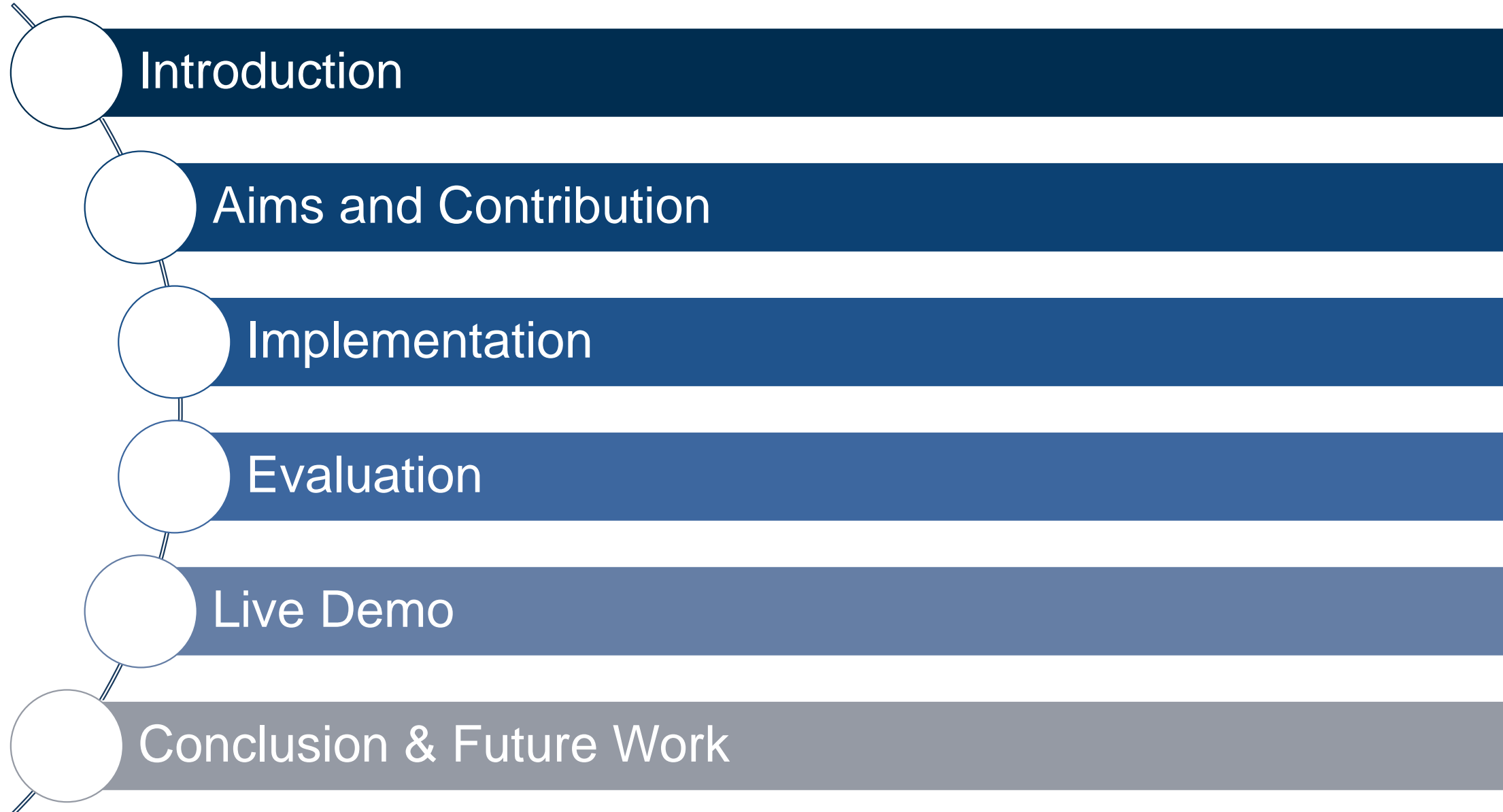


Computer-aided Analysis of Privacy Policies: Extraction of Data Subjects Rights According to GDPR

Sabrina Heinrich 19.08.2019

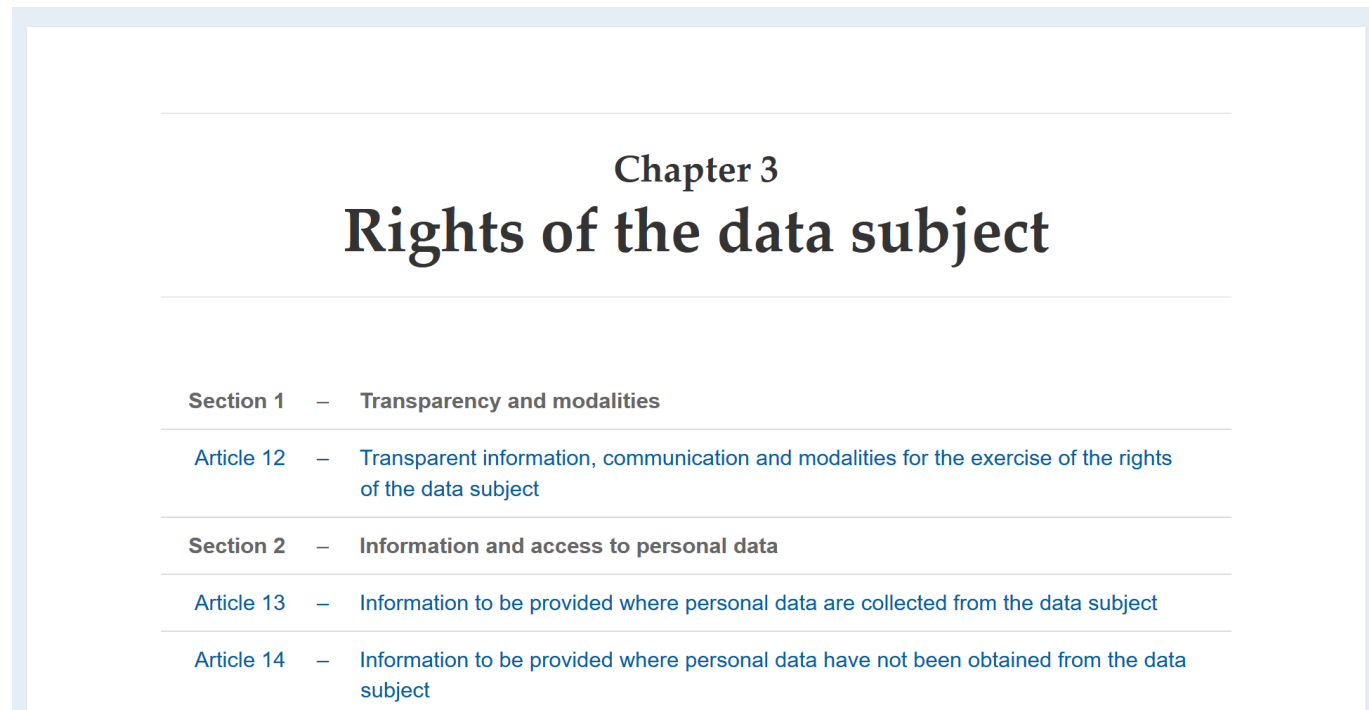
Chair of Software Engineering for Business Information Systems (sebis)
Faculty of Informatics
Technische Universität München
www.matthes.in.tum.de



Introduction

General Data Protection Regulation (GDPR)

- Implementation of GDPR in 2016 and applicable as of May 2018
 - In response to increasing digitalization
 - Regularization of processing of personal data within the EU
- One chapter deals with data subject rights



The image shows a table of contents for Chapter 3, titled 'Rights of the data subject'. The table is enclosed in a light blue border and contains the following entries:

| | |
|-----------------------------------|--|
| Chapter 3 | |
| Rights of the data subject | |
| Section 1 | – Transparency and modalities |
| Article 12 | – Transparent information, communication and modalities for the exercise of the rights of the data subject |
| Section 2 | – Information and access to personal data |
| Article 13 | – Information to be provided where personal data are collected from the data subject |
| Article 14 | – Information to be provided where personal data have not been obtained from the data subject |

Motivation – Automated Extraction of Data Subject Rights

Privacy policies are most important source of information for data collection and usage

Problem:

Studies show that:

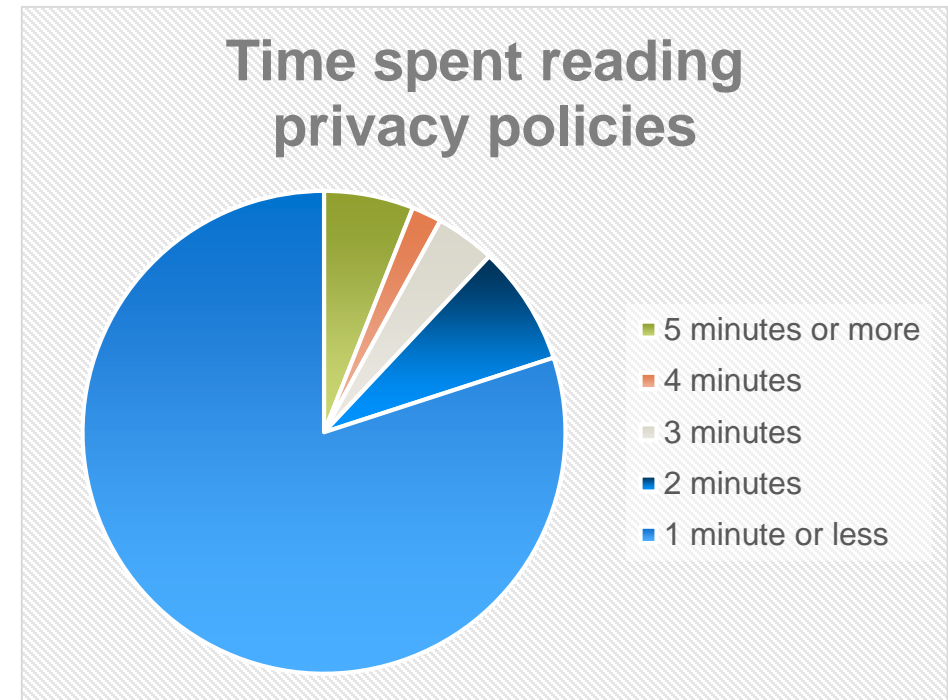
- Users are interested in protecting their data

But privacy policies...

