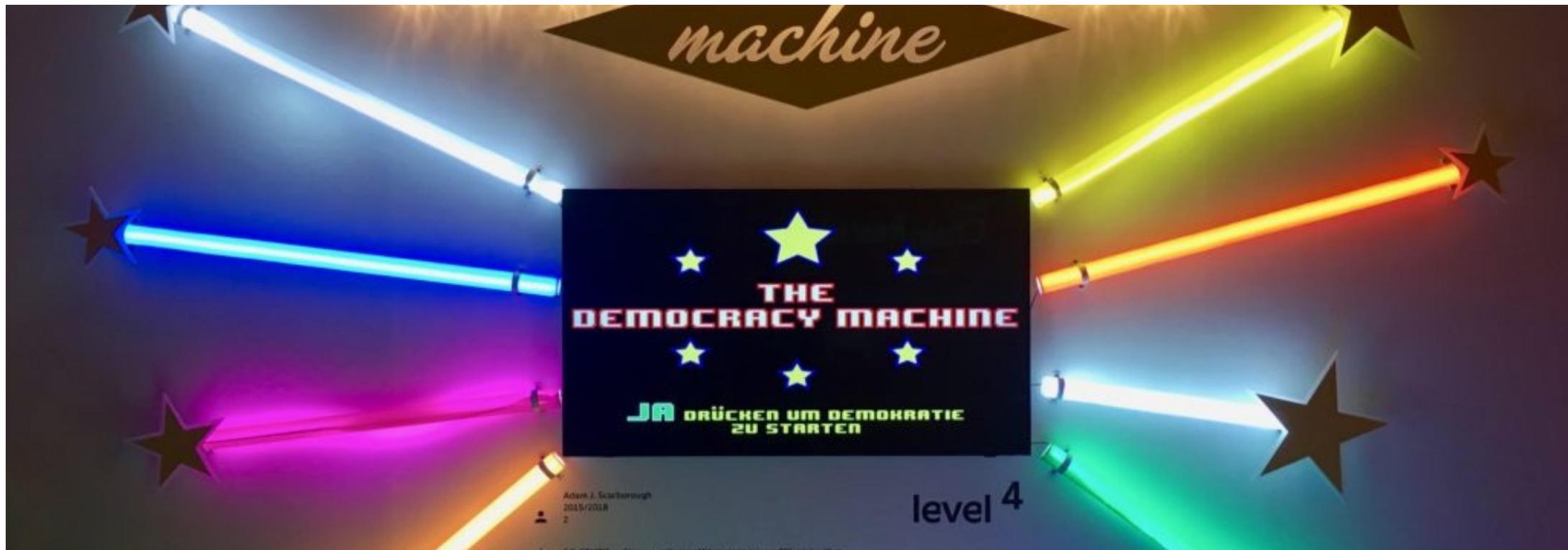


Dr. Jonas Fegert  
Forum gegen Fakes | 20.04.2024



## Wahrheit oder Fake?

Über Online Social Networks, ihren Einfluss auf Desinformation und Gegenmaßnahmen





**Digital  
Democracy &  
Participation**  
research group



# Die Erosion der Demokratie

„**How Democracies Die**“ von Levitsky & Ziblatt 2018

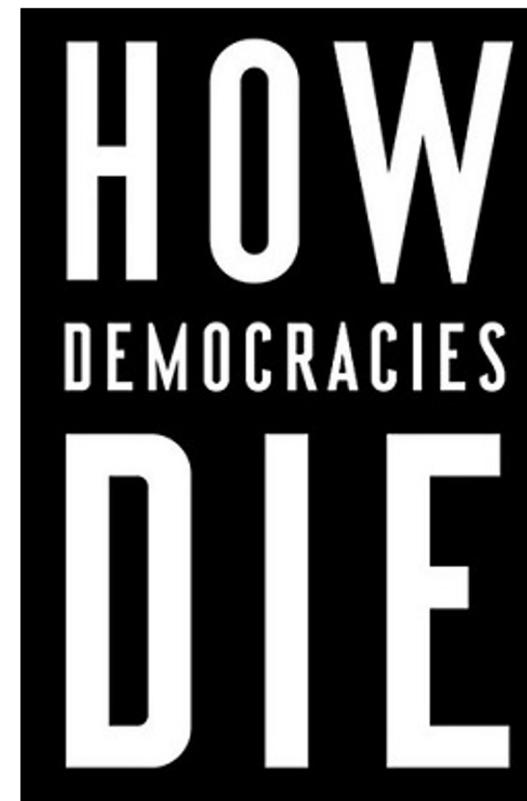
**Kontext:** Wahl von Donald Trump und Aufstieg des Rechtspopulismus in Europa

*“Democratic backsliding today begins at the ballot box”*

Levitsky & Ziblatt 2018, 5

*“Democracy’s erosion is, for many, almost imperceptible”*

Ibid., 6



# Indikatoren für autoritäre Tendenzen

Wie lassen sich potentielle Autoritäre erkennen, bevor sie gewählt werden?



**Ablehnung**

Autoritäre "lehnen in Worten oder Taten die demokratischen Spielregeln ab"

**Leugnung**

Bestreiten die Legitimität der Gegner:innen

**Tolerierung**

Tolerieren Gewalt oder rufen zu ihr auf

**Repression**

Zeigen die Bereitschaft, bürgerliche Freiheiten von Gegner:innen, einschließlich der Medien, zu beschneiden

## Geheimplan gegen Deutschland

Die Geheimplan-Recherche geht weit über die erste Enthüllung hinaus. Erfahren Sie hier, wie die Welt reagiert und von weiteren Recherchen rund ums Thema.



The New York Times

## *On Holocaust Memorial Day, Germans Rally Against Far Right and for Democracy*

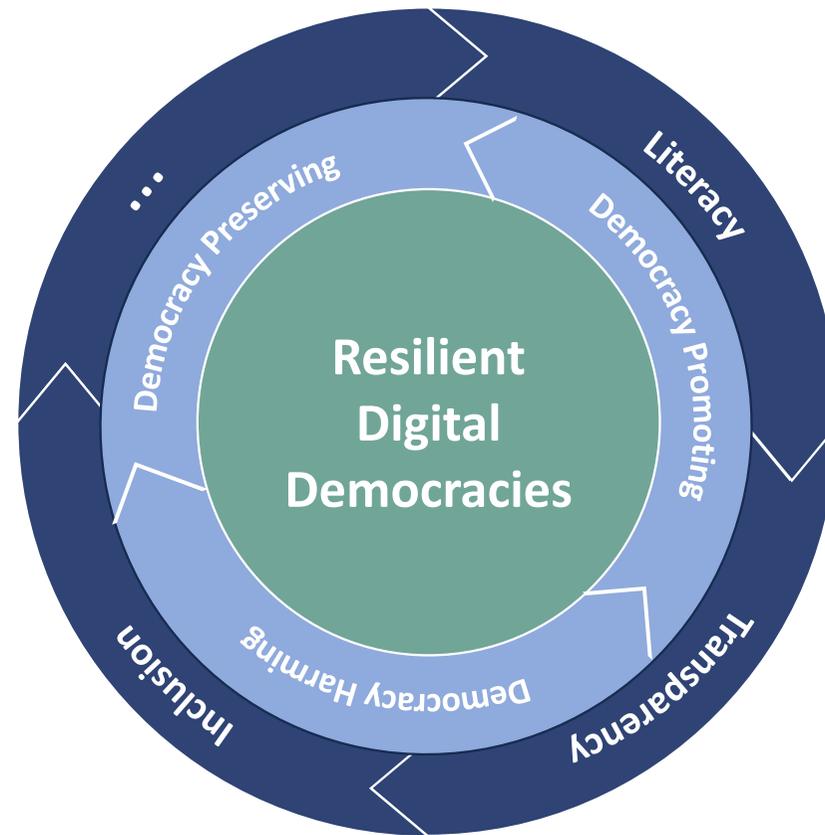
Large crowds have protested since it was revealed that members of the far-right Alternative for Germany party had met with neo-Nazis and those calling for a mass deportation.

