

# NLP Research Workshop

## NLP-Anwendung in NFP73 und NFP77

**PD Dr. Matthias Stürmer**  
**Dominic Schweizer**  
**Joel Niklaus**

Forschungsstelle Digitale Nachhaltigkeit  
der Universität Bern

Freitag, 14. August 2020  
Universität Bern



**Nachhaltige Wirtschaft**  
Nationales Forschungsprogramm



**Digitale Transformation**  
Nationales Forschungsprogramm

## NLP-Aktivitäten der FDN

Die Forschungsstelle Digitale Nachhaltigkeit beschäftigt sich im Bereich NLP aktuell mit folgenden Projekten:

- **NFP73:** Datenveredelung und Texterkennung
- **NFP77:** Topic Modelling, Named Entity Recognition etc.
- **Kanton Bern:** Chatbot-Analyse, Natural Language Understanding, Question & Answering Systems, Voicerecognition
- **SECO/ALV:** Matching zwischen offenen Stellen und Stellensuchen

# NFP73 «Nachhaltigkeit bei Beschaffungen»



Università della Svizzera Italiana



<sup>b</sup>  
UNIVERSITÄT  
BERN

University of Berne



**Prof. Dr. Dr. Peter Seele**  
CSR and business ethics professor



**Prof. Dr. Federica de Rossa**  
Law professor and judge at the Federal Supreme Court



**Sebastian Knebel**  
Corporate communication



**Clarissa David**  
Lawyer



**Dr. Matthias Stürmer**  
Computer scientist



**Tobias Welz**  
Industrial engineer

# SNSF NRP73 on Sustainable Procurement

- **Research question:**  
«How can public procurement become more sustainable?»
  - Public procurement accounts for CHF 41 bn spendings annually
  - Ecological and social criteria may make a big difference
  - In addition new procurement law in 2021 requires sustainability
- **Goal:** measuring sustainability criteria in public tender documents
- **Data and method:** crawling and refining of procurement data



**Nachhaltige Wirtschaft**  
Nationales Forschungsprogramm

# 1. Beschaffungsstatistik.ch

- *NLP light*
- Simap.ch Publikationen als unstrukturierte Daten
- Informationen aus Publikation extrahieren

12.08.2020 | Projekt-ID 194015 | Meldungsnummer 1097745 | Ausschreibungen

## Ausschreibung

Publikationsdatum Simap: 12.08.2020

### 1. Auftraggeber

#### 1.1 Offizieller Name und Adresse des Auftraggebers

Bedarfsstelle/Vergabestelle: Post CH AG / Postmail / ganzer Konzern Post

Beschaffungsstelle/Organisator: Post CH AG /

F41 / Beschaffung, zu Hdn. von CC-WTO - Projekt «Taxi», Wankdorfallee 4, 3030 Bern, Schweiz, E-Mail: [wto-taxi-dl@post.ch](mailto:wto-taxi-dl@post.ch)

#### 1.2 Angebote sind an folgende Adresse zu schicken

Adresse gemäss Kapitel 1.1

#### 1.3 Gewünschter Termin für schriftliche Fragen

26.08.2020

**Bemerkungen:** Fragen des Anbieters zur Ausschreibung und zu den Ausschreibungsunterlagen sind direkt in der Ausschreibungsplattform simap.ch bis spätestens 26.08.2020 zu stellen.

Es werden keine telefonischen oder mündlichen Auskünfte erteilt. Die Beantwortung der Fragen erfolgt direkt im Frageforum in der Ausschreibungsplattform simap.ch bis am 07.09.2020. Verspätet eintreffende Fragen werden nicht mehr beantwortet.

#### 1.4 Frist für die Einreichung des Angebotes

**Datum:** 28.09.2020 **Uhrzeit:** 12:00, **Spezifische Fristen und Formvorschriften:** Das Angebot muss zur Wahrung seiner Rechtzeitigkeit zum genannten Zeitpunkt vollständig am genannten Ort eingetroffen sein; Eingaben per E-Mail oder Fax genügen nicht. Die Aufgabe eines Angebotes per Post zum genannten Zeitpunkt (Poststempel) reicht zur Wahrung der Rechtzeitigkeit nicht. Zu spät eingetroffene Angebote werden nicht berücksichtigt.

