

Anonymität in Gefahr?

Gerichtsurteile mit Sprachmodellen re-identifizieren

Alex Nyffenegger

Softwareentwickler, APP Unternehmensberatung AG

Publikation von Gerichtsurteilen in der Schweiz

A.
X. [redacted] wird vorgeworfen, am 29. September 2014 auf einem Parkplatz in A. [redacted] linke Gesichtshälfte geschlagen zu haben, worauf dieser zu Boden gegangen sei und mit Gütergleis aufgeschlagen habe. Dem fallenden B. [redacted] habe X. [redacted] zusätzlich Faustschlages und des Fusstritts habe B. [redacted] diverse Verletzungen, darunter eine Zahnlockerung, als Folge des Sturzes zudem eine weitere Rissquetschwunde im Gesicht.

B.
Das Bezirksgericht Dielsdorf erklärte X. [redacted] am 17. Februar 2017 der versuchten Freiheitsstrafe von 18 Monaten und einer Busse von Fr. 1'000.--. Im Zivilpunkt stellte d B. [redacted] schadenersatzpflichtig ist. Überdies verpflichtete es X. [redacted], B. [redacted] X. [redacted] Berufung und die Staatsanwaltschaft Anschlussberufung.

Beispiel eines publizierten Urteils (entscheidsuche.ch)

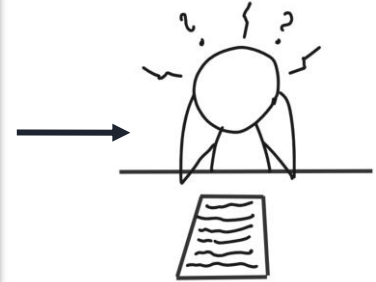
Was ist das Problem?

Heute

Sachverhalt:

A.
X. _____ wird vorgeworfen, am 29. September 2014 auf einem Park
linke Gesichtshälfte geschlagen zu haben, worauf dieser zu Boden gegar
Gütergleis aufgeschlagen habe. Dem fallenden B. _____ habe X. _____
Faustschlages und des Fusstritts habe B. _____ diverse Verletzungen,
Zahnlockerung, als Folge des Sturzes zudem eine weitere Rissquetschw

B.
Das Bezirksgericht Dielsdorf erklärte X. _____ am 17. Februar 2017
Freiheitsstrafe von 18 Monaten und einer Busse von Fr. 1'000.--. Im Ziv
B. _____ schadenersatzpflichtig ist. Überdies verpflichtete es X. _____
X. _____ Berufung und die Staatsanwaltschaft Anschlussberufung.

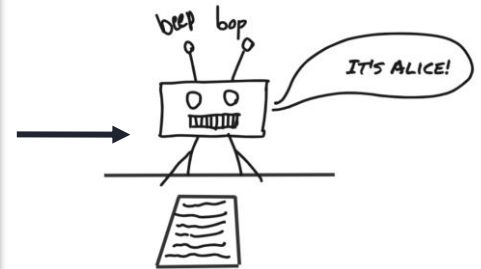


Zukunft?

Sachverhalt:

A.
X. _____ wird vorgeworfen, am 29. September 2014 auf einem Park
linke Gesichtshälfte geschlagen zu haben, worauf dieser zu Boden gegar
Gütergleis aufgeschlagen habe. Dem fallenden B. _____ habe X. _____
Faustschlages und des Fusstritts habe B. _____ diverse Verletzungen,
Zahnlockerung, als Folge des Sturzes zudem eine weitere Rissquetschw

B.
Das Bezirksgericht Dielsdorf erklärte X. _____ am 17. Februar 2017
Freiheitsstrafe von 18 Monaten und einer Busse von Fr. 1'000.--. Im Ziv
B. _____ schadenersatzpflichtig ist. Überdies verpflichtete es X. _____
X. _____ Berufung und die Staatsanwaltschaft Anschlussberufung.