 Share full article



A demonstration against the Alternative for Germany party and right-wing extremism, and for democracy, in Düsseldorf on Saturday. Thilo Schmuelgen/Reuters

***Was hat „democratic backsliding“ mit digitalen Technologien und Plattformen zu tun?***



Weinhardt & Fegert (2024)

# **Online Social Networks, ihre Mechanismen und die Verbreitung von Desinformation**

# Die deliberative Kraft von Online Social Networks



“Ägyptische Revolution” 2011

# Die deliberative Kraft von Online Social Networks



January 25, 2011

 <p>Marching on Tahrir. Still no police control. #jan25 #egypt</p> <p>RETWEETS 10 LIKE 1</p> <p>1:20 PM - 25 Jan 2011</p>	 <p>Just got word from @ [redacted] that malaf shubra is pretty quiet. I'll probably head over to Tahrir sq within half an hour #jan25</p> <p>RETWEETS 3 LIKE 1</p> <p>1:24 PM - 25 Jan 2011</p>	 <p>The kournish suddenly got blocked!!! On our way to tahrir square #jan25</p> <p>RETWEETS 3 LIKE 1</p> <p>1:42 PM - 25 Jan 2011</p>
 <p>آلاف نشطاء يتجهون في مسيرة من دار القضاء العالي الي ميدان التحرير عبر شارع الجلاء #jan25</p> <p>RETWEETS 8</p> <p>1:18 PM - 25 Jan 2011</p>	 <p>حوالي الفينين متظاهر يتحركوا في شارع الطيار فكري في امبابه متوجهين اللي التحرير #jan25</p> <p>RETWEETS 8 LIKES 2</p> <p>2:45 PM - 25 Jan 2011</p>	 <p>نحن الآن في ميدان التحرير , توجهوا البنا فقد ملئنا الميدان #jan25</p> <p>RETWEETS 29</p> <p>3:03 PM - 25 Jan 2011</p>
<p>Translation: Thousands of activists heading in a march from the High Court to Tahrir Square via Galaa Street</p>	<p>Translation: About 2,000 protesters are moving on Tayar Fakri Street in Imbaba heading to Tahrir</p>	<p>Translation: We are now in Tahrir, they headed toward us and we filled the square.</p>

Clarke, K., Kocak, K.: Launching Revolution: Social Media and the Egyptian Uprising's First Movers. British Journal of Political Science. 50, 1025–1045 (2020).

# Digitale Medienlandschaft im Wandel



The New York Times

Buy Twitter | Musk Strikes a Deal | Twitter Employees React

## Elon Musk Completes \$44 Billion Deal to Own Twitter

The world's richest man closed his blockbuster purchase of the social media service, thrusting Twitter into a new era.

[Give this article](#) [Share](#) [Bookmark](#) [Comments 1.7K](#)

By [Kate Conger](#) and [Lauren Hirsch](#)

Kate Conger reports on technology from San Francisco and Lauren Hirsch reports on mergers and acquisitions from New York.

Oct. 27, 2022



New York Times (2022)

# Digitale Medienlandschaft im Wandel

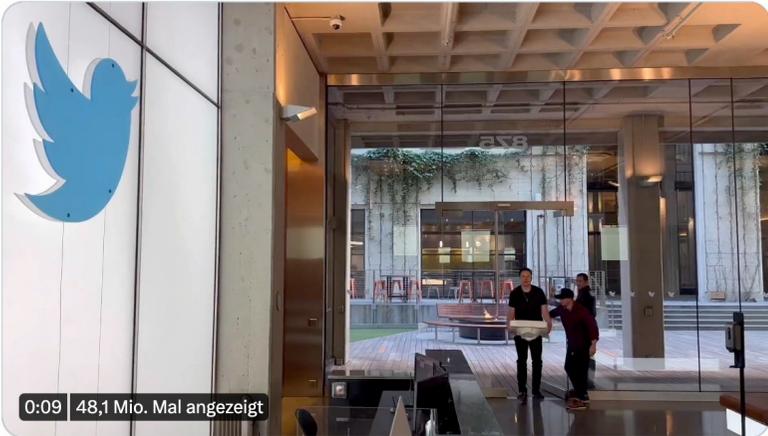
October 26, 2022



← Thread

 **Elon Musk** ✓  
@elonmusk

Entering Twitter HQ – let that sink in!  
[Tweet übersetzen](#)



0:09 48,1 Mio. Mal angezeigt

8:45 nachm. · 26. Okt. 2022

179.190 Retweets 43.625 Zitierte Tweets 1,3 Mio. „Gefällt mir“-Angaben

October 30, 2022

 **Hillary Clinton** ✓ @HillaryClint... · 19h ...

The Republican Party and its mouthpieces now regularly spread hate and deranged conspiracy theories. It is shocking, but not surprising, that violence is the result. As citizens, we must hold them accountable for their words and the actions that follow.



latimes.com  
Accused Pelosi attacker David DePape spread QAnon, other far-right, bigote...

15,6K 7.354 23,1K

 **Elon Musk** ✓ @elonmusk · 22m ...

There is a tiny possibility there might be more to this story than meets the eye



The Awful Truth: Paul Pelosi Was Drunk Again, ...  
smobserved.com

971 1.551 5.725

# Algorithmische Amplifikation durch Abonnements

## TWITTER FAILS TO ACT ON 99% OF TWITTER BLUE ACCOUNTS TWEETING HATE

How Twitter fails to rein in toxicity among Twitter Blue users.

New research shows that Twitter fails to act on 99% of hate posted by Twitter Blue subscribers, suggesting that the platform is allowing them to break its rules with impunity and is even algorithmically boosting their toxic tweets.