#### 1.5 Datum der Offertöffnung:

28.09.2020, **Uhrzeit:** 12:00, **Bemerkungen:** Die Öffnung der Angebote ist nicht öffentlich.

#### 1.6 Art des Auftraggebers

Bund (Dezentrale Bundesverwaltung - öffentlich rechtliche Organisationen)

#### 1.7 Verfahrensart

Offenes Verfahren

#### 1.8 Auftragsart

Dienstleistungsauftrag

#### 1.9 Gemäss GATT/WTO-Abkommen, resp. Staatsvertrag

Ja

## 2. Beschaffungsobjekt

### 2.1 Dienstleistungskategorie CPC:

[27] Sonstige Dienstleistungen

### 2.2 Projekttitel der Beschaffung

Taxi DL (Taxi Dienstleistungen)

# 1. Beschaffungsstatistik.ch

- Auftraggeber identifizieren
- Struktur erkennen
- Informationen aus String extrahieren
- Auftraggeber mit Trainingsdaten vergleichen

12.08.2020 | Projekt-ID 194015 | Meldungsnummer 1097745 | Ausschreibungen

## Ausschreibung

Publikationsdatum Simap: 12.08.2020

### 1. Auftraggeber

#### 1.1 Offizieller Name und Adresse des Auftraggebers

Bedarfsstelle/Vergabestelle: Post CH AG / Postmail / ganzer Konzern Post

Beschaffungsstelle/Organisator: Post CH AG /

F41 / Beschaffung, zu Hdn. von CC-WTO - Projekt «Taxi», Wankdorfallee 4, 3030 Bern, Schweiz, E-Mail: [wto-taxi-dl@post.ch](mailto:wto-taxi-dl@post.ch)

#### 1.2 Angebote sind an folgende Adresse zu schicken

Adresse gemäss Kapitel 1.1

#### 1.3 Gewünschter Termin für schriftliche Fragen

26.08.2020

**Bemerkungen:** Fragen des Anbieters zur Ausschreibung und zu den Ausschreibungsunterlagen sind direkt in der Ausschreibungsplattform simap.ch bis spätestens 26.08.2020 zu stellen.

Es werden keine telefonischen oder mündlichen Auskünfte erteilt. Die Beantwortung der Fragen erfolgt direkt im Frageforum in der Ausschreibungsplattform simap.ch bis am 07.09.2020. Verspätet eintreffende Fragen werden nicht mehr beantwortet.

#### 1.4 Frist für die Einreichung des Angebotes

**Datum:** 28.09.2020 **Uhrzeit:** 12:00, **Spezifische Fristen und Formvorschriften:** Das Angebot muss zur Wahrung seiner Rechtzeitigkeit zum genannten Zeitpunkt vollständig am genannten Ort eingetroffen sein; Eingaben per E-Mail oder Fax genügen nicht. Die Aufgabe eines Angebotes per Post zum genannten Zeitpunkt (Poststempel) reicht zur Wahrung der Rechtzeitigkeit nicht. Zu spät eingetroffene Angebote werden nicht berücksichtigt.

#### 1.5 Datum der Offertöffnung:

28.09.2020, **Uhrzeit:** 12:00, **Bemerkungen:** Die Öffnung der Angebote ist nicht öffentlich.

#### 1.6 Art des Auftraggebers

Bund (Dezentrale Bundesverwaltung - öffentlich rechtliche Organisationen)

#### 1.7 Verfahrensart

Offenes Verfahren

#### 1.8 Auftragsart

Dienstleistungsauftrag

#### 1.9 Gemäss GATT/WTO-Abkommen, resp. Staatsvertrag

Ja

## 2. Beschaffungsobjekt

### 2.1 Dienstleistungskategorie CPC:

[27] Sonstige Dienstleistungen

### 2.2 Projekttitel der Beschaffung

Taxi DL (Taxi Dienstleistungen)

# 1. Beschaffungsstatistik.ch

- Auftrag der Post CH AG
- Beschaffungsstelle:
  - Post CH AG / F41 / Beschaffung, Wankdorfallee 4 3030 Bern