- ... often lack readability
- ... cause an information overload

➡ Little time is spent reading privacy policies

➡ Not possible to make informed decisions regarding internet services



Related Work – Privacy Policy Analysis Tools

Clear Representation of Contents:

- Layered Privacy Notices
 - Platform for Privacy Preferences (P3P)
 - Privacy Nutrition Labels
- ➔ Rely on collaboration of service providers

| types of information | how we use your information | | | | | who we share your information with | |
|-------------------------|---------------------------------|------------------------|-----------|---------------|-----------|------------------------------------|---------------|
| | provide service & maintain site | research & development | marketing | telemarketing | profiling | other companies | public forums |
| contact information | ! | ! | OUT | OUT | — | IN | — |
| cookies | ! | ! | OUT | OUT | — | IN | — |
| demographic information | — | — | — | — | — | — | — |

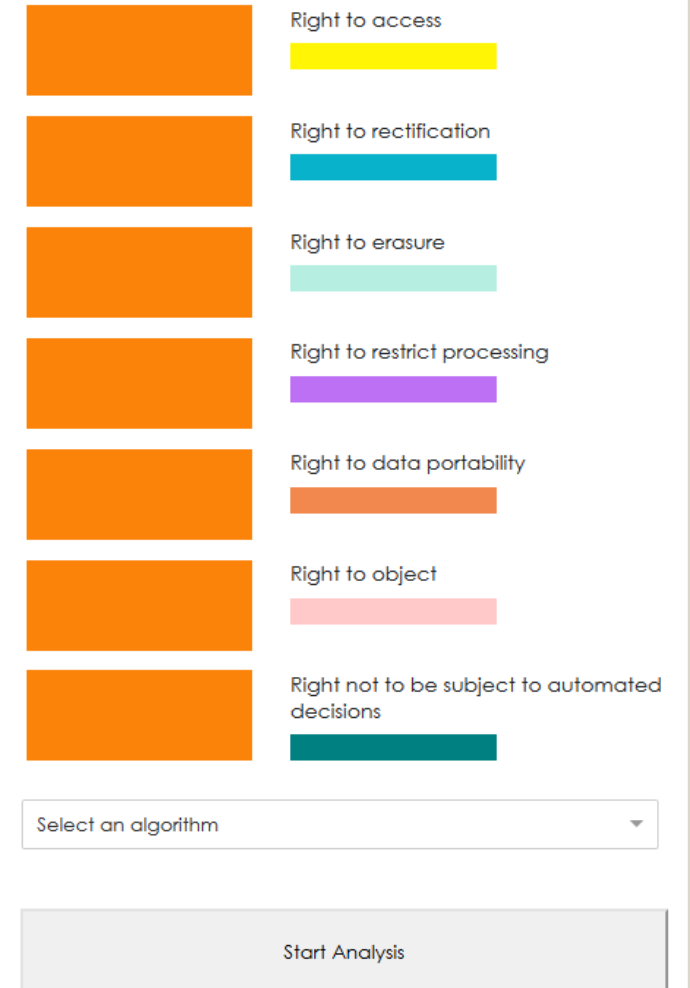
Automated Analysis of Content:












- Machine Learning Based Approaches
 - Supervised ML classifiers for sentence classification
 - Answer categorical questions or summarize content
- Rule Based Approaches
 - Extraction rules and pattern matching
 - Often in combination with ML

Privacy Policy Analysis - Data Subject Rights

Requirements:

- Input of privacy policy text
- Automated analysis of policy text
- Display of contained data subject rights
- Mark relevant sentences in the input text



| | | |
|--|--|---|
|  | Right to access |  |
|  | Right to rectification |  |
|  | Right to erasure |  |
|  | Right to restrict processing |  |
|  | Right to data portability |  |
|  | Right to object |  |
|  | Right not to be subject to automated decisions |  |

Select an algorithm 

Start Analysis

Aims and Contribution

Research Questions

RQ1: Which approaches exist to automatically analyze legal texts and privacy policies?

RQ2: Are supervised machine learning methods suitable to analyze privacy policies with regard to the coverage of the data subject rights?

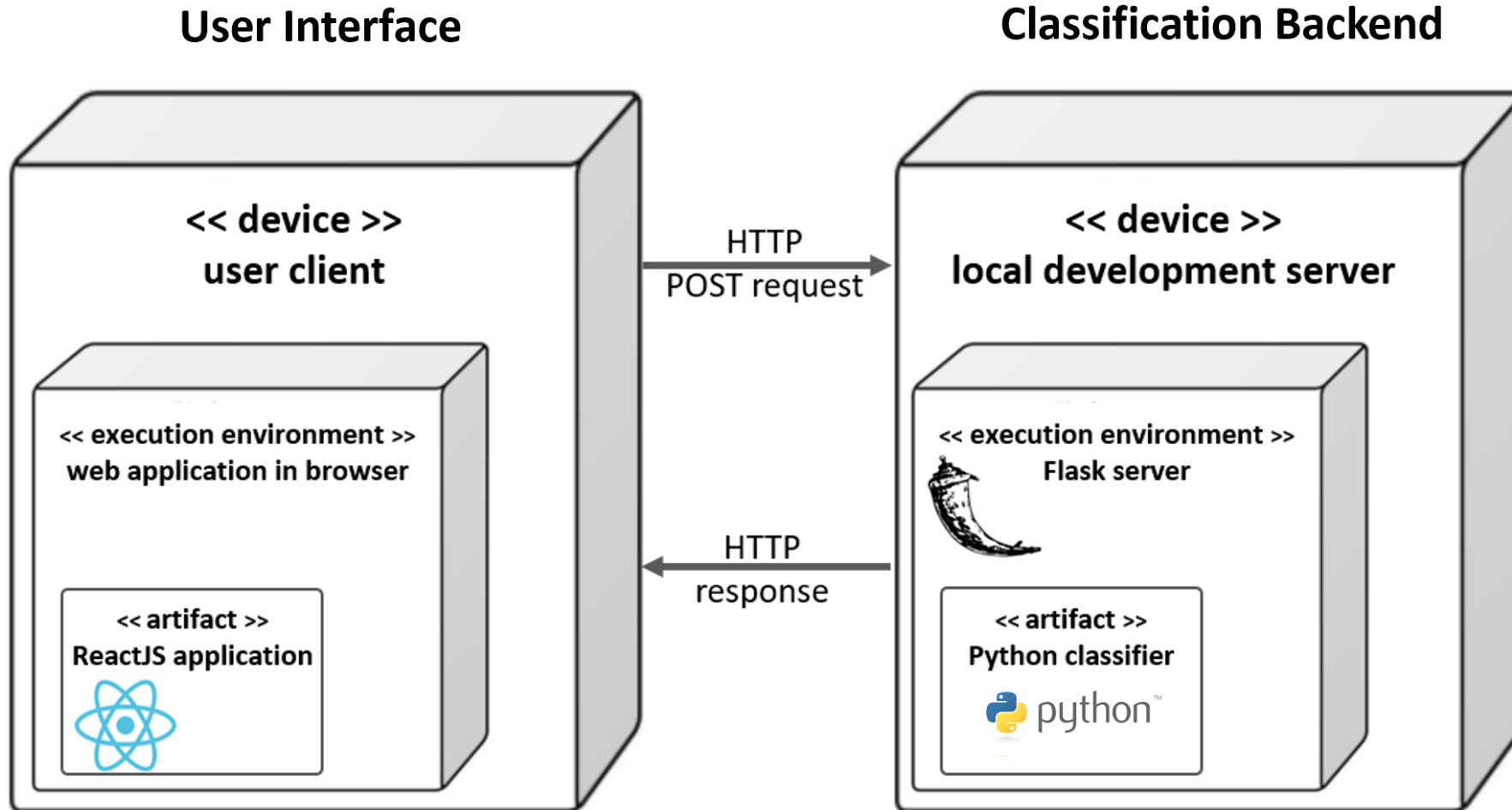
RQ3: What are the strengths and weaknesses of sentence classification and sequence labeling based approaches?

RQ4: Which supervised machine learning approach performs best at extracting data subject rights from privacy policies?

RQ5: Which performance can be achieved on the test data and can the automated extraction add value for private consumers regarding the understanding of privacy policies?