Forschungsfragen

1. **Wie gut können Sprachmodelle Re-Identifikation?**
2. **Was beeinflusst die Performance der Sprachmodelle?**
3. **Was heisst das für unsere Privatsphäre?**

Wie wir Urteile re-identifizieren können

Quiz: Wer wurde hier verurteilt?

Berufung gegen das Urteil des Kantonsgerichts Nidwalden, Zivilabteilung, Grosse Kammer I, vom 7. November 2001.

Sachverhalt:

A. _____ (Kläger), wohnhaft in Y. _____, tritt unter dem Pseudonym DJ Bobo weltweit als Popstar im Fernsehen sowie in Konzert-Tourneen auf. Ausserdem verkauft er Tonträger, insbesondere auch von ihm selbst geschaffene CD's sowie Zubehör von Discjockeys. Er ist Inhaber mehrerer national und international hinterlegter Wort- und Wort-/Bildmarken mit den Zeichen "D.J. Bobo" und "Bobo". Am 17. Juli 1996 liess er in der Schweiz den Internet-Domain-Namen "www.djbobo.ch" registrieren.

Wie wir Urteile re-identifizieren können

Quiz: Wer wurde hier verurteilt?

Berufung gegen das Urteil des Kantonsgerichts Nidwalden, Zivilabteilung, Grosse Kammer I, vom 7. November 2001.

Sachverhalt:

A.
A. _____ (Kläger), wohnhaft in Y. _____, tritt unter dem Pseudonym DJ Bobo weltweit als Popstar im Fernsehen sowie in Konzert-Tourneen auf. Ausserdem verkauft er Tonträger, insbesondere auch von ihm selbst geschaffene CD's sowie Zubehör von Discjockeys. Er ist Inhaber mehrerer national und international hinterlegter Wort- und Wort-/Bildmarken mit den Zeichen "D.J. Bobo" und "Bobo". Am 17. Juli 1996 liess er in der Schweiz den Internet-Domain-Namen "www.djbobo.ch" registrieren.

Wie wir Urteile re-identifizieren können

Search results for "dj bobo gericht internet adresse".

Ungefähr 291'000 Ergebnisse (0.55 Sekunden)

<https://www.nzz.ch> > ...
DJ BoBo gewinnt Streit um Internet-Adresse | NZZ
12.11.2002 — (ap) Der Schweizer Popstar **DJ BoBo** hat den Streit um die **Internetadresse** «www.djbobo.de» gegen seinen ehemaligen Produzenten gewonnen.

<https://www.bger.ch> > ext > live > php > aza > http > ...
4C.141/2002 07.11.2002 - Schweizerisches Bundesgericht
Berufung gegen das **Urteil** des **Kantonsgerichts** Nidwalden, Zivilabteilung, ... Oktober 1999, "den Link zwischen der **Internet-Adresse** www.djbobo.de und der ..."

DJ BoBo privat: Kinder, Frau, Karriere! SO lebt der "Sing meinen Song"-Star heute

Der Weg zum Erfolg war für René Peter Baumann alias DJ BoBo kein einfacher, doch heute gehört er zu den erfolgreichsten deutschsprachigen Musikern. Wie die Show-Legende heute mit Frau und Kindern lebt, das verraten wir Ihnen hier.



DJ Bobo alias René Baumann im Star-Portrait.
Bild: dpa

Und wie sollen Sprachmodelle das schaffen?

1. Lernen von Informationen im Training mit Zeitungsartikeln
2. Das gelernte Wissen zur Re-Identifikation nutzen

Simple, not true?

Simpel, nicht wahr?

Glücklicherweise nicht.

- Medienmitteilungen zum Fall?
- Sind darin genügend Informationen?
- Hat das Urteil differenzierbare Datenpunkte für die Verknüpfung?
 - Strafmass, Namen, Stichwörter (medizinisch, technische Begriffe, ...)
- Kann das Sprachmodell die Datenpunkte auch verknüpfen?