## 1. Auftraggeber

### 1.1 Offizieller Name und Adresse des Auftraggebers

Bedarfsstelle/Vergabestelle: Post CH AG / Postmail / ganzer Konzern Post

Beschaffungsstelle/Organisator: Post CH AG /

F41 / Beschaffung zu Hdn. von CC-WTO - Projekt «Taxi» Wankdorfallee 4 3030 Bern Schweiz E-Mail: [wto-taxi-dl@post.ch](mailto:wto-taxi-dl@post.ch)

# 1. Beschaffungsstatistik.ch

Post CH AG / Postmail / ganzer Konzern  
Post =

### Nicht eingeordnete Bedarfsstellen

Post CH AG / Postmail / ganzer Konzern Post [=>Google Suche](#)

Beschaffungsstelle Name: Post CH AG / F41 / Beschaffung

Beschaffungsstelle Strasse: Wankdorfallee 4

Auftraggeber Art: Bund (Dezentrale Bundesverwaltung - öffentlich rechtliche Organisationen)

Auftraggeber ID: 57022

Meldungsnummer: [1097745](#)

Projekt ID: [194015](#)

Anzahl betroffene Projekte: 2

Gemeindeverwaltung Eichberg SELECT

---

Chemins de fer fédéraux suisses CFF, Infrastructure, Projets, I-AEP-PJM-RWT-T1 SELECT

### Zuordnung

Die Schweizerische Post AG Post CH AG / PostMail	63% Wankdorfallee 4, Bern	<span style="border: 1px solid #ccc; padding: 2px 10px; background-color: #444; color: white;">OK</span>
Die Schweizerische Post AG Post CH AG / PostMail	63% Wankdorfallee 4, Bern	<span style="border: 1px solid #ccc; padding: 2px 10px; background-color: #444; color: white;">OK</span>
Die Schweizerische Post AG Post CH AG / PostMail	63% Wankdorfallee 4, Bern	<span style="border: 1px solid #ccc; padding: 2px 10px; background-color: #444; color: white;">OK</span>

Suchen

OK

Neu eingeben

NEW Post CH AG / Postmail / ganzer Konzern Post OK

# 1. Beschaffungsstatistik.ch

## Ausschreibung

(20 Einträge von 78'138 angezeigt)

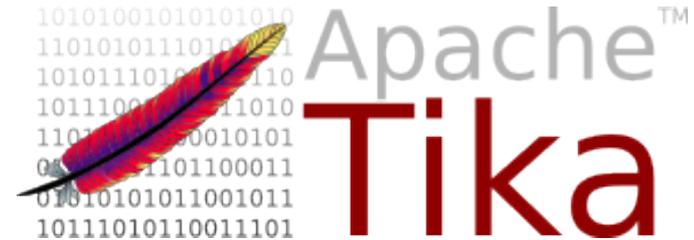
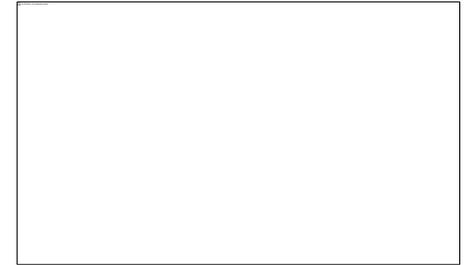
Volltextsuche in Tabelle Ausschreibung



▼ Datum	Projekt	Meldung	Projekttitel	IT Projekt	CT projekt	Verfahrensart	Auftraggeber	Auftraggeber-Art
2020-08-12	<a href="#">194015</a>	<a href="#">1097745</a>	Taxi DL (Taxi Dienstleistungen)	Nein	Nein	Offenes Verfahren	<a href="#">Die Schweizerische Post AG</a>	Bund (Dezentrale Bundesverwaltung - öffentlich rechtliche Organisationen)
2020-08-12	<a href="#">204584</a>	<a href="#">1135925</a>	KHK St. Gallen, Erneuerung der Rauchgasreinigung	Nein	Nein	Offenes Verfahren	<a href="#">Stadt St. Gallen</a>	Gemeinde/Stadt
2020-08-12	<a href="#">206995</a>	<a href="#">1145087</a>	IWB, Fachplaner GP Bau Retrofit	Nein	Nein	Offenes Verfahren	<a href="#">IWB - Industrielle Werke Basel</a>	Andere Träger kantonaler Aufgaben
2020-08-12	<a href="#">207883</a>	<a href="#">1148399</a>	Umbau und Sanierung Wirtschaftsgymnasium, BKP 283.2 Deckenverkleidung aus Gips	Nein	Nein	Offenes Verfahren	<a href="#">Kanton Basel-Stadt</a>	Kanton
2020-08-12	<a href="#">207934</a>	<a href="#">1148679</a>	Appel d'offres pour des tunnels de lavage, désinfection à la vapeur et séchage de lits hospitaliers	Nein	Nein	Offenes Verfahren	<a href="#">Universitätsspital Lausanne CHUV</a>	Kanton