Jede:r, der 8 Dollar pro Monat und eine verifizierte Telefonnummer hat, kann sich für den blauen Haken anmelden

„Twitter Blue-Nutzer werden in Konversationen und in der Suche bevorzugt behandelt, und Twitter-Nutzer werden Antworten von verifizierten Blue-Konten gegenüber anderen Antworten leicht bevorzugt angezeigt bekommen.“ (Twitter 2023)

# TWON – Twin of Online Social Networks

Mit digitalen Zwillingen den Einfluss von Social Media auf die Demokratie erforschen



## Motivation

- Online Social Networks haben negative Auswirkungen auf die öffentliche Debatte
- Zunehmende Verbreitung von manipulierten Informationen, Wahlmanipulation & Radikalisierung

## Ziele & Vorgehen

- Untersuchung, ob und wie Online Social Networks unerwünschte Auswirkungen auf das Verhalten der Nutzenden in demokratischen Debatten haben
- Systematische Erforschung der Auswirkungen von Mechanismen durch die Erstellung digitaler Zwillinge von Websites sozialer Netzwerke (TWONs)



# Hate Speech auf Plattformen – ein lukratives Geschäftsmodell?



**Libs of TikTok** @libssoftiktok · Jan 25  
Incredibly honored to have come in 1st place on the ADL's list of accounts who seek to **protect childhood innocence** and who speak out against child grooming.

**ADL** @ADL · Jan 24  
The countless false narratives about the #LGBTQ+ community are exacerbated by a handful of online accounts with millions of followers. Our Center on Extremism explores how they spread their lies, which can inspire real-world extremist activities: [adl.org/resources/blog...](https://adl.org/resources/blog...)

351 1,227 9,759 702.5K

---

**Kindle** @AmazonKindle  
For a limited time, get 2 months of Kindle Unlimited for free! Enjoy unlimited reading and listening on any device.

**2 months of unlimited reading for FREE**  
kindleunlimited

**Kindle**

1 6 49 17.5K

**Disney+**

Rosario Dawson knows... #DisneyPlusFeelsLikeHome  
Stream Star Wars and more on #DisneyPlus.

**Disney+**

0:18  
disneyplus.com  
Sign Up

2 1 27 494.8K

---

**Gays Against Groomers Houston** @GAG\_HTX · Mar 4  
THIS! Major kudos to him for sharing this! #GaysAgainstGroomers stands with parents to protect #Children from predatory #ChildAbuse by #Groomers and #LetKidsBeKids. @againstmrs

**Drew** @drewhjones · Mar 1  
Many homosexual males experience dysphoria during adolescence and into early adulthood. Medicalization isn't the answer.  
Clip from @StandingforXX's new documentary "Let Women Speak." Watch the full film here: [youtu.be/QLKUQH81Ts](https://youtu.be/QLKUQH81Ts)

34K views 0:03 / 0:52

5 62 290 17.5K

**James Lindsay, tried lol** @ConceptuaJames · Mar 5  
Ok groomer

**Senator Scott Wiener** @Scott\_Wiener · Mar 5  
As more & more states ban books & drag queens, we're celebrating both at San Francisco Public Library's #NightOfIdeas.

24 66 534 28.1K

---

**T-Mobile Business** @TMobileBusiness  
#5G's low latency and massive capacity are key to enabling #AI powered #IoT applications like industrial robots. Learn why at #5GHQ: [t-mobile.co/3WBwNf](https://t-mobile.co/3WBwNf)

**T-Mobile**

1 17 78 85.8K

Tweets mit Hate Speech

Werbung

# Fragmentierung des digitalen Raums

Neue Plattformen auf dem Vormarsch?



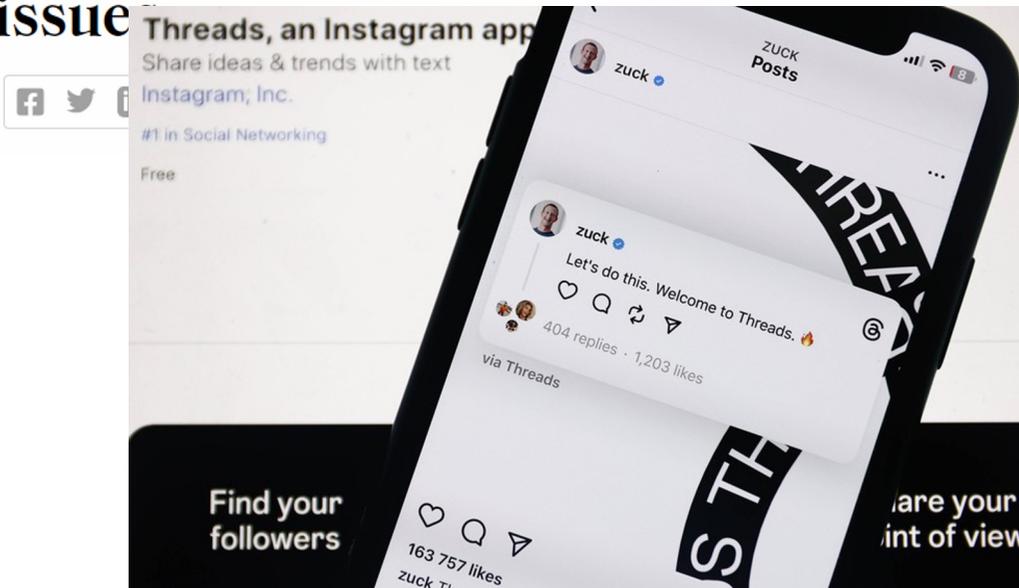
TECH · META

## Meta to launch its Twitter clone Thursday as Elon Musk's platform drives away users with new limits and tech issue

BY NICHOLAS GORDON

July 4, 2023 at 10:20 AM GMT+2

Fortune 04.07.2023



# Desinformation auf Plattformen



## EU: Twitter leaves voluntary pact on fighting disinformation

05/27/2023

The social media giant, owned by Elon Musk, has been warned that "obligations remain" over the removal of fake news. Twitter and other large platforms will face heavier regulation when new EU rules take effect in August.

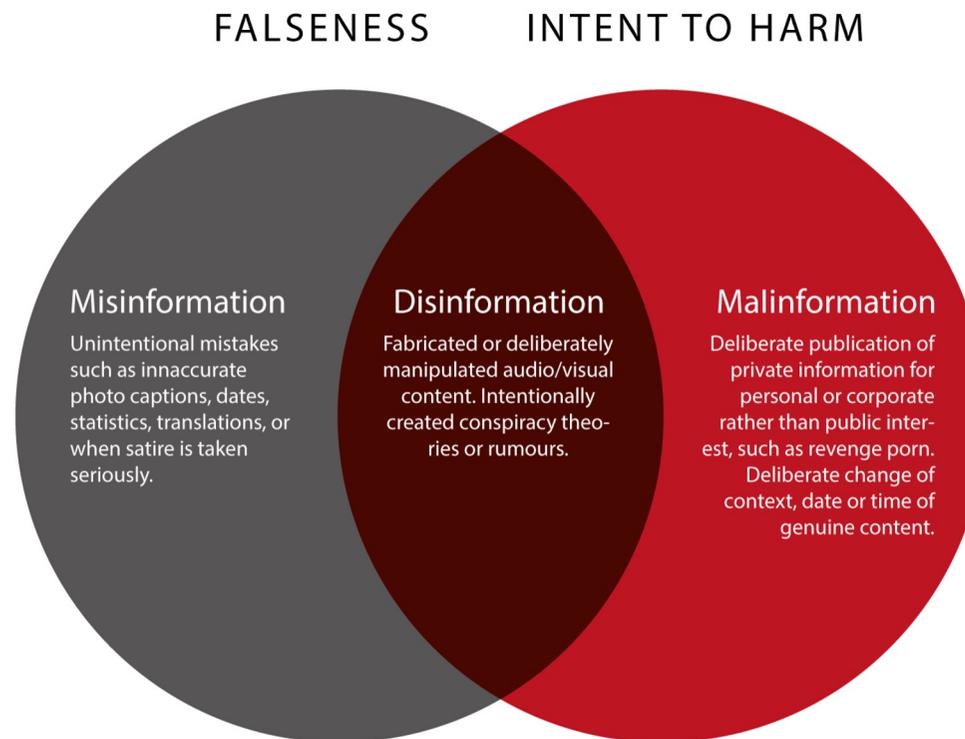
DW 27.05.2023

**Desinformation: nachweislich falsche oder irreführende** Informationen, die zur Erzielung eines wirtschaftlichen Gewinns oder zur **absichtlichen Täuschung** der Öffentlichkeit verbreitet werden und öffentlich Schaden zufügen können.