Methoden

3 Datensets

1. Gerichtsurteile

- 7'600 Urteile vom Bundesgericht
- Nur 2019

2. Gerichtsurteile mit News

- 7 händisch re-identifizierte Urteile
- Passende News von Swissdox.ch

3. Wikipedia

- 70'000 biographische Artikel
- Min. 2'250 Wörter / Artikel
- Paraphrasiert und Personen maskiert

"<mask> (; February 12, 1809 - April 15, 1865) was an American lawyer and statesman who served as the 16th president of the United States from 1861 until his assassination in 1865. <mask> led the nation through the American Civil War and succeeded in preserving the Union, abolishing slavery, bolstering the federal government, and modernizing the U.S. economy. <mask> was born into poverty in a log cabin in Kentucky and was raised on the frontier primarily in Indiana.

Example extract for Abraham Lincoln

Metriken

Was gilt als re-identifiziert?

1. Partial Name Match Score (**PNMS**)
2. Weighted PNMS (**W-PNMS**)
3. Normalized Levenshtein Distance (**NLD**)

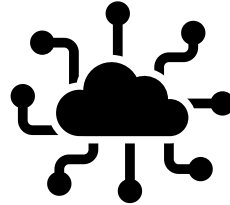
*Max **Muster** passt zu Hans **Muster** mit 65%*

Re-Identifizierungsmethoden: direkt

Wer ist Mr. X? Hier ist das Urteil: „....“

Wer ist <mask>?

<mask> (Kläger), wohnhaft in Y. _____, tritt unter dem Pseudonym **DJ Bobo** weltweit als Popstar im Fernsehen sowie in Konzert-Tourneen auf. Ausserdem verkauft er Tonträger, insbesondere auch von ihm selbst geschaffene CD's sowie Zubehör von Discjockeys.