## 2. Intelliprocare.ch

- Grösserer Datensatz
  - 490'991 indexierte Ausschreibungsdokumente
  - 171'297 Simap Publikationen
- Durchsuchen der Unterlagen
  - Volltextsuche



## 2. Intelliprocedure

SUCHEN 

Nachhaltigkeit

Total Resultate:	6738	Anzahl indexierte Projekte:	19'495
Zeitdauer der Suche:	0.184 Sekunden	Total indexierte Ordner:	137'679
		Total indexierte Dokumente:	490'991 
		Grösse der Daten:	1'665.132 GB
		Neuste Ausschreibung:	13.08.2020

0-10/6738

ZURÜCK

WEITER

Erneuerung und Erweiterung Campus Horw Das Verfahren besteht aus einer öffentlich ausgeschriebenem Präqualifikation (Phase 1) und einem nachfolgenden, einstufigen und nicht anonymen Studienauftrag (Phase 2) für 4 bis 5 Generalplanerteams. [.pdf](#)

Projekt ID: 188232 | Datum: 01.06.2019

Auftraggeber: Kanton Luzern (Simap:Kanton Luzern / Finanzdepartement)

... ja / nein Bauherr Adresse, PLZ Ort Kurzbeschreibung Adresse, PLZ Ort Name Adresse, PLZ Ort Funktion Beschreibung in Stichworten: - Projektbeschreibung - Innovationen bezüglich **Nachhaltigkeit**, Raumgestaltung, Gemeinschaft und Unterricht Referenz 2 Vorname, Name im Betrieb tätig seit Funktion Referenz 1 Adresse, PLZ Ort Name Adresse, PLZ Ort Funktion Beschreibung in Stichworten:...

Datei Score: 62.045002 | Dateiname: | Grösse: 0.25 MB

Ausschreibung Anzahl Dokumente: 4 | Anzahl Ordner: 0 | Grösse: 2.959 MB

-  In Beschaffungsstatistik ansehen
-  Auf Simap.ch ansehen
-  Dokument ansehen
-  Alle Dokumente zu diesem Projekt ansehen
-  Projekt als Zip herunterladen

## 2. Intelliprocure.ch: Ausblick

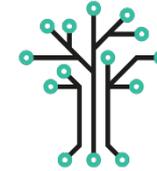
- Identifizierung von EK & ZK
  - Topic Modelling mit Random Forest & Active Learning
- Erstellung eines Kataloges für EK & ZK

Rang	Kombiniertes Modell		Einfaches Modell	
	Wort	Quelle	Wort	Quelle
1	zuschlagskriterium	Überschrift	zuschlagskriterium	Überschrift
2	zk	Überschrift	preis	Text
3	gewichtung	Text	zuschlagskriterium	Text
4	zuschlagskriterium	Text	zk	Text
5	angebot	Text	punkt	Text
6	punkt	Text	gewichtung	Text
7	zuschlagskriterium	Eltern-Überschrift	angebot	Text
8	schlüsselperson	Text	nr	Text
9	qualität	Text	angabe	Text
10	preis	Text	table	Text