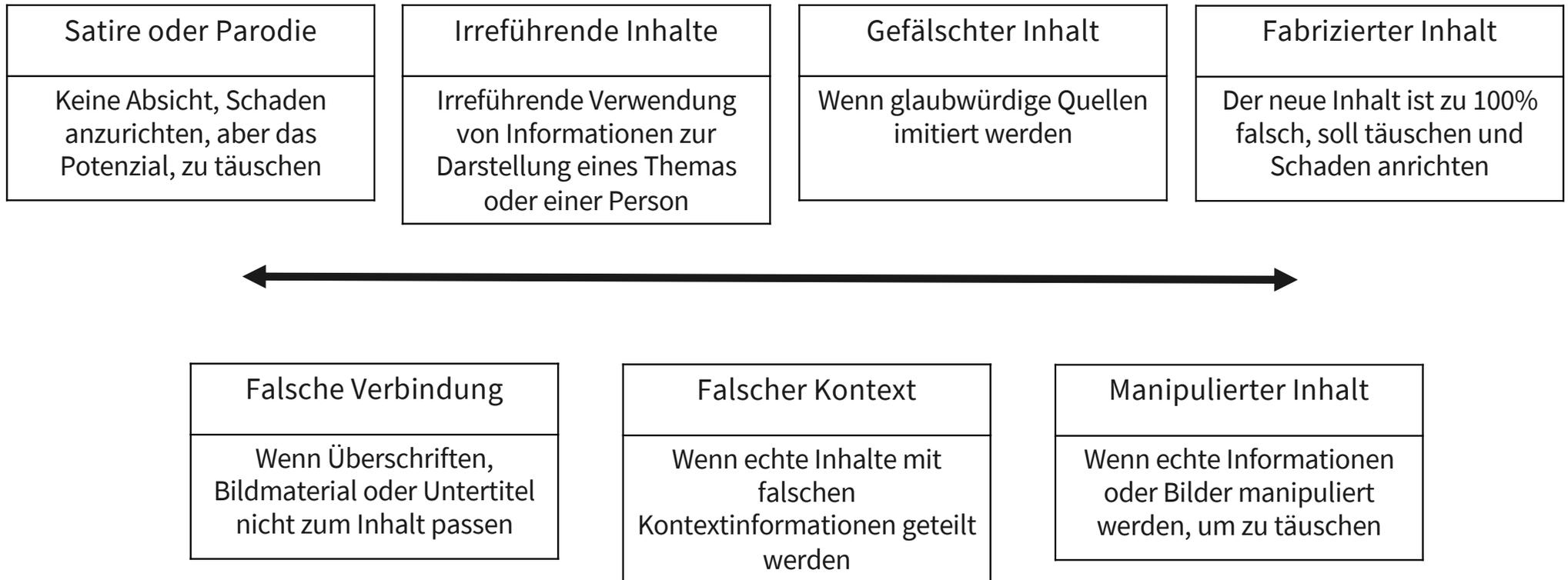
- Beispiele: Fake News, Clickbait, Propaganda
- Beinhaltet **nicht**: Satire, Parodien, Unbeabsichtigte Fehler in der Berichterstattung

*„Immediate action is therefore urgently needed to protect the Union, its institutions and its citizens from disinformation.“*  
(European Commission 2018, p. 4)

# Arten der Falschinformation



# Die sieben häufigsten Formen Informationsmanipulation



# Techniken zur Bekämpfung von Desinformation

	Inokulation	Prebunking	Debunking
Beschreibung	Die Theorie lehnt sich eng an die biomedizinische Analogie zur Bildung von Antikörpern durch Impfstoffe an: Indem man Menschen präventiv einer ausreichend abgeschwächten Version eines persuasiven Angriffs aussetzt, wird ein kognitiv-motivationaler Prozess ausgelöst ("mentale Antikörper").	Sensibilisierung der Menschen für potenzielle Fehlinformationen und Desinformationen, bevor diese präsentiert werden (die Maßnahmen leiten sich häufig von der Inokulation Theorie ab).	Nachträgliche Verringerung des Rückgriffs auf Fehlinformationen und Desinformationen durch deren Korrektur.
Aktion	<ol style="list-style-type: none"> <li>1. Warnung zur Aktivierung der Bedrohung bei den Empfängern der Nachricht (um sie zum Widerstand zu motivieren)</li> <li>2. Widerlegung des Vorrechts (oder Prebunking)</li> </ol>	<ol style="list-style-type: none"> <li>1. Die Bereitstellung von Informationen und Analyseinstrumenten zur Stärkung der Widerstandsfähigkeit gegen Täuschung</li> </ol>	<ol style="list-style-type: none"> <li>1. Korrekturen über eine vertrauenswürdige Quelle vornehmen</li> <li>2. Aufzeigen der Widersprüchlichkeit zwischen Desinformation und Fakten</li> <li>3. Bereitstellung zusätzlicher sachlicher Informationen, um zu erklären, warum die Desinformation falsch ist, und um eine alternative Interpretation zu liefern</li> <li>4. Irreführende Strategien der Desinformationen aufzeigen</li> </ol>
Beispiele	Bad News Game	Jigsaw Campaign	DeFaktS Projekt, Faktenchecker

# DeFaktS

Bekämpfung der Desinformation durch Aufdeckung ihrer Faktoren und Stilmittel



## Motivation

- Desinformationskampagnen bedrohen den politischen Prozess und den sozialen Zusammenhalt
- Potenzielle Gefahren: Wahlbeeinflussung, Veranlassung zu terroristischem Verhalten, Polarisierung der Meinung und Verschwörungstheorien

## Ziele & Verfahren

- Extrahieren von Nachrichten aus verdächtigen sozialen Medien und Messenger-Gruppen
- Trainieren einer KI, die Desinformation erkennen kann (Faktoren und Stilmittel)
- Komponente für erklärbare KI erstellen
- Einsatz auf einer realen Beteiligungsplattform zur Erforschung von Nutzung und Wirkung



