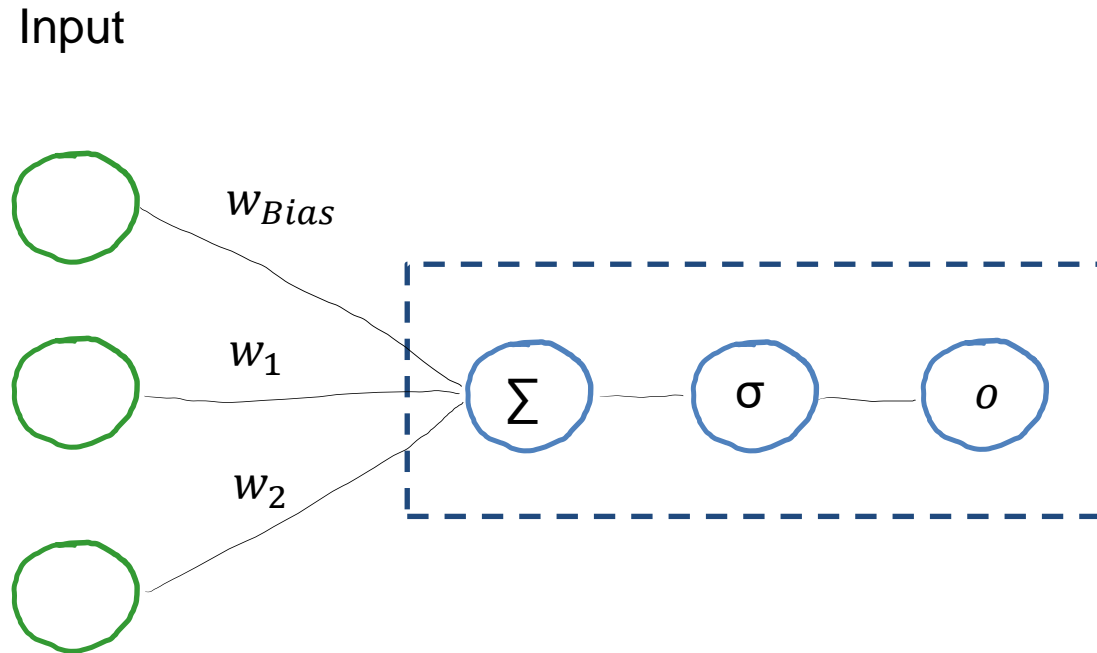
@teenybiscuit



Deep Learning



Deep Learning



Reinforcement Learning



Aktion

Neue Zustand

Belohnung / Bestrafung







ChatGPT

ChatGPT



Examples

"Explain quantum computing in simple terms" →

"Got any creative ideas for a 10 year old's birthday?" →

"How do I make an HTTP request in Javascript?" →



Capabilities

Remembers what user said earlier in the conversation

Allows user to provide follow-up corrections

Trained to decline inappropriate requests



Limitations

May occasionally generate incorrect information

May occasionally produce harmful instructions or biased content

Limited knowledge of world and events after 2021

GPT-3.5
Generative Pre-trained
Transformer 3.5

Language Model

Das Wetter in Österreich wird ?.

$$P_{\theta}(? = \textit{schön} | h) = 0.15$$

$$P_{\theta}(? = \textit{bewölkt} | h) = 0.2$$

$$P_{\theta}(? = \textit{regnerisch} | h) = 0.4$$

$$P_{\theta}(? = \textit{hungrig} | h) = 0.02$$

...