# NFP77 «Open Justice vs. Privacy»



Institut für  
Wirtschaftsinformatik



Forschungsstelle Digitale  
Nachhaltigkeit am  
Institut für Informatik



**Prof. Dr. Andreas Lienhard**  
Rechtsprofessor und  
Dekan Rechtswissenschaftliche Fakultät



**Daniel Kettiger**  
Rechtsanwalt



**Nathalie Schwager**  
Doktorandin  
Rechtswissenschaften



**Prof. Dr. Thomas Myrach**  
Professor für  
Wirtschaftsinformatik



**PD Dr. Matthias Stürmer**  
Leiter Forschungsstelle  
Digitale Nachhaltigkeit



**Joel Niklaus**  
Doktorand  
Informatik

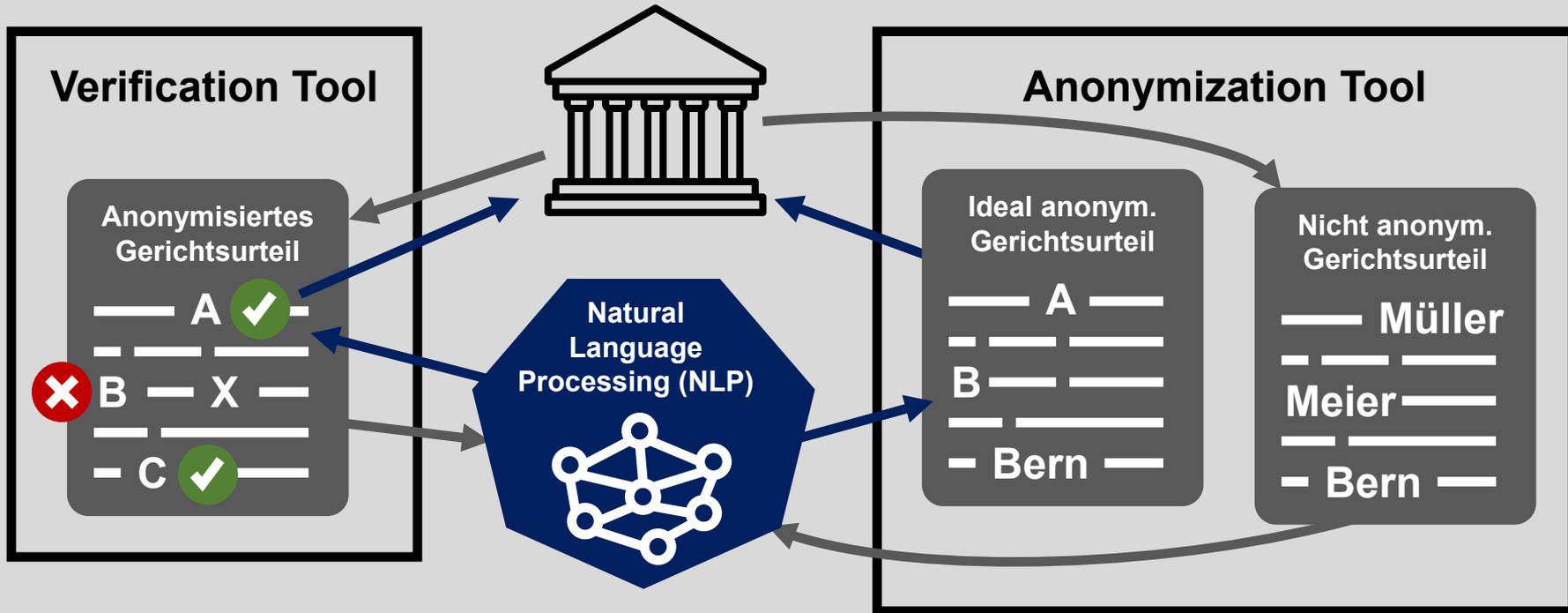
# SNSF NRP77 on Open Justice vs. Privacy

- **Research question:** «How can court decisions be released as open data without compromising privacy?»
  - Transparency in justice expects court decisions to be published openly
  - Usually plaintiffs in a court decision are pseudonymized (e.g. «A vs. B»)
  - But big data and artificial intelligence may deanonymize court decisions
- **Goal:** conduct ethical deanonymization for improving anonymization
- **Data and method:** crawling and NLP analysis of Swiss court data



**Digitale Transformation**  
Nationales Forschungsprogramm

# Verification and Anonymization Tool



# Plan

1. Datenset vorbereiten mit Entscheidsuche.ch
2. Topic Modeling für Gerichtsurteile (wichtig für Linkage)
3. Externe Daten für Linkage crawlen (bspw Zeitungen, Social Media)
4. Ein anonymisiertes Gerichtsurteil auswählen und re-identifizieren
5. Anonymisierungssystem (mit NER) für ein Gerichtsurteil bauen
6. Die beiden Systeme für eine Kategorie (topic) ausbauen
7. Die beiden Systeme für alle Kategorien ausbauen

# 1. Aktueller Stand der Datensammlung

- Im Gespräch für Kooperation mit Verein **entsuchedsuche.ch**
- Hat 2018 rund **500'000 Gerichtsurteile** aus rund 80 Gerichtsdatenbanken gecrawlt
- **Wiederkehrendes Crawling** durch neue Programmierung der Scraper

Search results for "beschaffung" (7744 results found).