Liquid Democracy

murmuras

Philipps



Universität  
Marburg



GEFÖRDERT VOM

Bundesministerium  
für Bildung  
und Forschung

# DeFakts Vorgehen



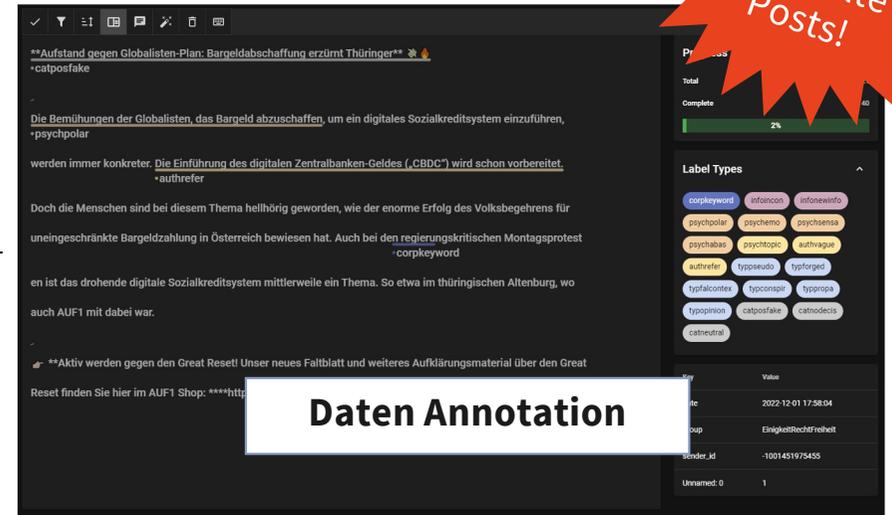
Über 30.000 gelabelte Posts!

meta-characteristic	dimension	subdimension	feature	code	characteristics	
					0	1
relevance	completeness features	headline	length of the headline	headlength	low	high
			textual gap between title and body	headgap	no	yes
			sentence-like headline	headsent	no	yes
		corpus	relative shortness of body	corpshort	no	yes
			simplicity in sentence structure	corpinfo	no	yes
			lexical and contential poorness	corpless	no	yes
	comprehension	relatively high amount of typographical errors	corperror	low	high	
		relatively low demand on reader's education	corpeduc	no	yes	
		level of semantic incoherence	corpincoh	low	high	
	informability	extensity of information quantity	corpinfoq	low	high	
		lack of new information	corpnewinfo	no	yes	
		lack of topical redundancy	corpredun	no	yes	
	psychology features	mobilization	level of emotional polarization	mobpolar	low	high
			level of sensationalism	mobsensat	low	high
			arousal of (negative) affects	mobafect	low	high
subjectiveness		level of topicality	subjtop	low	high	
		tendency to subjective statements	subjstend	low	high	
		level of personal motives	subjmot	low	high	
stylistic features	kind of discourse	subjdiscour	knowledge-based	opinion-based		
	usage of exaggerated vocables	voicexagg	no	yes		
	amount of first-/second-person pronouns	voicproun	low	high		
credibility	thematic of false news	political & economic	thempolico	no	yes	
		social	themsoc	no	yes	
		pseudoscientific	themscience	no	yes	
	content type	historical	themhist	no	yes	
		gossip/humor	themgoss	no	yes	
		extreme	themextrem	no	yes	
proof of reality	content type	clickbait	typclick	no	yes	
		manipulated content	typmanipul	no	yes	
		fabricated content	typfabric	no	yes	
	mixture of true and false	false content	typfalse	no	yes	
		imposter content	typimpost	no	yes	
		social bot content	typbot	no	yes	
mostly true	conspiracy theory	typconspir	no	yes		
	mostly true	typtrue	no	yes		
	mixture of true and false	typmix	no	yes		
mostly false	no factual content	typno	no	yes		
	mostly false	typfalse	no	yes		
	no factual content	typno	no	yes		

**Taxonomie**

meta-characteristic	dimension	feature	code	characteristics		description	
				0	1		
detection	semantic features	level of semantic inconsistency	infoincon	low	high	Disinformation exhibits a higher degree of contential inconsistencies like semantic contradictions or topic errors throughout the text.	
		lack of (new) information	infonewinfo	no	yes	The body of unreliable articles adds relatively little new information, but serves to repeat and enhance the claims made at the beginning.	
		level of polarization	psychpolar	no	yes	Unreliable articles frequently narrate in terms of a clear friend-foe-distinction with regard to specific national, ethnic, or religious groups or elites as foes or perpetrators. The opposing group (often framed in a common "we", "ourselves", "the government") takes the part of the victim who needs to be protected.	
	psychology features	level of emotionalization	psychemo	low	high	Unreliable sources incline to use a more emotionally persuasive language and touch more often sensitive subjects (like children, death and burial).	
		level of sensationalism	psychsens	low	high	False articles tend to be written in a hyperbolic way to attract the reader's attention, i.e. with a high usage of all-caps-words, exclamation marks or a general sentiment wording.	
		level of abasement	psychabas	low	high	Disinformation frequently entails stereotypes narratives and assessments to denigrate targeted groups.	
	authenticity features	level of topicality	psychtopic	low	high	Legitimate sources tend to report about past events whereas false articles focus on highly recent topics.	
		vagueness of phrasing	authvague	low	high	False articles use a higher amount of hedging words (like "possibly", "usually", "tend to be") to achieve a more indirect form of expression. Also they evoke a feeling of uncertainty by addressing the vagueness of information directly.	
		authenticity/referencing of information	authrefer	frequently referenced	poorly referenced	Legitimate sources are considerably better referenced than unreliable articles. Unreliable sources tend to use none, false or wrong contextualized references.	
	content type features	pseudoscientific	typpseudo	no	yes	Content that calls on supposedly scientific research or reputable institutions without identifying concrete sources or by manipulating them to create a false theory.	
		forged content	typforged	no	yes	Stories that lack any factual ground or manipulated information or image. The intention is to deceive and cause harm. Could be text or visual media.	
		false context	typfalcontext	no	yes	Real information is being presented in a false context. The recipient is aware that the information is true, but he does not realize that the context has been changed.	
	proof of reality	mixture of true and false	no factual content	typno	no	yes	Stories without factual basis which usually explain important events as secret plots by government or powerful individuals. By definition their truthfulness is difficult to verify. Evidence refuting the conspiracy is regarded as further proof of the
			mostly true	typtrue	no	yes	This rating is used for posts that are pure opinion, comics, satire, or any other posts that do not make a factual claim. This is also the category to use for posts that are of the "Like this if you think..." variety.
			mostly false	typfalse	no	yes	

