

# CAPPA: A Collective Awareness Platform for Privacy Policy Annotations

*Giorgos Hompis*

Thesis submitted in partial fulfillment of the requirements for the  
*Masters' of Science degree in Computer Science and Engineering*

University of Crete  
School of Sciences and Engineering  
Computer Science Department  
Voutes University Campus, 700 13 Heraklion, Crete, Greece

Thesis Advisor: Prof. *Dimitris Plexousakis*



UNIVERSITY OF CRETE  
COMPUTER SCIENCE DEPARTMENT

**CAPPA: A Collective Awareness Platform for Privacy Policy  
Annotations**

Thesis submitted by  
**Giorgos Hompis**  
in partial fulfillment of the requirements for the  
Masters' of Science degree in Computer Science

THESIS APPROVAL

Author: \_\_\_\_\_  
Giorgos Hompis

Committee approvals: \_\_\_\_\_  
Dimitris Plexousakis  
Professor, Thesis Supervisor

\_\_\_\_\_  
Evangelos Markatos  
Professor, Committee Member

\_\_\_\_\_  
Giorgos Flouris  
Researcher, Committee Member

Departmental approval: \_\_\_\_\_  
Antonios Argyros  
Professor, Director of Graduate Studies

Heraklion, March 2018



# CAPPA: A Collective Awareness Platform for Privacy Policy Annotations

## Abstract

The huge expansion of digital products and the corresponding user generated and gathered data have raised the importance of users privacy and privacy concerns. Currently, businesses and organizations around the world are enforced by law (e.g. the EU General Data Protection Regulation - GDPR) to provide information about how customers' data is being treated, usually in the form of privacy policy documents. Despite the fact that such regulations try primarily to give control back to citizens over their personal data, it's a common case that users are not engaged in this process. Since such documents are usually long and hard to read, users are not willing to spend a lot of time to read and understand them.

A current direction for addressing this problem is towards enriching privacy policy documents with annotations either through expert users or machine learning algorithms. In this thesis, we designed and implemented an on-line crowd-sourced platform that allows users to explore, annotate and review privacy policies of any kind of digital product (e.g. mobile applications, websites, appliances, etc.) in a friendly way. The platform is part of the tools designed for the CAPrice community, a collective awareness platform for privacy concerns and expectations.

Privacy policies are being annotated using a predefined set of tags, designed to address user concerns about what data are being collected and processed, by whom, for how long they are retained, how they are secured, and other privacy concerns. Users can contribute by adding entities like digital products, privacy policies, and annotations to documents or by reviewing entities added by other users. The platform helps and engages users towards this quest through various engagement tools (e.g. user scores) and document analysis tools (e.g. readability scores of privacy policies). An annotation can be considered as a valid or invalid one, based on the votes of the users and their aggregated score obtained using the Wilson score interval. The aim is to provide a collaborative crowdsourcing platform that will be considered the reference for user annotated privacy policy documents, for users, developers, researchers and policy makers. Towards this direction we have designed a ReST API that provides access to the database of digital products and their annotated privacy policies. As a result, this information can be exploited for the development of third party tools and algorithms.

We conducted a user-based evaluation of our platform, where users were split in two groups. Each group was asked to annotate a specific set of privacy policies obtained from the OPP-115 dataset, which is an expert-based annotated collection of privacy policies. Then each group had to review/vote the annotations of the other group and fill in the corresponding questionnaire. The analysis of the results shows the user friendliness of our platform and that the gathered crowd-sourced privacy policy annotations are of high importance and quality, comparable to annotations created by expert users.



# CAPPA: Μία Πλατφόρμα Συλλογικής Επίγνωσης για Επισημειώσεις Πολιτικών Απορρήτου

## Περίληψη

Η μαζική εξάπλωση των ψηφιακών προϊόντων και των αντίστοιχων παραγόμενων και συλλεγμένων δεδομένων χρηστών έχουν αυξήσει τη σημασία της προστασίας της ιδιωτικότητας. Επί του παρόντος, οι επιχειρήσεις και οι οργανισμοί σε διάφορα μέρη του κόσμου υποχρεούνται από το νόμο να παρέχουν πληροφορίες σχετικά με τον τρόπο επεξεργασίας των δεδομένων των πελατών τους, συνήθως με τη μορφή εγγράφων πολιτικής απορρήτου (π.χ. ο κανονισμός γενικής προστασίας δεδομένων της ΕΕ - GDPR). Παρά το γεγονός ότι οι κανονισμοί αυτοί προσπαθούν να δώσουν τον έλεγχο των προσωπικών δεδομένων πίσω στους πολίτες, συνήθως οι χρήστες δεν εμπλέκονται σε αυτή τη διαδικασία. Δεδομένου ότι τα έγγραφα αυτά είναι συνήθως μακριά και δύσκολο να διαβάσουν, οι χρήστες δεν είναι διατεθειμένοι να αφιερώσουν πολύ χρόνο για να τα διαβάσουν και να τα κατανοήσουν.