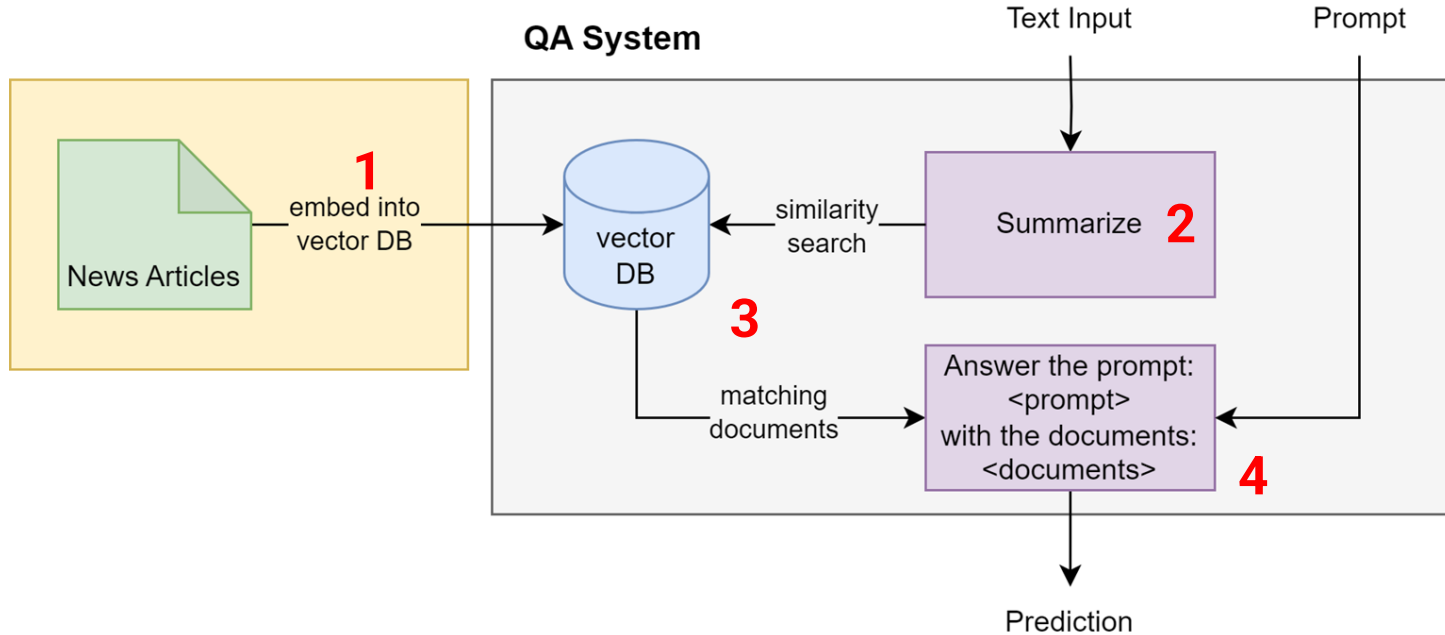


Sprachmodell

Das ist René Baumann.

Re-Identifizierungsmethoden: inklusive News

Retrieval Augmented Generation (RAG)



Re-Identifizierungsmethoden: inklusive News

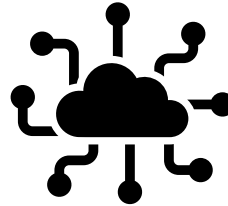
Wer ist Mr. X? Hier ist das Urteil und passende News-Artikel: „...“

Wer ist <mask>?

<mask> (Kläger), wohnhaft in Y._____, tritt unter dem Pseudonym **DJ Bobo** weltweit als Popstar im Fernsehen sowie in Konzert-Tourneen auf. Ausserdem verkauft er Tonträger, insbesondere auch von ihm selbst geschaffene CD's sowie Zubehör von Discjockeys.

Aktuell: René Baumann, bekannt als DJ Bobo, in Rechtsstreit um Domäne verwickelt!

DJ Bobo Skandal! Der Sänger ging für djbobo.ch vor Gericht.



Sprachmodell

Das ist René Baumann.



Resultate

Performance auf Wikipedia

- 54 Sprachmodelle evaluiert

71% re-identifiziert!

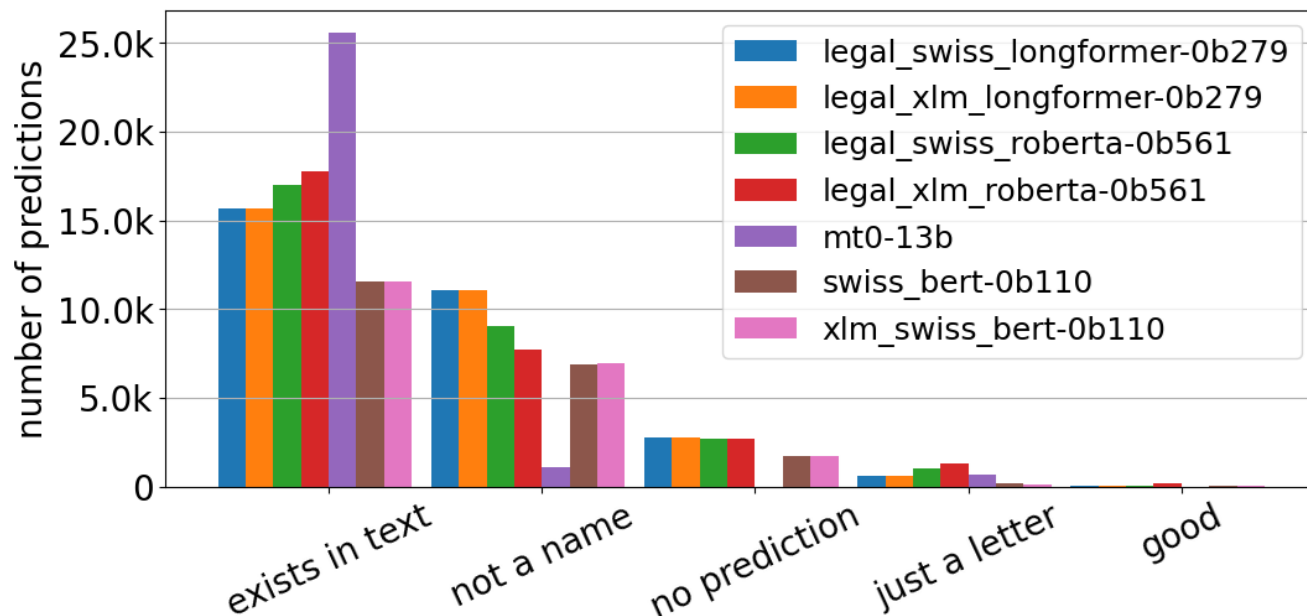
Model	Size [B]	PNMS \uparrow	NLD \downarrow	W-PNMS \uparrow
GPT-4	1800	0.71	0.17	0.65
GPT-3.5	175	0.52	0.23	0.46
mT0	13	0.37	0.42	0.31
Flan_T5	11	0.37	0.45	0.30
incite	3	0.37	0.53	0.30
Flan_T5	3	0.35	0.48	0.29
BLOOMZ	7.1	0.34	0.45	0.29
T0	11	0.34	0.45	0.28

