# A Human-in-the-Loop Approach for Information Extraction from Privacy Policies under Data Scarcity

Michael Gebauer, Faraz Maschhur, Nicola Leschke, Elias Grünewald, and Frank Pallas

Information Systems Engineering
TU Berlin

**ISE**ngineering

https://www.nytimes.com/interactive/2019/06/12/opinion/facebook-google-privacy-policies.html

https://www.designboom.com/readers/dima-yarovinsky-visualizes-facebook-instagram-snapchat-terms-of-service-05-07-2018/

?

Browser-Plugins

Visualising Data-Sharing-Networks

Privacy Icons

Grünewald et al. (2023). „Enabling Versatile Privacy Interfaces Using Machine-Readable Transparency Information". In Privacy Symposium 2023

# Indispensable: machine-readable representation

Transparency Information
Language & Toolkit



Browser-Plugins

Visualising Data-Sharing-Networks

Grünewald and Pallas (2021). „TILT: A GDPR-Aligned
Transparency Information Language and Toolkit for
Practical Privacy Engineering ". In ACM FAccT 2021.

Privacy Icons

# How to get TILT?



textual privacy policy → ? → Structured, machine-readable **transparency information**

# General Approach



*TILTer*

**Manual annotation**

textual privacy policy

Structured, machine-readable **transparency information**

**Privacy Icons**
**Chatbot**
**Browser extension**
**Dashboard**
**Disclosure analysis**
**...**

ISEngineering

# TILTer

# Annotation Interface

Hierarchically structured **labels** for:
- Tokens, e.g.
  - DPO (name, e-mail, address, ...)
  - categories
  - purpose
- Sentences, e.g.
  - data subject rights
  - automated decision making
  - changes

Multiple **tasks** per policy:
- Identify labels for a specific hierarchical level

## → still, very tedious work!

# NLP Coming to the Rescue!

# ML-supported Annotation of Privacy Policies

OPP-115 data set (Wilson et al. 2016)

- Polisis (Harkous et al. 2018)

- TLDR (Alabduljabbar et al. 2021)

✓ reliable classification (F1: 0.83/0.91)
✓ improving navigability

✗ data scarcity
✗ creating detailed and accurate machine-readable representations

# Refining the Classification Task

**Annotations** for:
- Tokens, e.g.
  - DPO
  - categories
  - purpose
- Sentences, e.g.
  - data subject rights
  - automated decision making
  - changes

"We therefore would like to inform you about data privacy for websites under *.tu.berlin."

"You can obtain information about the data we have stored about you free of charge at any time without having to give a reason."

"You also have the right to object at any time to the processing of personal data concerning you that is carried out on the basis of Article 6 paragraph 1 lit. e and f EU-GDPR."

https://www.tu.berlin/en/data-protection

# Sentence-based Information Retrieval

set of potential candidates

perfect retrieval distribution

$$S = \left\{ b_i : \max_{\forall b_i \in B} P(c_i = 1 \| b_i, \theta) \right\}$$

paragraphs (blobs)

textual privacy policy

candidate

candidate is data subject right

true parameter

# Sentence-based Information Retrieval

perfect retrieval distribution

$$S = \left\{ b_i: \max_{\forall b_i \in B} \boxed{P(c_i = 1 | b_i, \theta)} \right\}$$