**Filters:**

- Jahr:** 2000 - 2018
- Monat:** 1 - 12
- Sprache:**
  - de: 5397
  - fr: 47
  - it: 8
- Level:**
  - bund: 5874
  - kantone: 1870
- Kanton:**
  - bund: 5874
  - zh: 734
  - gr: 255
  - bs: 149
  - sg: 139
  - bi: 90
  - lu: 83
  - ti: 79
  - be: 74
  - ag: 63
  - so: 44
  - tg: 32
  - fr: 24
  - vs: 23
  - sh: 22
  - zg: 15
  - ne: 14
  - ur: 11
  - gl: 9
  - ar: 3
  - sz: 2

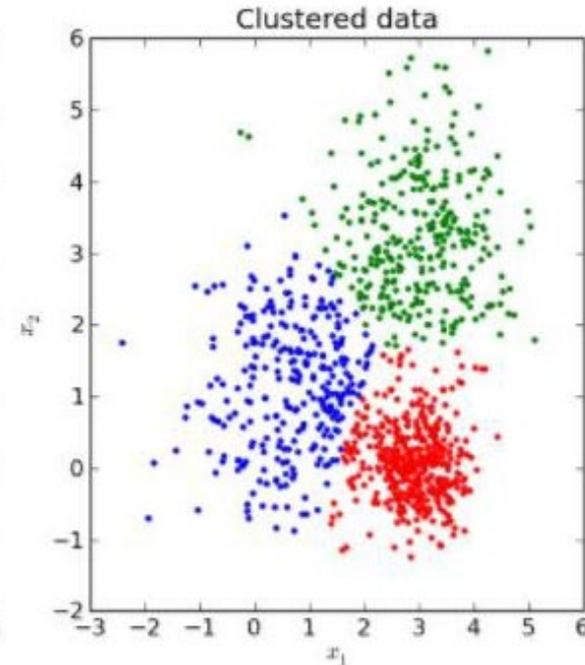
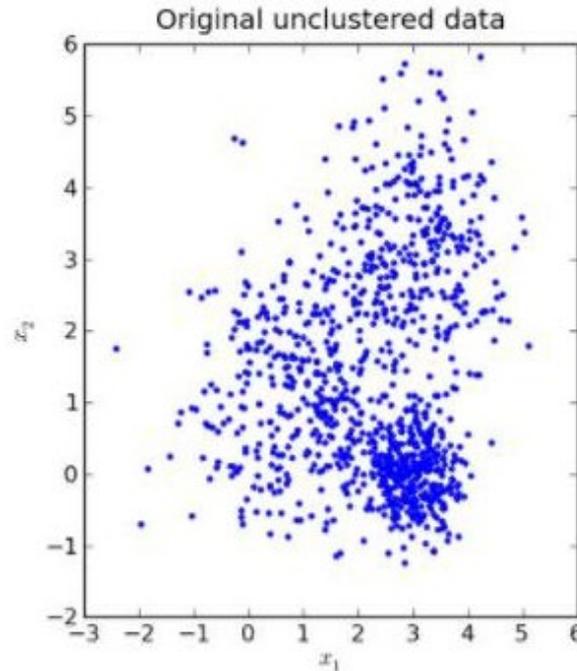
**Results:**