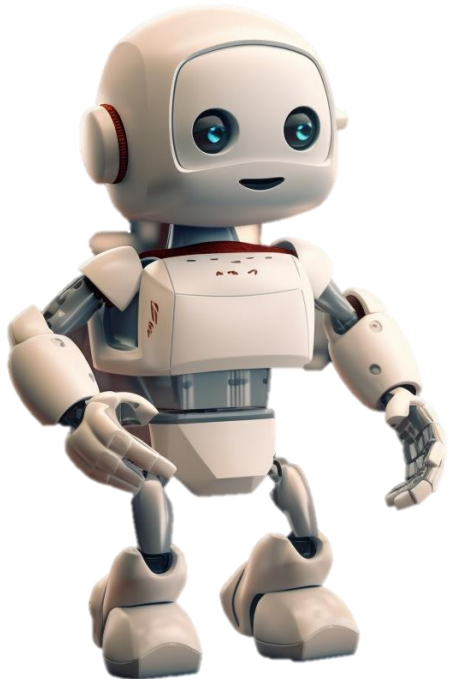
Transformer

+

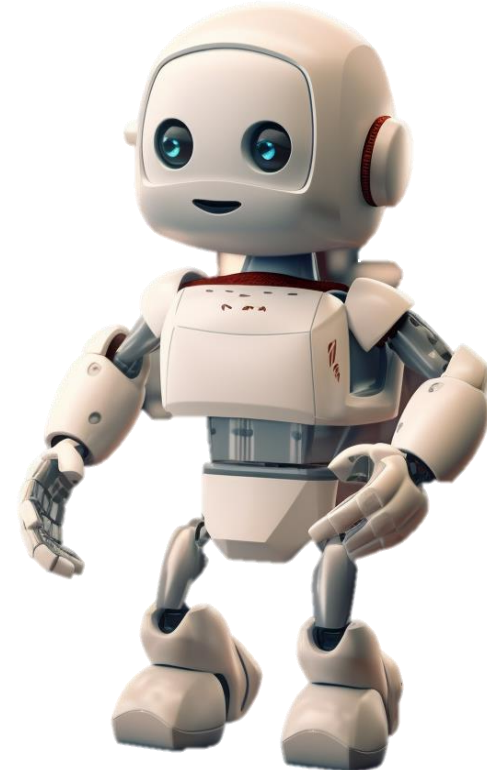
Reinforcement
Learning from Human
Feedback

ChatGPT

Interaktionen werden möglich!



ChatGPT ist ein
Sprachmodell und
arbeitet mit
Wahrscheinlichkeiten.



A photograph of two young girls sitting on a paved path in a park. The girl on the left is wearing a blue t-shirt and jeans, sitting cross-legged and holding a bunch of white daisies. The girl on the right is wearing a white t-shirt with a zebra print and dark jeans, sitting cross-legged. They are both looking at each other and appear to be in conversation. The background shows trees and a path.

Es fühlt sich so real an...



Ziel ist auch, dass
ChatGPT
„menschelt“.

Halluzinationen



[Digital Life](#)

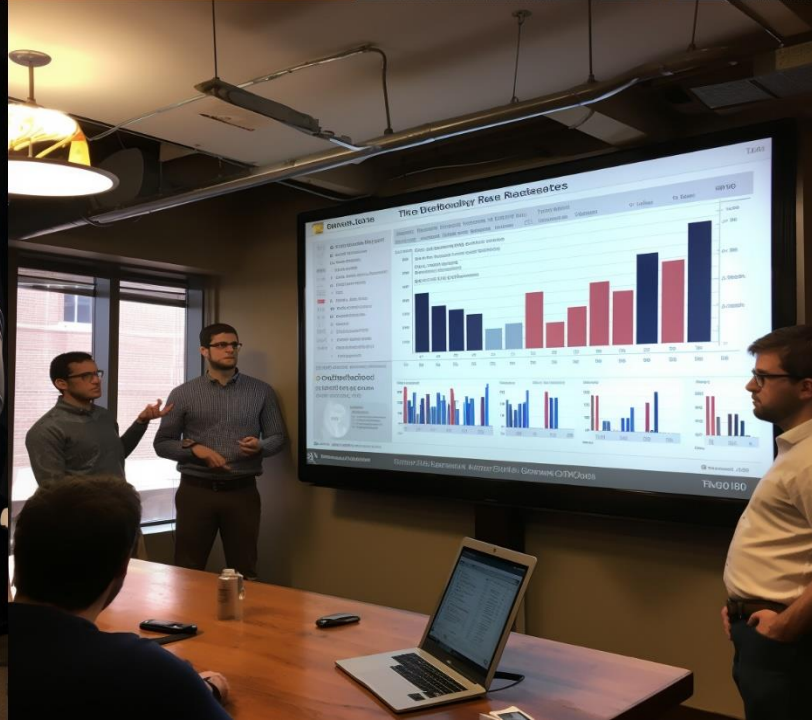
Anwalt reicht ChatGPT-formulierte Klage mit erfundenen Fällen ein

28.05.2023

Gericht in New York deckt auf, dass Anklage mit gefälschten Fällen und Zitaten gespickt ist. Anwalt schiebt Schuld auf ChatGPT.