$$\hat{f}(b_i, \hat{\theta})$$

extraction model

ISEngineering

# Extraction Model

Static Word Embeddings
& Naive Bayes

BinaryBERT

SentenceBERT

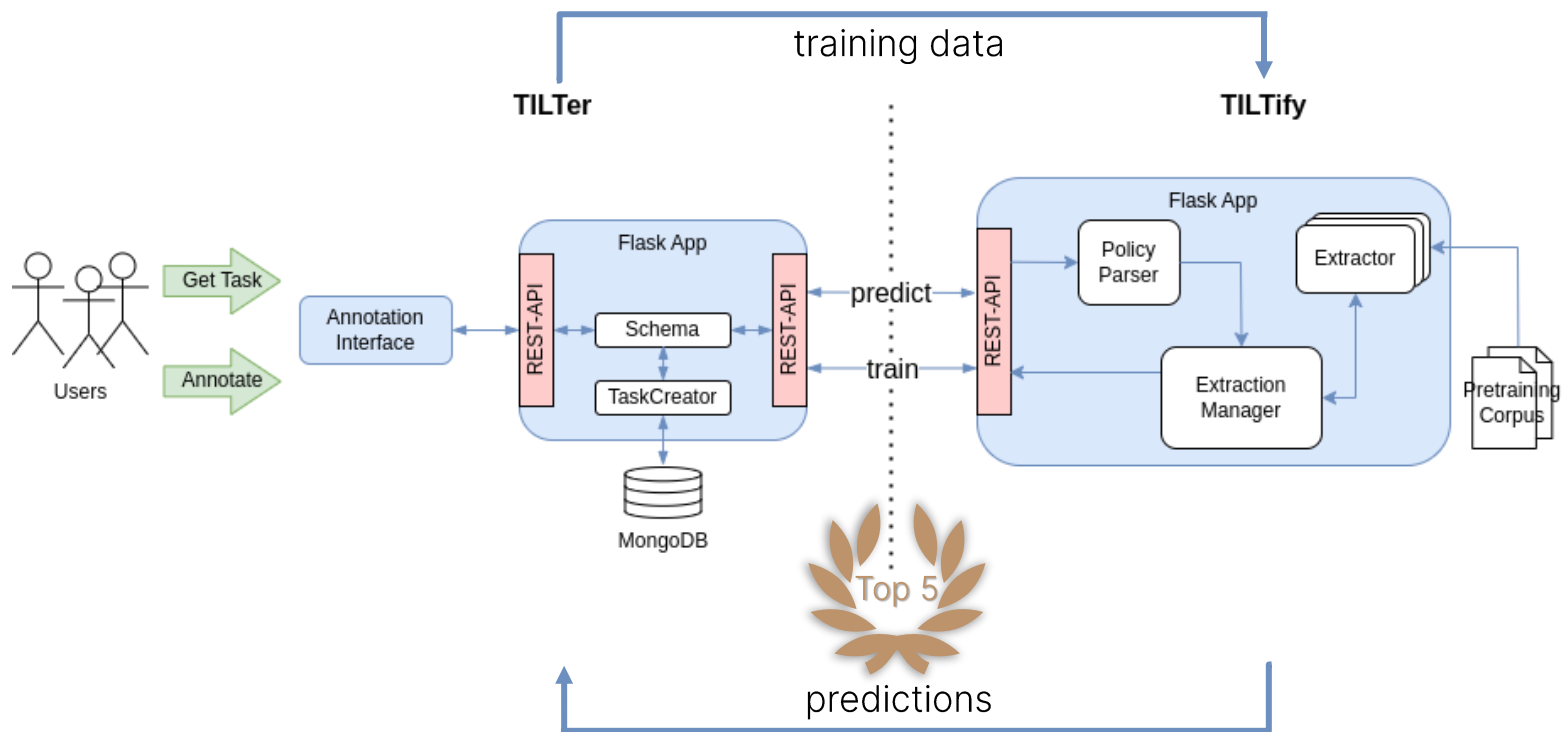# Human-in-the-Loop

60 annotated policies

🇩🇪 German language

🖐 data scarcity

▦  differing amounts of data subject rights

⚖ data imbalance:
>16600 paragraphs
<450  containing data subject rights

https://github.com/DaSKITA/tiltify/tree/main/data/annotated_policies

# Evaluation
## Models

| | Static Word Embeddings & Naive Bayes | BinaryBERT | SentenceBERT |
|---|---|---|---|
| Right to Information | 0.18 | 0.0 | 0.93 |
| Right to Deletion | 0.0 | 0.0 | 0.86 |
| Right to Data Portability | 0.0 | 0.0 | 0.86 |
| Right to Complain | 0.29 | 0.15 | 0.93 |
| Right to Withdraw Consent | 0.0 | 0.08 | 0.90 |

F1 Score (for 5-rank, if not stated otherwise)

# Take-aways

We tackled problem of retrieving machine-readable transparency information from privacy policies under the constraint of data scarcity.



Adopted from
https://datascientest.com/de/bert
and  https://tinyurl.com/2bhe3brk

S-BERT - in conjunction with k-rank approaches - is extremely promising to simplify the annotation of privacy policies, even with few training data.

# A Human-in-the-Loop Approach for Information Extraction from Privacy Policies under Data Scarcity

**Nicola Leschke**

nl@ise.tu-berlin.de

tu.berlin/ise/nl

Nicola Leschke

tu.berlin/ise/daskita

ISEngineering