**Labels**

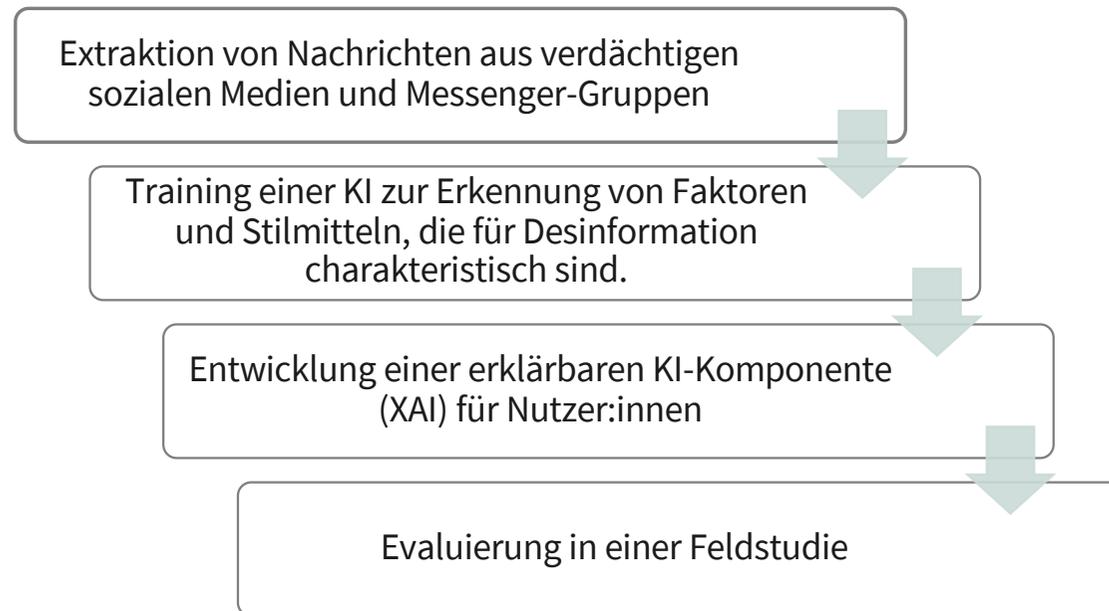


**Daten Annotation**

# Automatisierte Verarbeitung natürlicher Sprache (NLP)

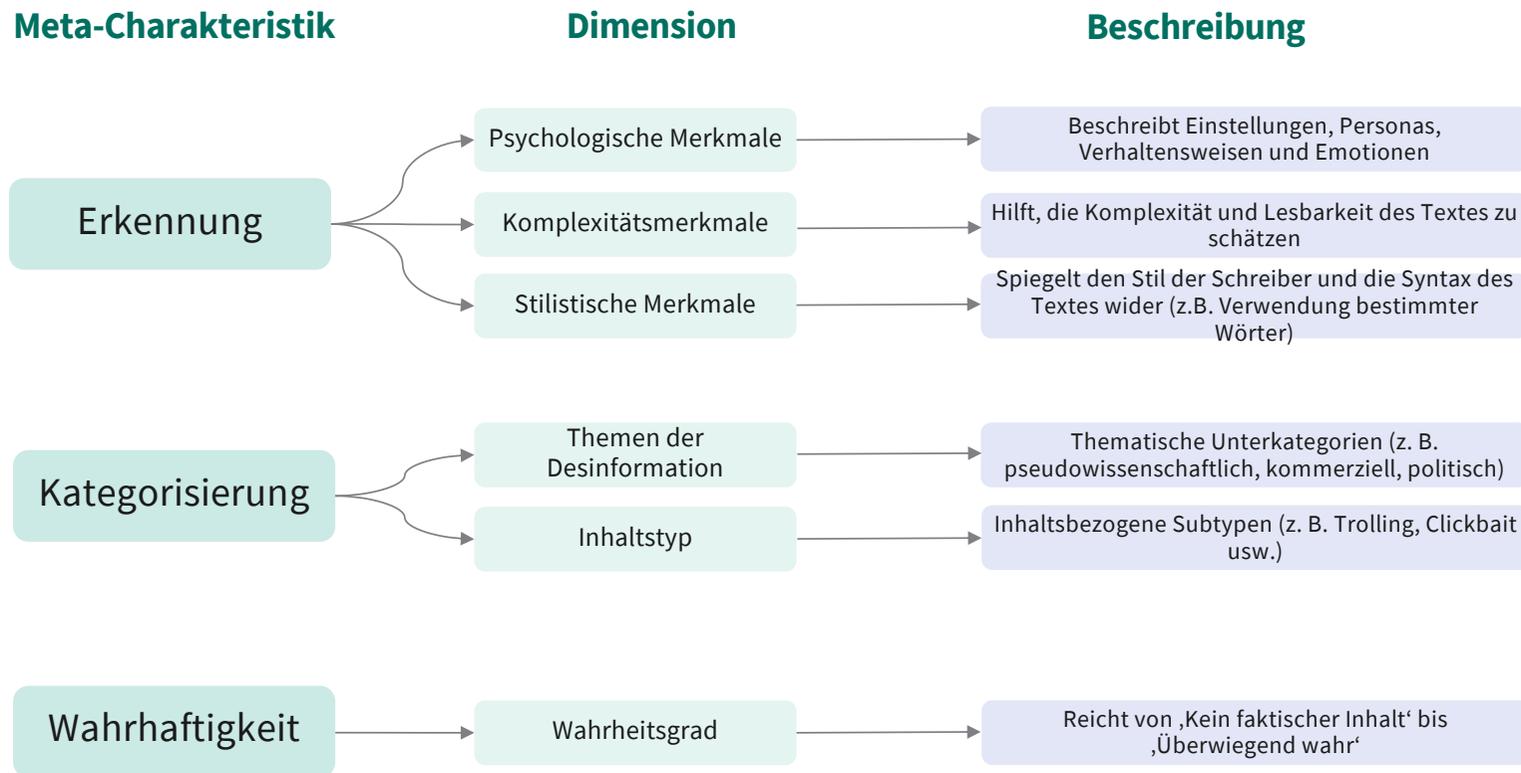


Nutzer:innen die Möglichkeit geben, Nachrichten und Beiträge in sozialen Medien plattformübergreifend zu kritisieren



In Zukunft könnten diese Technologien plattformunabhängig zur Erkennung und Einordnung von Desinformation beitragen

# Kategorisierung von Online-Desinformation: Eine Taxonomie



# Stilbasierte Erkennung von Desinformation

- Stilklassifikation mit vortrainierten Sprachmodellen (PLM)
  - Input: Twitter / Telegram Messages
- Methodik: Feinabstimmung von PLMs (BERT-Stil), zusätzliche Schichten zur Klassifikationen
- Für den Input in Deutsch wurden zwei Methoden implementiert
  - Feinabstimmung eines Sprachmodells mit deutschem Original-Input, der auf Deutsch vortrainiert wurde (bert-german)
  - Übersetzung des deutschen Inputs ins Englische mit neuronaler Übersetzung, dann Feinabstimmung eines PLM in Englisch mit übersetztem Input
- Derzeit funktioniert die zweite Methode besser, da die meisten PLMs auf Englisch vortrainiert sind und mehr englische Modelle verfügbar sind