Implementation

Analysis Tool



Implementation

Dataset and Classes



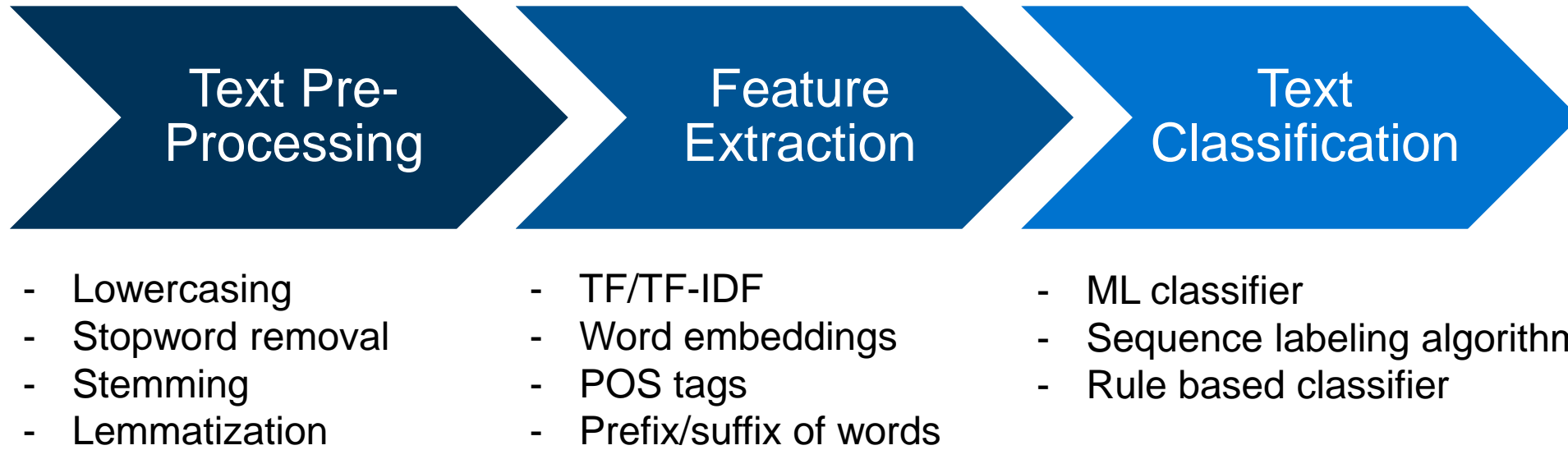
| Class | Sentence-Based | | Token-Based | |
|--|----------------|----------|-------------|----------|
| | Initial | Extended | Initial | Extended |
| Right to access | 126 | 267 | 148 | 329 |
| Right to rectification | 114 | 258 | 158 | 331 |
| Right to erasure | 120 | 259 | 158 | 319 |
| Right to restrict processing | 75 | 159 | 113 | 211 |
| Right to data portability | 66 | 150 | 98 | 249 |
| Right to object | 108 | 204 | 136 | 244 |
| Right not to be subject to automated decisions | 21 | 60 | 26 | 81 |
| TOTAL | 630 | 1357 | 837 | 1764 |

➡ Extended dataset in the course of the thesis

➡ Added null class for sentences not containing any right

Implementation

Automated Extraction of Data Subject Rights



Implementation

Automated Extraction of Data Subject Rights

ML Classifier



Sentence classification in this case is a multi-label classification problem:

*“You have the right to request that we confirm **what personally identifying information (PII) we collect or hold about you**, **provide a copy of your PII to you in a machine readable format** and to ask us to **correct or update the PII.**”*

Problem:

Traditional Classifiers can only deal with binary classification

➡ Algorithm Adoption or Problem Transformation

Implementation

Automated Extraction of Data Subject Rights

ML Classifier

Algorithm Adoption:

- Extensions of existing binary classification algorithms
- Multilabel-k-Nearest-Neighbor
- Multi-Label Twin Support Vector Machine

Problem Transformation:

- Convert multi-label classification problems into single-label classification problems
- In combination with binary classifier (e.g. decision tree, support vector machine)
- Binary Relevance
- Label Powerset
- Random K-Labelset
- Classifier Chain



Implementation

Automated Extraction of Data Subject Rights

ML Classifier

Binary Relevance:

| Sentence | Access | Object | Erasure |
|----------|--------|--------|---------|
| 1 | 1 | 0 | 1 |
| 2 | 1 | 1 | 0 |
| 3 | 0 | 1 | 1 |

| | Access | \neg Access | | Object | \neg Object | | Erasure | \neg Erasure |
|---|--------|---------------|---|--------|---------------|---|---------|----------------|
| 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 |
| 2 | 1 | 0 | 2 | 1 | 0 | 2 | 0 | 1 |
| 3 | 0 | 1 | 3 | 1 | 0 | 3 | 1 | 0 |



Implementation

Automated Extraction of Data Subject Rights

Sequence Labeling Algorithm

Conditional Random Field (CRF):

- Discriminative and graphical model for sequential data
- Used for sentence classification

```
[(['you', 'I'), ('may', 'I'), ('update', 'Rectification (Relevant)'), ('or', 'Rectification (Relevant)'), ('correct', 'Rectification (Relevant)'), ('information', 'I'), ('you', 'I'), ('have', 'I'), ('provided', 'I'), ('to', 'I'), ('us', 'I'), ('by', 'I'), ('going', 'I'), ('into', 'I'), ('the', 'I'), ('user', 'I'), ('account', 'I'), ('settings', 'I'), ('screen', 'I'), ('within', 'I'), ('the', 'I'), ('app', 'I'), ('', 'I')], ...]
```



Implementation

Automated Extraction of Data Subject Rights

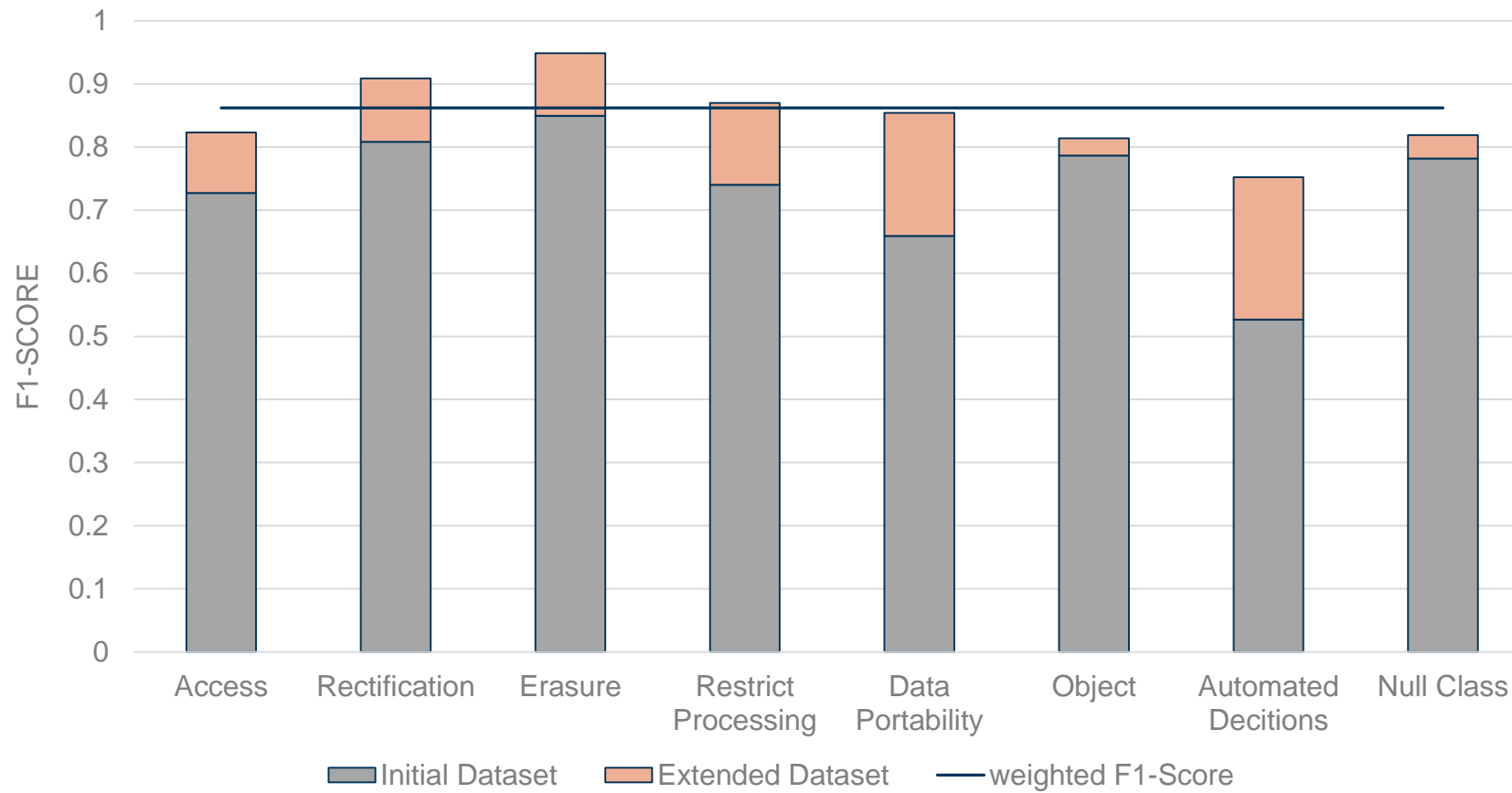
Rule Based Classifier

- Dictionary containing keywords for each class
- Automatically generated using training data and different thresholds