Table 1: Models w/ W-PNMS > 0.28 on Wikipedia dataset

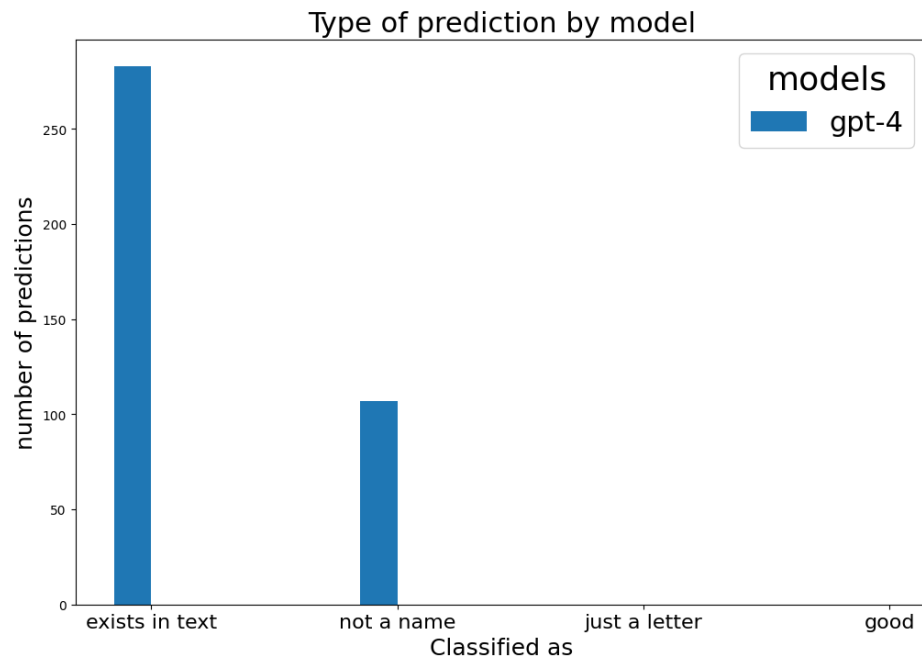
Re-Identifikation auf Sammlung von Urteilen

- 1 / 7'600 Urteilen teilweise re-identifiziert

0.014% re-identifiziert!