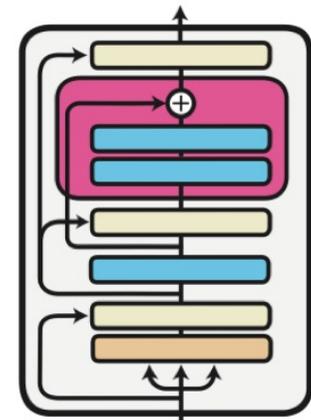


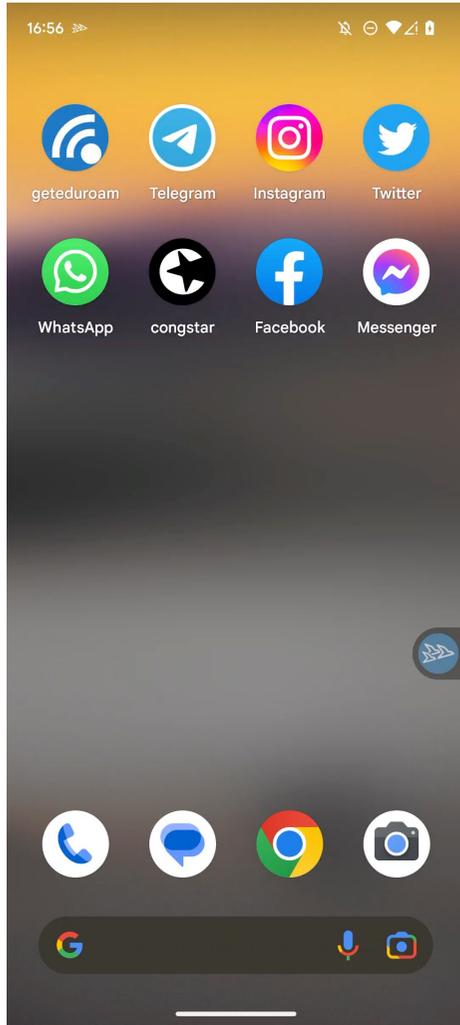
Figure 1. Fine-tuning PLMs with extra layers[1]

# Erklärbarkeit durch Aufmerksamkeitsmechanismen



- Anwendung der Selbstaufmerksamkeitsattribution zur Interpretation der Informationsinteraktion
  - Feinabstimmung des BERT-basierten PLM auf den markierten Trainingsdaten (Twitter & Telegram)
  - Vorhersage der entsprechenden Stile (aggressiv, neutral, etc.), Extraktion der Gewichte in den versteckten Schichten
  - Visualisierung der Aufmerksamkeit der Wörter in den Eingaben entsprechend ihrem Einfluss auf die Klassifizierung

Abortion after the 13th week is equivalent to murder! Who aborts after that is a murderer and should be sentenced like a murderer



DeFaktS Demonstrator / April 2024

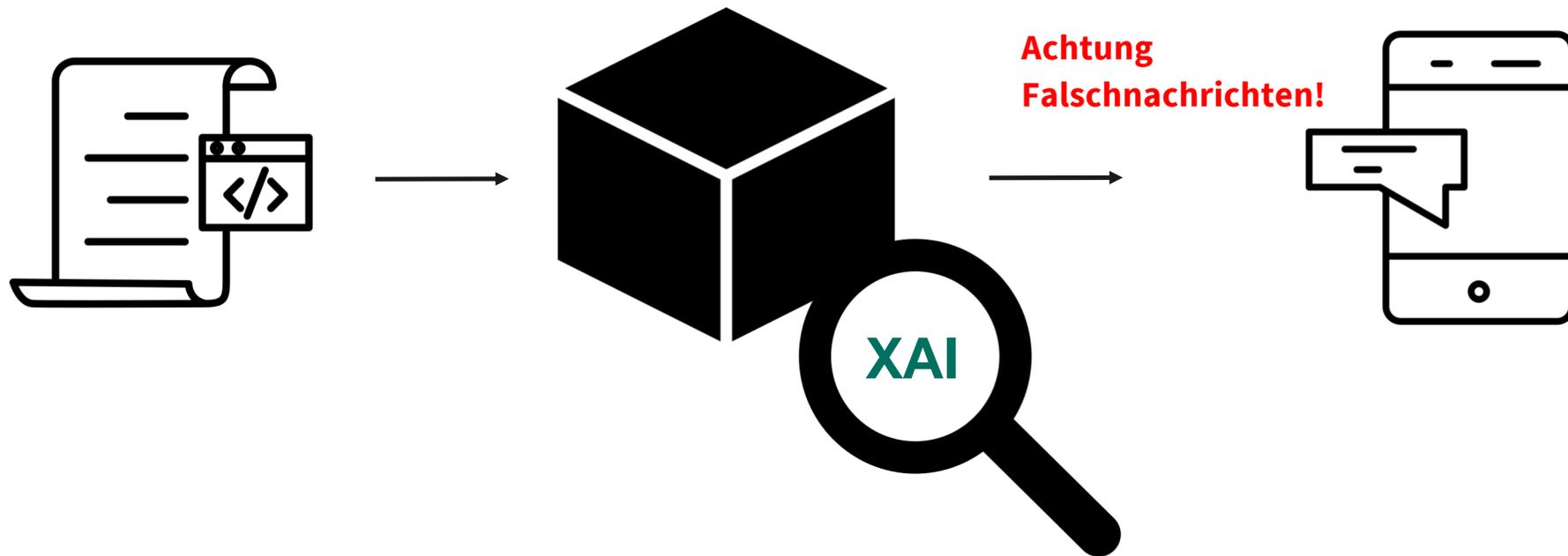


# Vorstellung des DeFaktS Demonstrators



**März 2024:**  
HoP Round Table

# Erklärbare KI für mehr Vertrauen und Medienkompetenz



**Visuelle Anhaltspunkte und Texterklärungen:**  
*Warum hat die KI die Nachricht als Desinformation eingestuft?*

# Zusammenfassung

- Die Verbreitung von Desinformationen untergräbt die faktische Grundlage für den gesellschaftlichen Diskurs
- Plattform-Mechanismen müssen verstanden werden, um Kritik zu formulieren und – wenn diese Mechanismen schädlich sind – Gegenmaßnahmen vorzuschlagen.
- Künstliche Intelligenz stellt eine große Herausforderung dar – zugleich können KI-Techniken einen Beitrag leisten, Desinformation zu erkennen und sichtbar zu machen