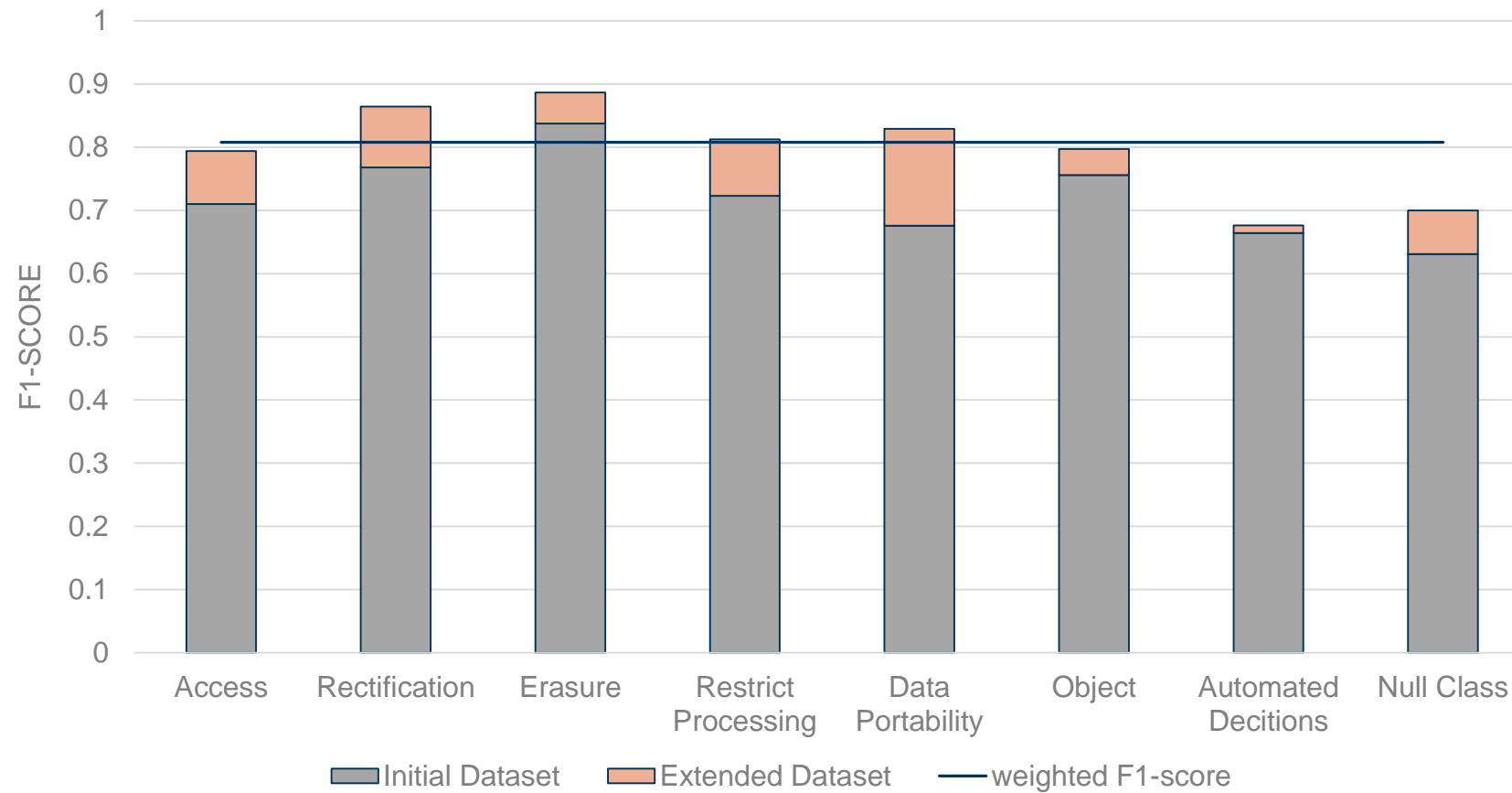
```
{Access:          {"access": 58,    "review": 8},
Rectification:    {"updat": 17,    "correct": 40,
                  "rectif": 18,    "edit": 5,
                  "incorrect": 7,  "rectifi": 7,
                  "incomplet": 6},
Erasure:         {"delet": 57,    "erasur": 16,
                  "account": 8,    "eras": 13,
                  "remov": 6},
Restrict Processing: {"restrict": 42, "limit": 5}
Data Portability: {"format": 14,   "portabl": 16,
                  "structur": 8,   "commonli": 5,
                  "machin": 10,    "readabl": 10, ...}
```