Sensible Informationen

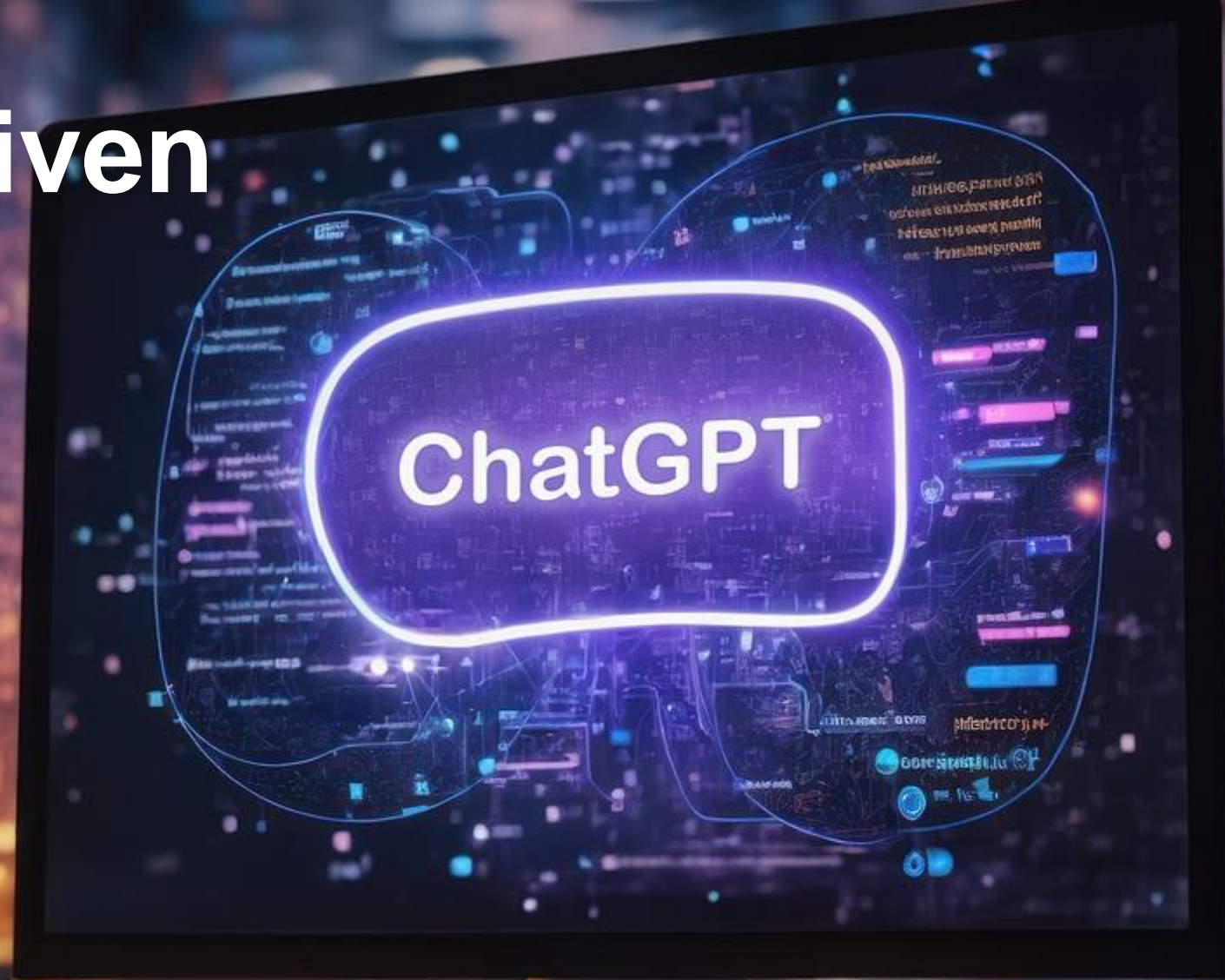


Bias

presenting data visualizations to board of directors - @marlies (fast)

Threats

Alternativen



Human (in the Loop)



Expertise

/fh///
st.pölten



**Herzlichen Dank für die
Aufmerksamkeit!**