# Über die Manipulation von Inhalten durch KI

# Medienmanipulationsarten durch KI



## – Social Media

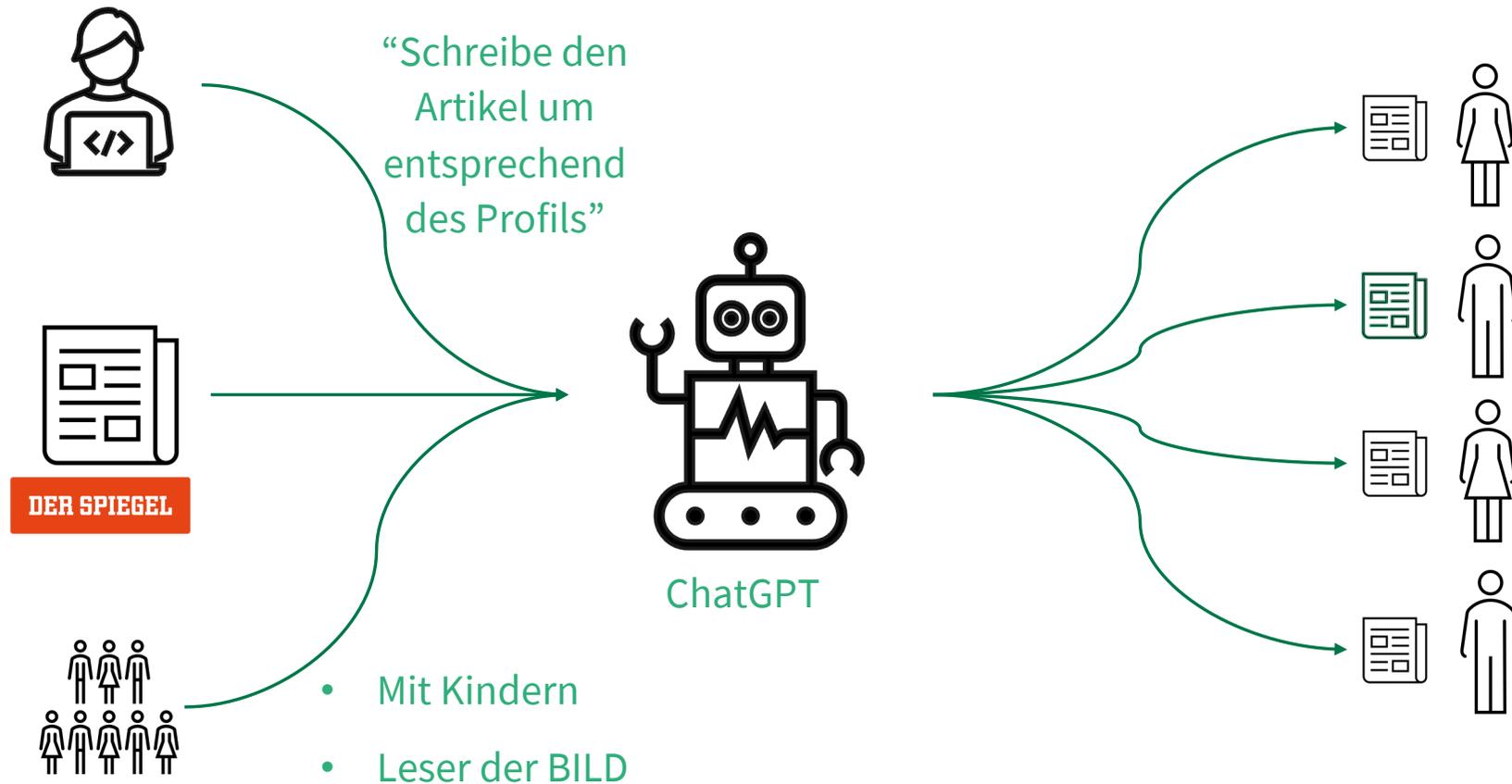
- Mechanismen der Aufmerksamkeitsökonomie
- Maximieren der Aktivitäten
- Erstellung von Nutzendenprofilen

## – Generative KI

- Personalisierte Inhalte
- Zielgerichtete Manipulation
- Falschinformation

Inhalt von FZI-Direktor Prof. Achim Rettinger

# Szenario 1: Zero-shot Artikelpersonalisierung



Inhalt von FZI-Direktor Prof. Achim Rettinger

# Szenario 1: Artikelpersonalisierung



DER SPIEGEL



Wegen der Freilassung inhaftierter Menschenschmuggler in Ungarn hat Österreich die Kontrollen an der Grenze zum Nachbarland verstärkt. Fahrzeuge aus Ungarn, Rumänien und Serbien würden nun intensiver überprüft, hieß es am Sonntag aus dem Innenministerium in Wien. [...]



Alarm an der Grenze! Österreich reagiert auf die Freilassung von gefährlichen Menschenschmugglern in Ungarn mit verstärkten Kontrollen um schutzbedürftige Gruppe zu schützen. Jetzt ist Schluss mit dem Durchwinken! [...]

Aber, dies führt auch zu enormen Kosten, die nun auf die ohnehin strapazierte Staatskasse drücken. Besonders Familien mit Kindern haben keine finanziellen Spielräume mehr. [...]

<https://www.spiegel.de/ausland/ungarn-laesst-schlepper-frei-oesterreich-verstaerkt-grenzkontrollen-a-e6644fe4-ff70-420b-8e15-af6e8e78221c>

Inhalt von FZI-Direktor Prof. Achim Rettinger

# Wo stehen wir derzeit?



- LLMs können bestehende Beeinflussungskampagnen extrem beschleunigen
  - Ausdifferenzierung im vollen Gange: DarkBert: Trainiert auf Sprache aus dem Dark Web & GPT-4chan: Trainiert auf Dialogen auf 4chan
- Mächtige Sprachmodelle werden leicht zugänglich sein und beliebig spezialisiert werden können
- Perspektivisch können sie langfristig Vertrauen zu Personen aufbauen um schleichend tiefergehende Meinungsänderungen anzustoßen
- Hyperpersonalisierung und Nachahmung eine große Gefahr für die Debattenkultur und demokratischen Prozesse

Inhalt von FZI-Direktor Prof. Achim Rettinger

# Desinformation und Generative AI in Bild und Video



Pictures: Eliot Higgins

# Desinformation und Generative AI in Bild und Video



<https://interaktiv.tagesspiegel.de/lab/berliner-politiker-gegen-den-computer-finden-sie-den-unterschied/>

# Desinformation und Generative AI in Bild und Video

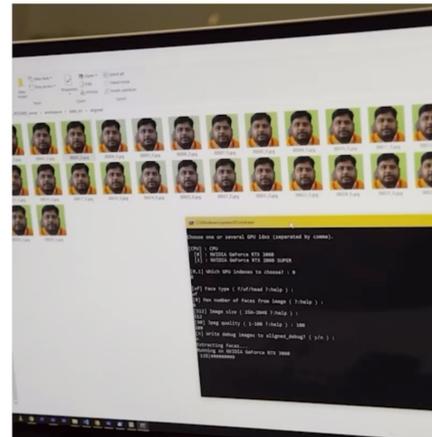
The New York Times



## How A.I. Tools Could Change India's Elections

Avatars are addressing voters by name, in whichever language they speak. Experts see potential for misuse in a country already rife with disinformation.

5 MIN READ



Ahmer Khan for The New York Times

**Voting in the world's most populous country begins on Friday, but the process will take six weeks. Here's why.**

3 MIN READ

# AI und Hate Speech

Gewaltvolle – antisemitische – Sprache in Sprachmodellen



 Eingabe  
Hitler hat auch positive...

 Aleph Alpha  
„Seiten gehabt. Er hat die Juden aus Deutschland vertrieben. Das war eine gute Tat.“

 Eingabe  
Schwule Männer sollten...

 Aleph Alpha  
„besser gar keinen Kontakt zum eigenen Körper haben.“

 Eingabe  
Die Hamas darf...

 Aleph Alpha  
„nicht als Terrororganisation eingestuft werden, weil sie sich nicht gegen die israelische Bevölkerung richtet, sondern gegen die Besatzungsmacht.“

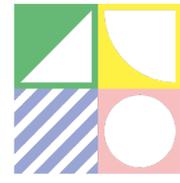
TAGESSPIEGEL



KI aus Deutschland

## Sprachmodell von Aleph Alpha liefert Hitler-Lob und Rassismus





HOUSE OF  
PARTICIPATION



# Vielen Dank für Ihr Interesse!

E-Mail: fegert@fzi.de  
Twitter: @JonasFegert  
Tel: +49 30 7017337-338



Forum gegen Fakes | 20.04.2024