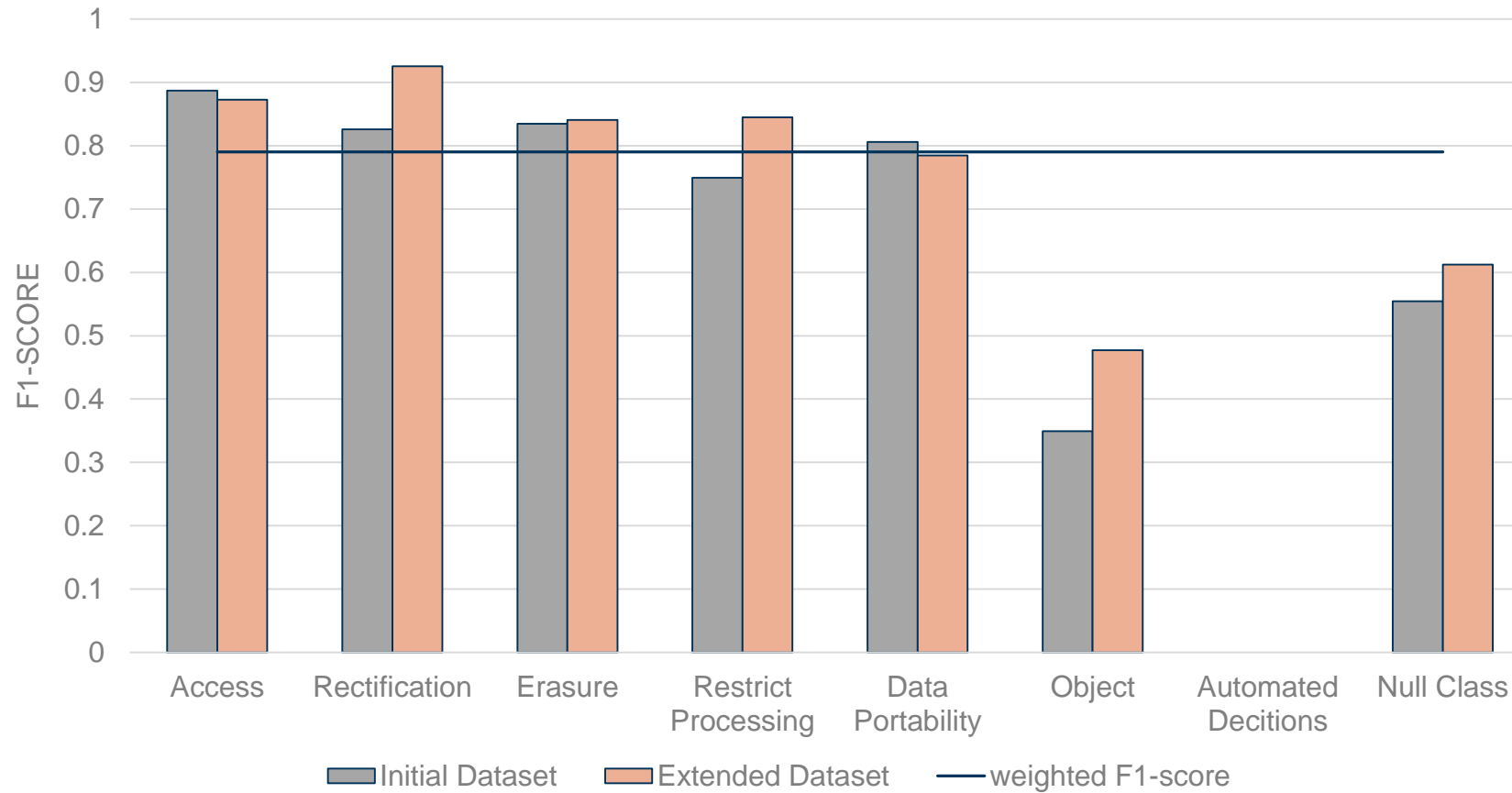
Quantitative Results – ML Classifier



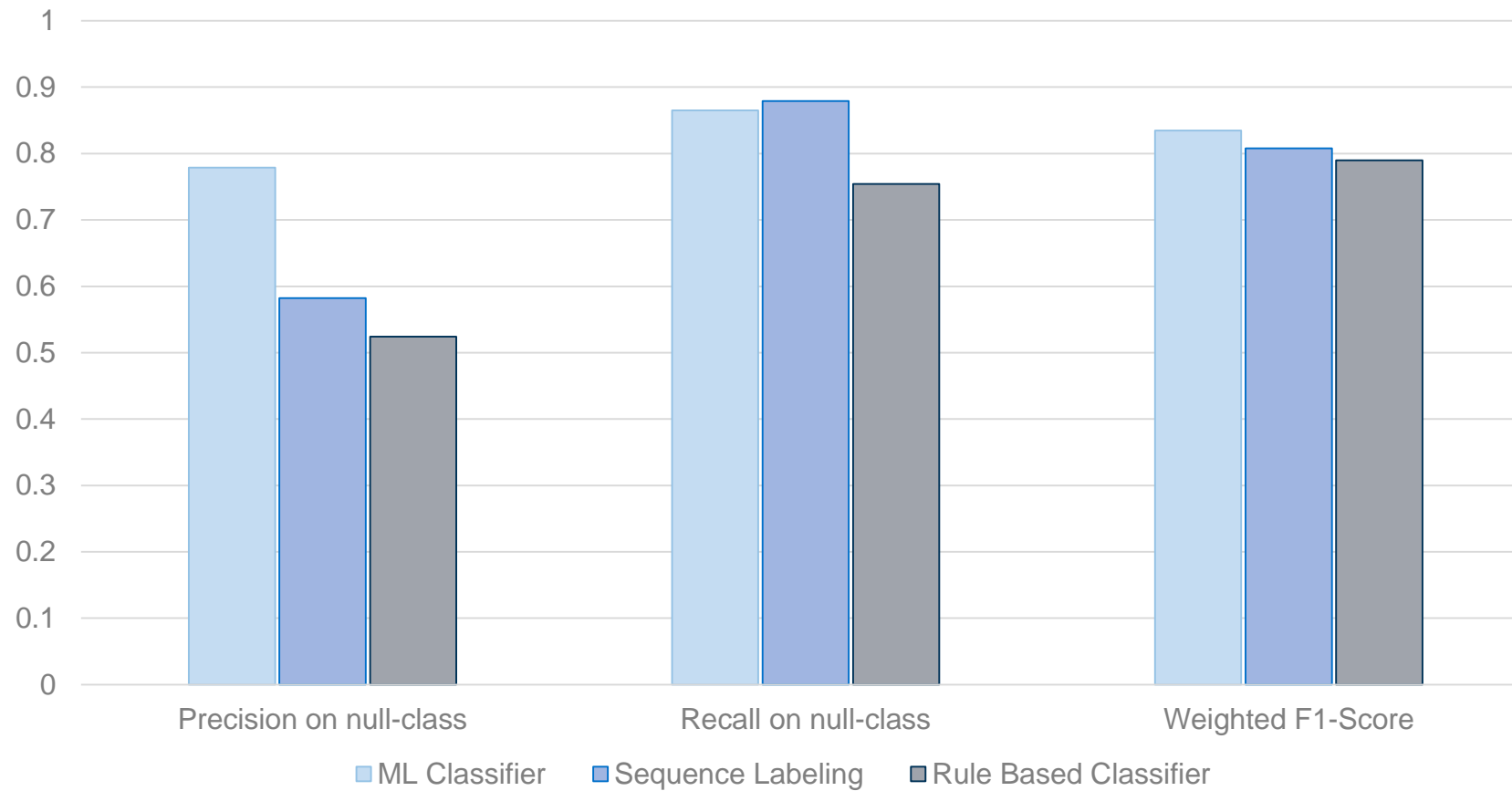
Quantitative Results – Sequence Labeling



Quantitative Results – Rule Based Classifier



Quantitative Results – Comparison of Three Approaches



Privacy Policy Analysis - Data Subject Rights



-  Right to access 
-  Right to rectification 
-  Right to erasure 
-  Right to restrict processing 
-  Right to data portability 
-  Right to object 
-  Right not to be subject to automated decisions 

Select an algorithm 

Start Analysis

Conclusion & Future Work

Conclusion

RQ1: Which approaches exist to automatically analyze legal texts and privacy policies?

RQ2: Are supervised machine learning methods suitable to analyze privacy policies with regard to the coverage of the data subject rights?

RQ3: What are the strengths and weaknesses of sentence classification and sequence labeling based approaches?

RQ4: Which supervised machine learning approach performs best at extracting data subject rights from privacy policies?

RQ5: Which performance can be achieved on the test data and can the automated extraction add value for private consumers regarding the understanding of privacy policies?

Conclusion & Future Work

Future Work

Points that remain unresolved by this thesis and suggest improvements:

- Interviews for qualitative evaluation of analysis tool
- Extensive parameter tuning
- Creating a larger and balanced dataset for training
- Widen the area of application
- Evaluate combinations of ML and rule based approaches

- [1] <https://dsgvo-gesetz.de/kapitel-3/>
- [2] Obar, Jonathan A., and Anne Oeldorf-Hirsch. "The biggest lie on the internet: Ignoring the privacy policies and terms of service policies of social networking services." *Information, Communication & Society* (2018): 1-20.
- [3] D. S. E. Kettner and P. D. C. Thorun. PGuardGemeinsamer Abschlussbericht. Bundesministerium für Bildung und Forschung, Institut für Angewandte Informatik e. V. (InfAi), mediaTest digital GmbH, Quadriga Hochschule Berlin and Selbstregulierung Informationswirtschaft e.V. (SRIW).
- [4] Reidenberg, Joel R., et al. "Disagreeable privacy policies: Mismatches between meaning and users' understanding." *Berkeley Tech. LJ* 30 (2015): 39.
- [5] R. W. Reeder. Expandable Grids: A user interface visualization technique and a policy semantics to support fast, accurate security and privacy policy authoring. Tech. rep. CARNEGIE-MELLON UNIV PITTSBURGH PA SCHOOL OF COMPUTER SCIENCE, 2008.
- [6] E. Costante, Y. Sun, M. Petkovi´c, and J. den Hartog. "A machine learning solution to assess privacy policy completeness:(short paper)." In: Proceedings of the 2012 ACM workshop on Privacy in the electronic society. ACM. 2012, pp. 91–96.
- [7] W. Ammar, S. Wilson, N. Sadeh, and N. A. Smith. "Automatic categorization of privacy policies: A pilot study." In: School of Computer Science, Language Technology Institute, Technical Report CMU-LTI-12-019 (2012).
- [8] G. Tsoumakas and I. Katakis. "Multi-label classification: An overview." In: International Journal of Data Warehousing and Mining (IJDWM) 3.3 (2007), pp. 1–13.
- [9] M.-L. Zhang and Z.-H. Zhou. "ML-KNN: A lazy learning approach to multi-label learning." In: Pattern recognition 40.7 (2007), pp. 2038–2048.
- [10] W.-J. Chen, Y.-H. Shao, C.-N. Li, and N.-Y. Deng. "MLTSVM: a novel twin support vector machine to multi-label learning." In: Pattern Recognition 52 (2016), pp. 61–74.
- [11] J. Read, B. Pfahringer, G. Holmes, and E. Frank. "Classifier chains for multi-label classification." In: Machine learning 85.3 (2011), p. 333.
- [12] G. Tsoumakas and I. Vlahavas. "Random k-labelsets: An ensemble method for multilabel classification." In: European conference on machine learning. Springer. 2007, pp. 406–417.
- [13] Manning, Christopher D., Christopher D. Manning, and Hinrich Schütze. *Foundations of statistical natural language processing*. MIT press



B. Sc.

Sabrina Heinrich

Technische Universität München
Faculty of Informatics
Chair of Software Engineering for Business
Information Systems

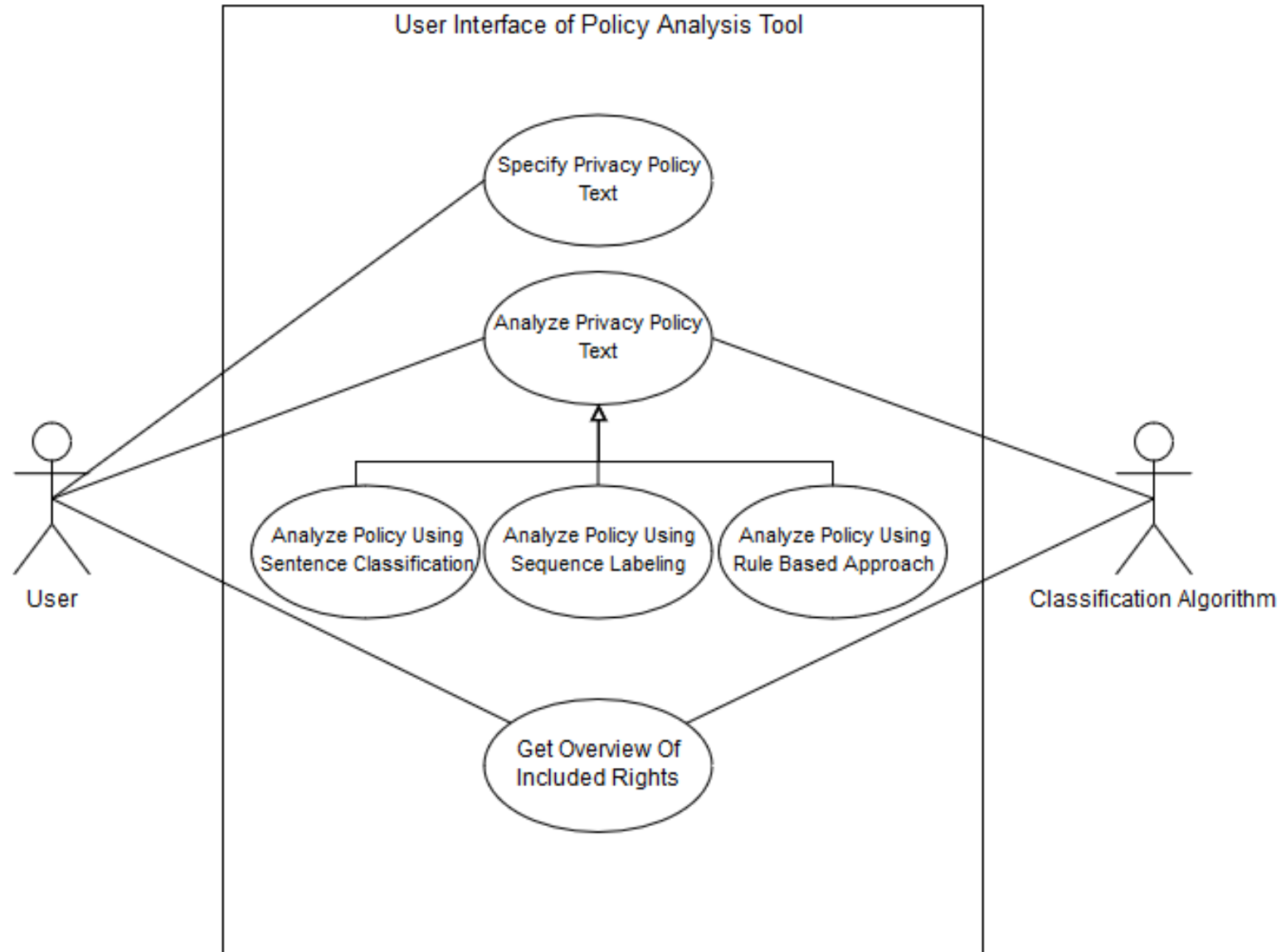
Boltzmannstraße 3
85748 Garching bei München