- Result 1:** vs - vs, KGVS-20170209-A1-16-234-20170612-A73.pdf. Text: treffen: 3. es sei der Beschwerdeführerin umfassende Akteneinsicht in die amtlichen Akten des Beschaffungsverfahrens wiederholt werden darf (ZWR 2000 S. 52; Peter Galli/Daniel Lehmann/Peter Rechsteiner, Das öffentliche Beschaffung. Für die Beschaffungen hält die IVöB fest, dass Verfahrensabbruch und Wiederholung des Verfahrens nur brauch vor, der nicht zu schützen ist. 4.2 Provisorischer Natur ist jener Abbruch, bei dem die Beschaffungsabsicht Dies führt zu einer wesentlichen Änderung des Beschaffungsgegenstands.
- Result 2:** bund - bvger, Urteil B-536-2013.pdf. Text: Weil damit der strittige Beschaffungsgegenstand entfalle und insofern auch eine Diskriminierung von hätten die im Beschaffungsprojekt involvierten Bundesämter zum weiteren Vorgehen noch keinen Beschluss "In-House"-Geschäft), was die Anwendbarkeit des Beschaffungsrechts ausschliesse. Eigen- beschaffung), jedoch unter definitivem Verzicht auf den Bezug der im (abgebrochenen) Beschaffungsverfahren Lässt sich somit der am (...) 2013 verfügte Abbruch des Beschaffungsver- fahrens im Projekt X.....
- Result 3:** bund - bvger, B-6177-2008.pdf. Text: In diesem Zusammenhang weisen sie darauf hin, dass die EU-Beschaffungsrichtlinien, welche den



## 2. Topic Modeling

- Unsupervised oder Supervised (Topic Classification)
- Methoden:
  - Latent Dirichlet Allocation (LDA)
  - Latent Semantic Analysis (LSA)



## 2. Topic Modeling

car, power, light,  
drive, mount,  
controller,  
cool, engine, back,  
turn

**AUTOMOBILE**

patient,  
study, slave,  
wing, disease, food,  
eat, pain, treatment,  
syndrome

**MEDICAL**

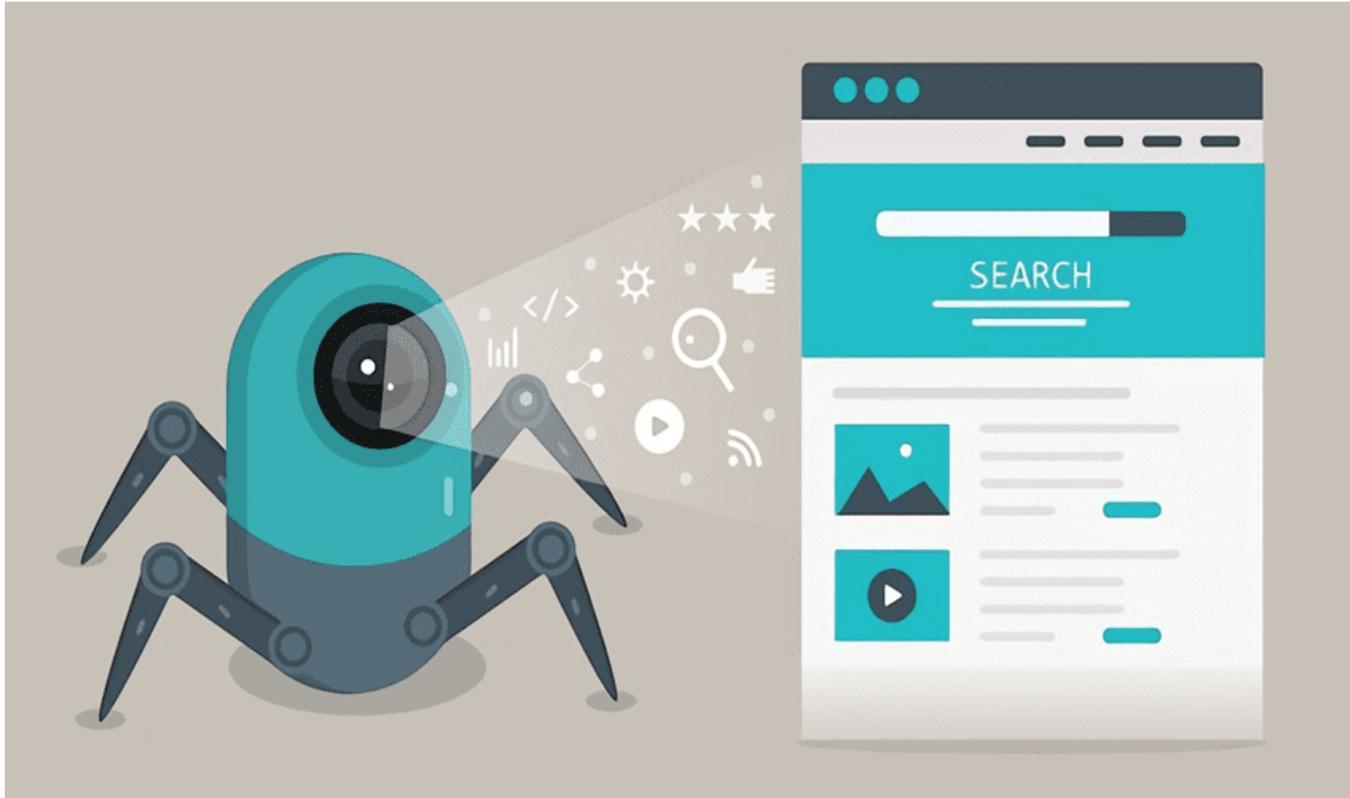
god, evidence,  
christian,  
believe, reason, faith,  
exist, bible, religion,  
claim