GPT-4 verweigert die Re-Identifikation



Re-Identifikation auf händisch gesammelten Urteilen

Retrieval Augmented Generation:

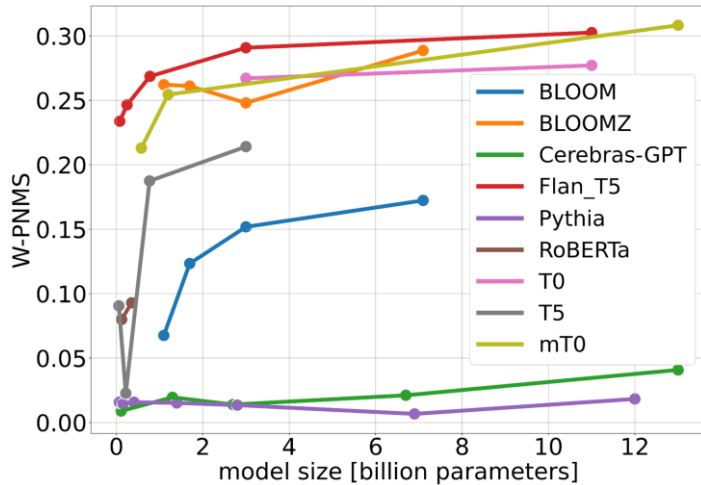
- full name for 1 example

GPT-3.5-turbo-16k	4 / 7
GPT-4	5 / 7

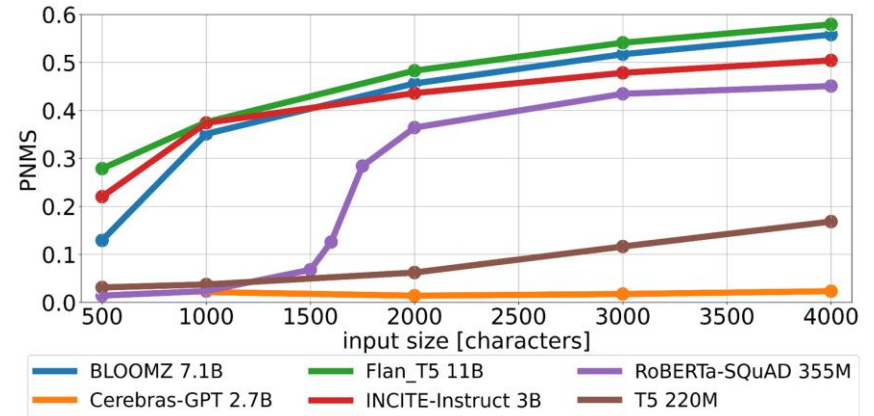
5 von 7 re-identifiziert!

Beeinflussende Faktoren der Re-Identifikation

Grösse der Modelle

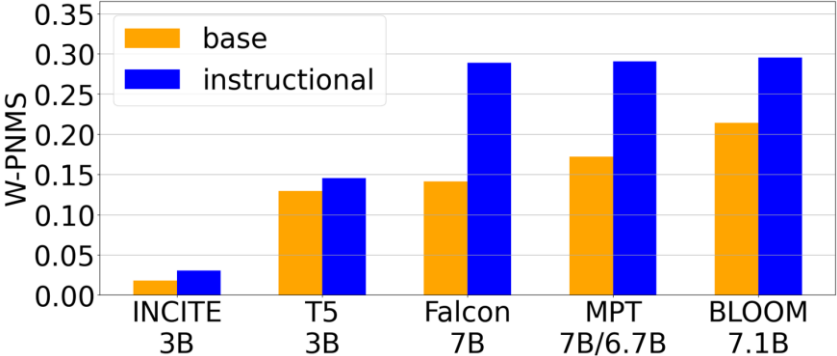


Länge vom genutzten Textabschnitt



Influential Factors on Re-identification

Instruction tuning



Fazit

1. Können Sprachmodelle re-identifizieren?

Gerichtsurteile



Wikipedia



2. Beeinflussende Faktoren?

*Modell-Grösse, Textlänge,
Instruction Tuning*

3. Privatsphäre gefährdet?

*Anonymität ist (noch) nicht in
Gefahr*