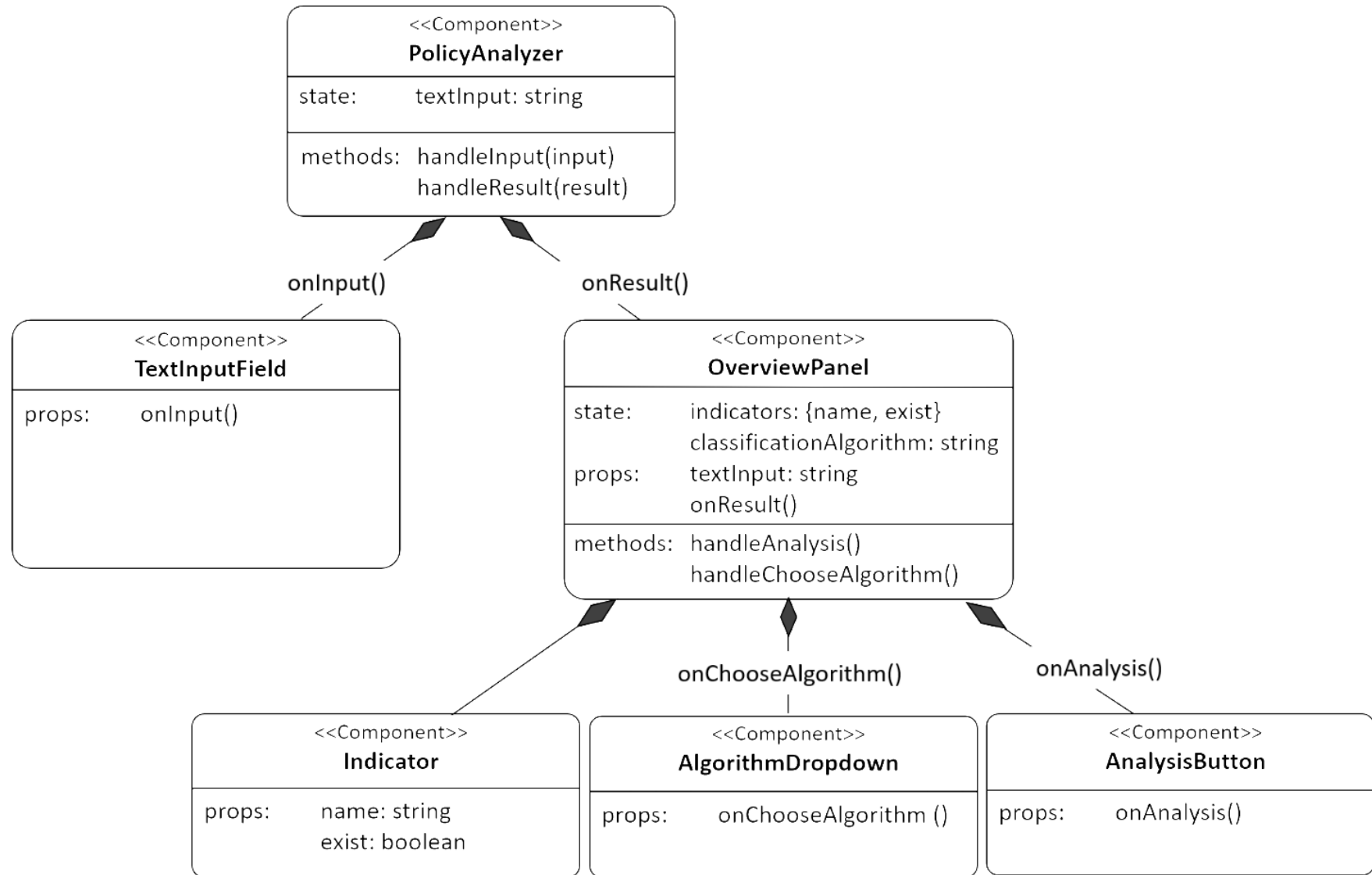
Tel +49.89.289.17132
Fax +49.89.289.17136