**RELIGION**

game, team, year,  
play, win, good, season  
fan, run, score

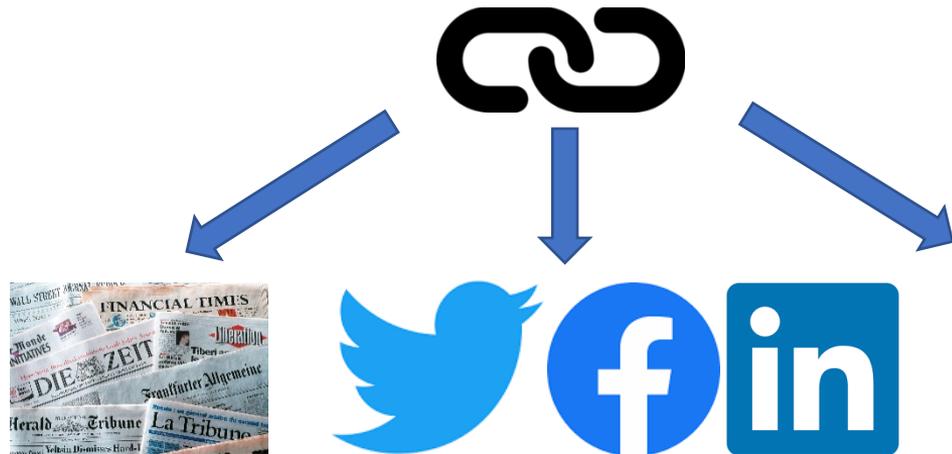
**SPORTS**

## 3. Externe Daten Crawlen



## 4. Re-Identifizierung

- Linkage mit externen Informationsquellen (Knowledge Base) mit
  - Identifikatoren
  - Quasi-Identifikatoren (nur zusammen mit anderen Informationen aussagekräftig)
- Entity Linking (NER + NED)



## 5. Anonymisierung: Named Entity Recognition

In fact, the Chinese NORP market has the three CARDINAL most influential names of the retail and tech space – Alibaba GPE, Baidu ORG, and Tencent PERSON (collectively touted as BAT ORG), and is betting big in the global AI GPE in retail industry space. The three CARDINAL giants which are claimed to have a cut-throat competition with the U.S. GPE (in terms of resources and capital) are positioning themselves to become the ‘future AI PERSON platforms’. The trio is also expanding in other Asian NORP countries and investing heavily in the U.S. GPE based AI GPE startups to leverage the power of AI GPE. Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one CARDINAL, with an anticipated CAGR PERSON of 45% PERCENT over 2018 - 2024 DATE.

To further elaborate on the geographical trends, North America LOC has procured more than 50% PERCENT of the global share in 2017 DATE and has been leading the regional landscape of AI GPE in the retail market. The U.S. GPE has a significant credit in the regional trends with over 65% PERCENT of investments (including M&As, private equity, and venture capital) in

## 5. Anonymisierung

- Vier Methoden:
  - Suppression  
(Peter → XXX)
  - Tagging  
(Peter → Vorname\_1)
  - Random substitution  
(Peter → Hans)
  - Generalization  
(Peter → Name)



# Diskussion

## Anregende Fragen:

- Bei welchen NLP Techniken ist der state-of-the-art noch nicht gut genug?
- Welche anderen NLP Techniken könnten relevant werden?