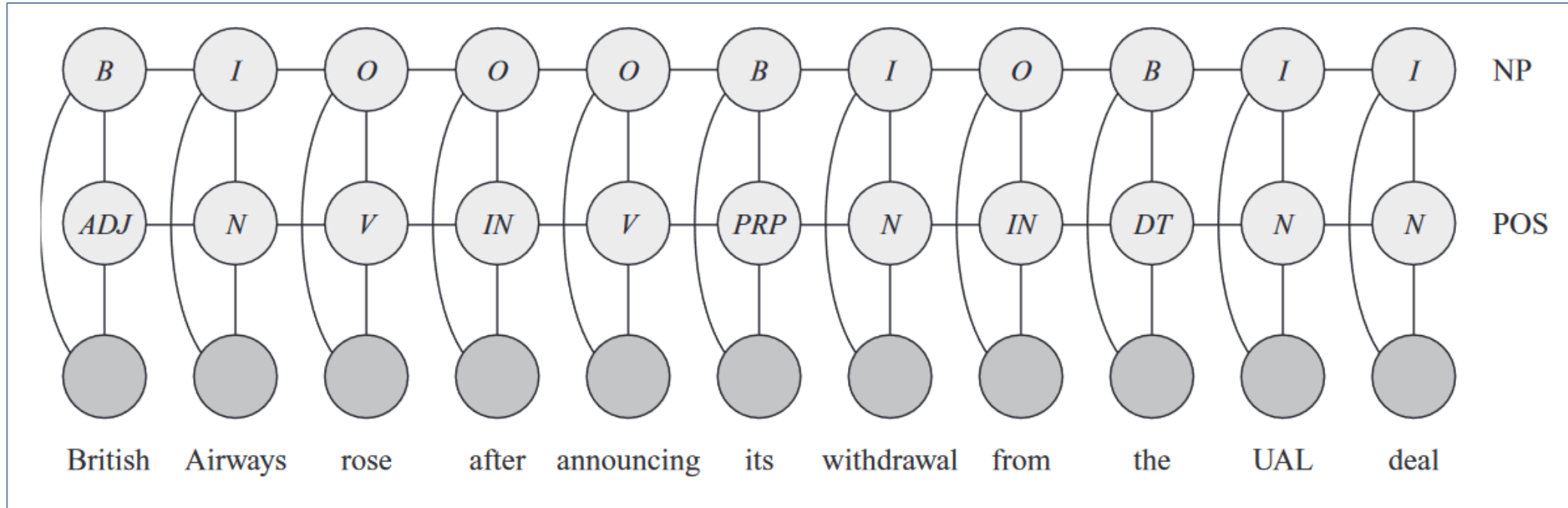
matthes@in.tum.de
www.matthes.in.tum.de



Use-Case Diagram



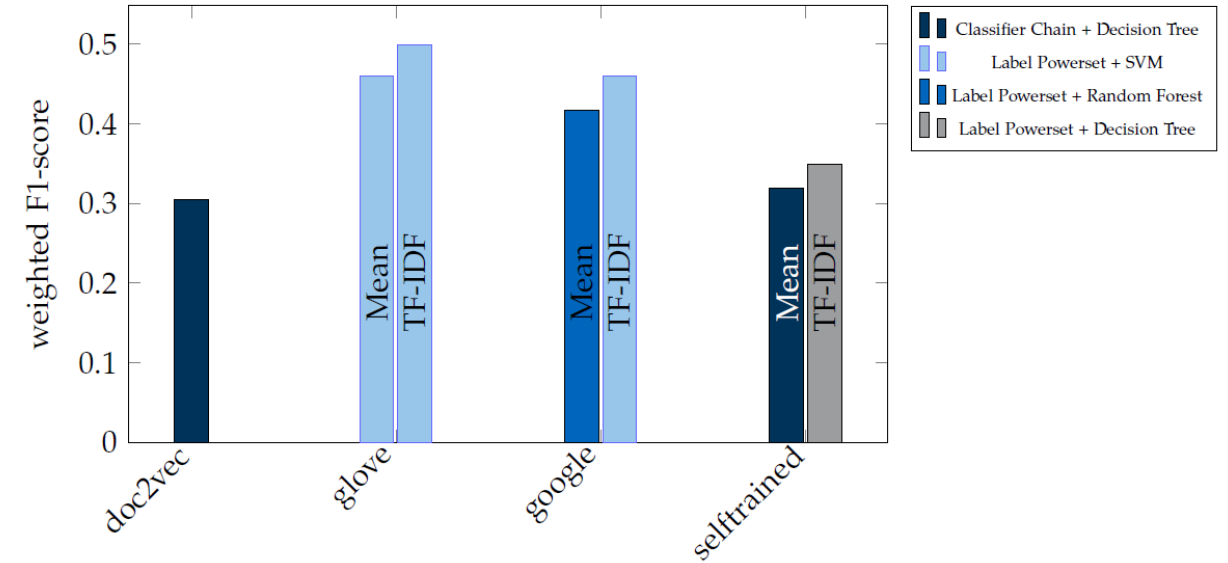
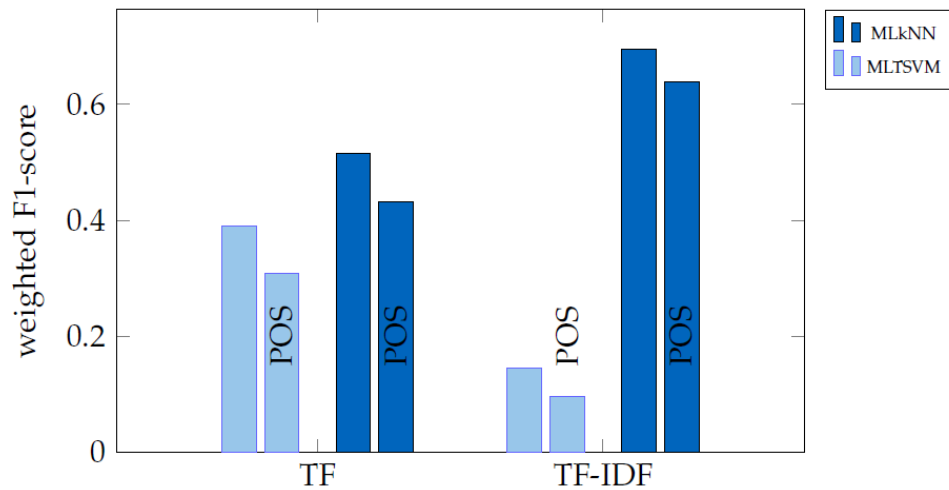
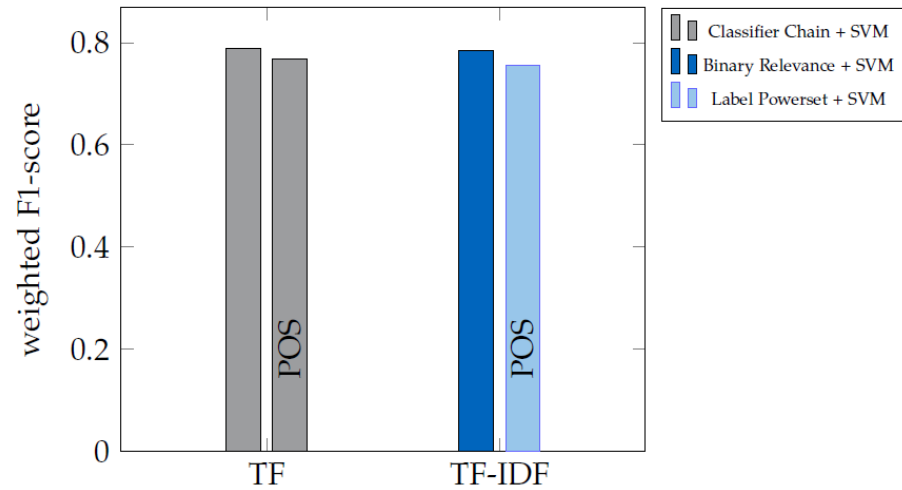





```
{Access: {"access": 58, "review": 8},
Rectification: {"updat": 17, "correct": 40,
"rectif": 18, "edit": 5,
"incorrect": 7, "rectifi": 7,
"incomplet": 6},
Erasure: {"delet": 57, "erasur": 16,
"account": 8, "eras": 13,
"remov": 6},
Restrict Processing: {"restrict": 42, "limit": 5}
Data Portability: {"format": 14, "portabl": 16,
"structur": 8, "commonli": 5,
"machin": 10, "readabl": 10,
"transmit": 4, "transfer": 4},
Object: {"object": 56, "withdraw": 14,
"consent": 9, "opt": 16}
Automated Decisions: {"decis": 9, "subject": 4,
"base": 6, "sole": 5,
"autom": 8}}
```

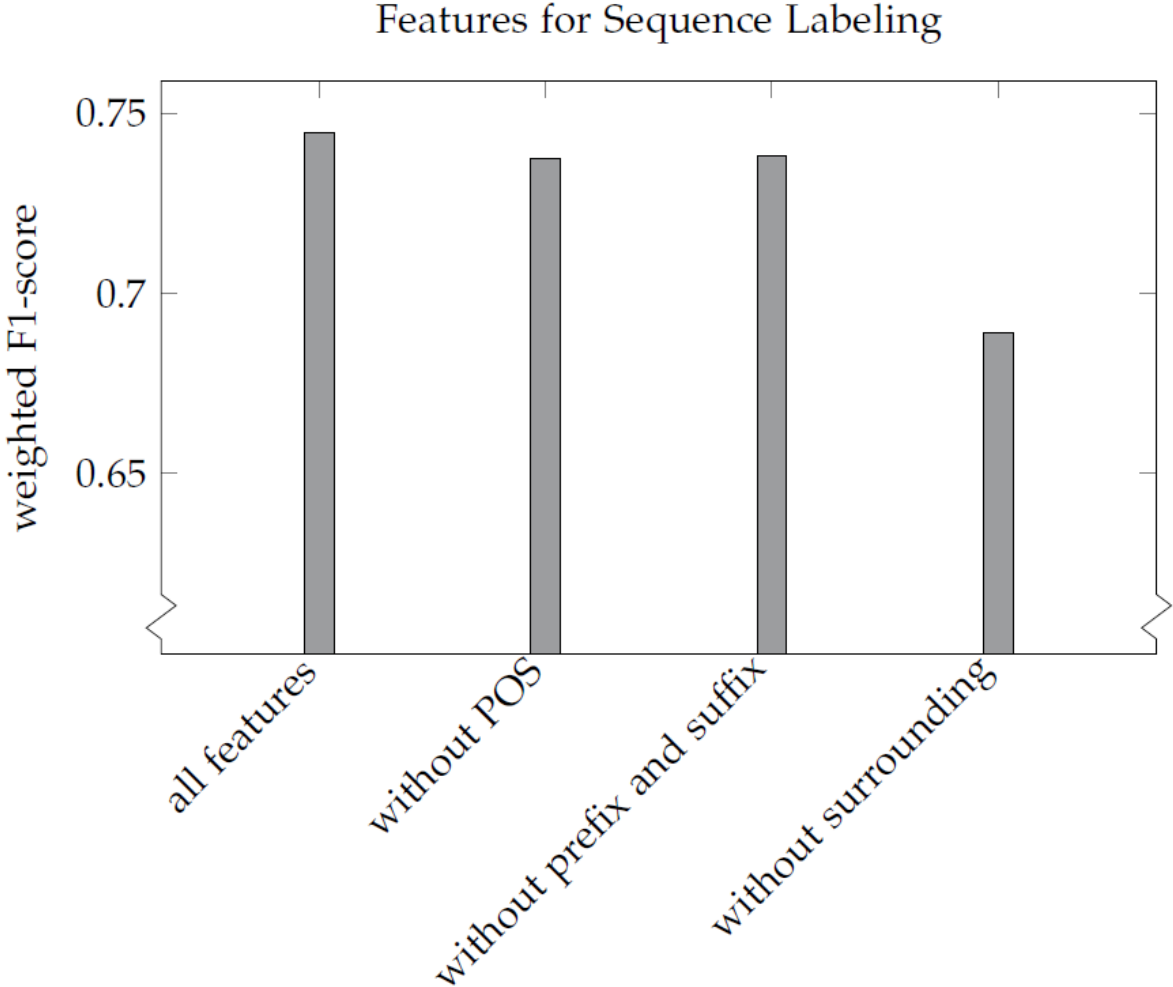
Sentence Classification

Embeddings and Sparse Representations

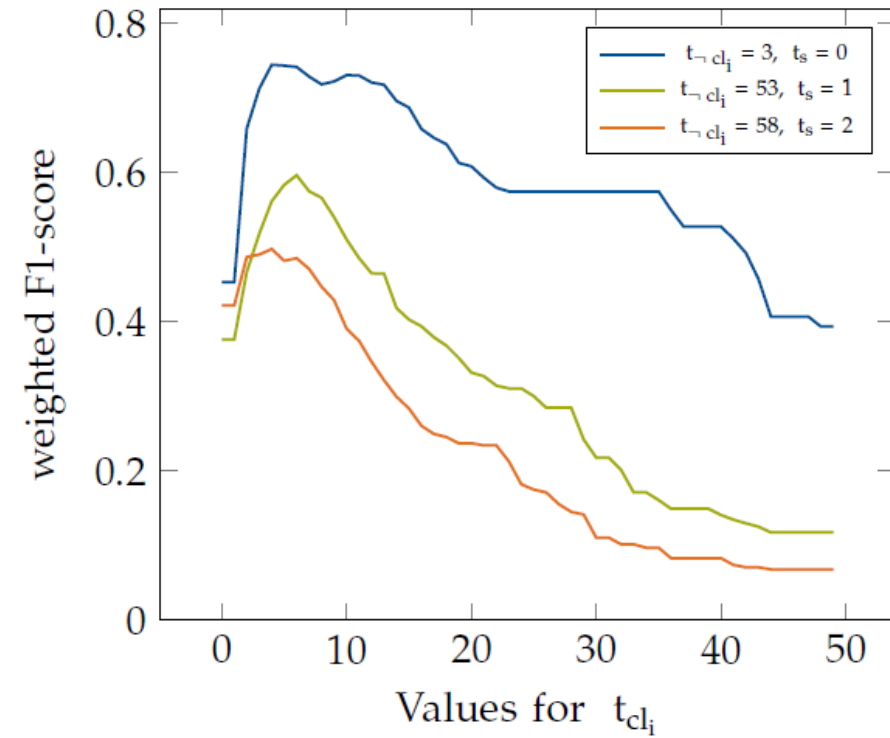
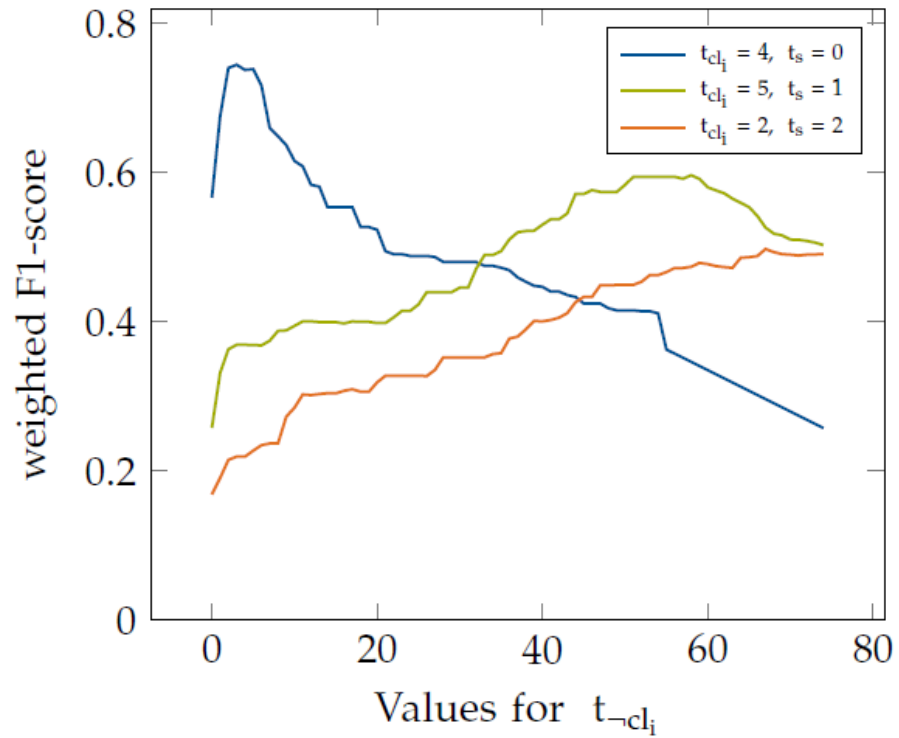


Sequence Labeling

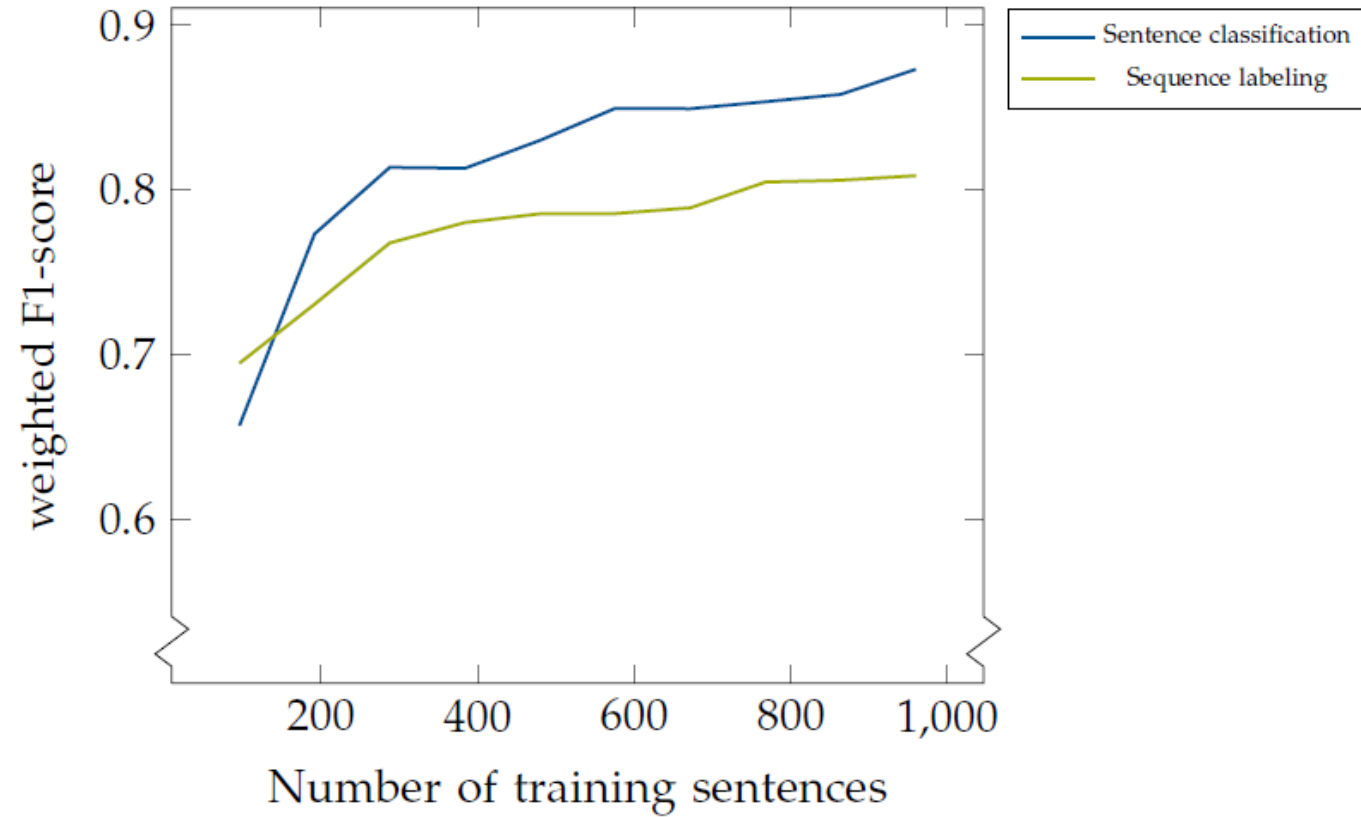
Impact of Different Features



Combinations of Threshold for Dictionary Generation



Learning Curve of ML Methods



Privacy Policy Analysis - Data Subject Rights

7\.. What rights do I have?

Access my data:

You may be able to receive a report about the data Blizzard stores about you. To learn more about accessing your data, please click on this [link](<https://eu.battle.net/support/help/product/services/1327/1329>) .

Erasure:

You may have a right to request the deletion of some your data. Blizzard is required by law to retain certain data. To learn more about having your data erased or anonymized, please click on this [link](<https://eu.battle.net/support/help/product/services/1327/1328>) .

Object:

You may object to the processing of your data if you believe Blizzard doesn't have the appropriate rights to engage in that processing or if you wish to ask Blizzard not to process your personal data for direct marketing purposes. To learn more about making an objection, please click on this [link](<https://eu.battle.net/support/help/product/services/1327/1330>) .

Rectification:

You can correct your personal data if you feel it's inaccurate. Please log in to your Blizzard account and use the account management settings to make the changes.

Restrict Processing:

You can request that your data no longer be processed in certain cases, for instance, when you exercise your right to object as described above. When objecting, you will be given an option to restrict processing.

Portability:

Data portability is the ability to move data from one company to another (for instance, when you transfer your mobile phone service to another carrier). While this right is not very applicable to Blizzard's current business, we will provide you with an electronic file with your most basic account information, should you request it. To obtain this file, please click this [link](<https://eu.battle.net/support/help/product/services/1327/1342/solution>) .

8\.. What is RealID?

- Right to access
- Right to rectification
- Right to erasure
- Right to restrict processing
- Right to data portability
- Right to object
- Right not to be subject to automated decisions

ML classifier

Start Analysis

Privacy Policy Analysis - Data Subject Rights

7\ What rights do I have?

Access my data:

You may be able to receive a report about the data Blizzard stores about you. To learn more about accessing your data, please click on this [link] (<https://eu.battle.net/support/help/product/services/1327/1329>) .

Erasure:

You may have a right to request the deletion of some your data. Blizzard is required by law to retain certain data. To learn more about having your data erased or anonymized, please click on this [link] (<https://eu.battle.net/support/help/product/services/1327/1328>) .

Object:

You may object to the processing of your data if you believe Blizzard doesn't have the appropriate rights to engage in that processing or if you wish to ask Blizzard not to process your personal data for direct marketing purposes. To learn more about making an objection, please click on this [link] (<https://eu.battle.net/support/help/product/services/1327/1330>) .

Rectification:

You can correct your personal data if you feel it's inaccurate. Please log in to your Blizzard account and use the account management settings to make the changes.














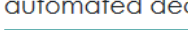
Restrict Processing:

You can request that your data no longer be processed in certain cases, for instance, when you exercise your right to object as described above. When objecting, you will be given an option to restrict processing.

Portability:

Data portability is the ability to move data from one company to another (for instance, when you transfer your mobile phone service to another carrier). While this right is not very applicable to Blizzard's current business, we will provide you with an electronic file with your most basic account information, should you request it. To obtain this file, please click this [link] (<https://eu.battle.net/support/help/product/services/1327/1342/solution>) .

8\ What is RealID?

-  Right to access 
-  Right to rectification 
-  Right to erasure 
-  Right to restrict processing 
-  Right to data portability 
-  Right to object 
-  Right not to be subject to automated decisions 

Sequence labeling with CRF

Start Analysis