Μια τρέχουσα κατεύθυνση επίλυσης του προβλήματος αυτού είναι ο εμπλουτισμός των εγγράφων πολιτικής απορρήτου με επισημειώσεις είτε μέσω εμπειρογνομόνων είτε μέσω αλγορίθμων μηχανικής μάθησης. Σε αυτή την εργασία, σχεδιάσαμε και υλοποιήσαμε μια διαδικτυακή συλλογική πλατφόρμα που επιτρέπει στους χρήστες να εξερευνούν, να επισημαίνουν και να εξετάζουν τις πολιτικές απορρήτου οποιουδήποτε ψηφιακού προϊόντος (π.χ. κινητές εφαρμογές, ιστότοποι, έξυπνες συσκευές κ.λπ.) με φιλικό τρόπο. Αυτή η πλατφόρμα αποτελεί μέρος των εργαλείων που έχουν σχεδιαστεί για την κοινότητα CAPrice, μια συλλογική πλατφόρμα ευαισθητοποίησης για την προστασία της ιδιωτικότητας.

Οι πολιτικές απορρήτου επισημαίνονται χρησιμοποιώντας ένα προκαθορισμένο σύνολο ετικετών, σχεδιασμένο για να αντικατοπτρίζουν τις ανησυχίες των χρηστών σχετικά με το τι προσωπικά δεδομένα συλλέγονται και επεξεργάζονται, από ποιους, για πόσο καιρό διατηρούνται, πώς προστατεύονται και άλλες ανησυχίες γύρω από την ιδιωτικότητα. Οι χρήστες μπορούν να συνεισφέρουν προσθέτοντας οντότητες όπως ψηφιακά προϊόντα, έγγραφα πολιτικής απορρήτου και επισημάνσεις στα έγγραφα ή αξιολογώντας οντότητες που προστέθηκαν από άλλους χρήστες. Η πλατφόρμα βοηθά και εμπλέκει τους χρήστες προς αυτή την αναζήτηση μέσω διαφόρων εργαλείων εμπλοκής (π.χ. βαθμολογίες χρηστών) και εργαλεία ανάλυσης εγγράφων (π.χ. βαθμός αναγνωσιμότητας εγγράφων πολιτικής απορρήτου). Μια επισημείωση μπορεί να θεωρηθεί έγκυρη ή λανθασμένη βάσει των ψήφων των χρηστών και του συγκεντρωτικού τους σκορ που αποκτήθηκε χρησιμοποιώντας το Wilson score interval. Ο στόχος είναι να δημιουργηθεί μια συνεργατική πλατφόρμα crowdsourcing που θα αποτελέσει σημείο αναφοράς για επισημειωμένα έγγραφα πολιτικής απορρήτου από χρήστες, για χρήστες, προγραμματιστές, ερευνητές και δημιουργούς πολιτικών απορρήτου. Προς αυτή την κατεύθυνση έχουμε σχεδιάσει ένα ReST API που παρέχει πρόσβαση στη βάση δεδομένων των ψηφιακών προϊόντων και τις επισημειωμένες πολιτικές απορρήτου τους, επιτρέποντας την εκμετάλλευση αυτών των πληροφοριών για την ανάπτυξη τρίτων εργαλείων και αλγορίθμων.

Πραγματοποιήσαμε μια βασισμένη σε χρήστες αξιολόγηση της πλατφόρμας μας, όπου είχαμε δύο ομάδες χρηστών. Σε κάθε ομάδα ζητήθηκε η επισημείωση ενός συγκεκριμένου πλήθους πολιτικών απορρήτου από το σύνολο δεδομένων OPP-115, οι οποίες είναι επισημειωμένες από έμπειρους χρήστες. Στη συνέχεια από κάθε ομάδα ζητήθηκε η αξιολόγηση των επισημειώσεων της άλλης ομάδας χρηστών και η συμπλήρωση ενός ερωτηματολογίου. Η ανάλυση των αποτελεσμάτων αναδεικνύει τη φιλικότητα της πλατφόρμας μας προς το χρήστη και ότι οι επισημειώσεις σε πολιτικές απορρήτου που συλλέχθηκαν είναι υψηλής σημασίας και ποιότητας, συγκρίσιμες με εκείνες που δημιουργήθηκαν από έμπειρους χρήστες.



## Ευχαριστίες

Αρχικά, θα ήθελα να ευχαριστήσω τον καθηγητή του τμήματος Επιστήμης Υπολογιστών του Πανεπιστημίου Κρήτης και επιβλέποντα καθηγητή μου κ. Δημήτρη Πλεξουσάκη, για την υποστήριξη του και το χρόνο που μου αφιέρωσε αλλά και γιατί μου έδειξε εμπιστοσύνη και δέχτηκε να γίνει ο επόπτης μου. Θα ήθελα επίσης να ευχαριστήσω τον καθηγητή του τμήματος Επιστήμης Υπολογιστών του Πανεπιστημίου Κρήτης κ. Ευάγγελο Μαρκάτο και τον ερευνητή του IT-ITE κ. Γεώργιο Φλουρή για την προθυμία τους να συμμετάσχουν στην τριμελή επιτροπή.

Στη συνέχεια θα ήθελα να ευχαριστήσω τα μέλη του εργαστηρίου του των Πληροφοριακών Συστημάτων του Ινστιτούτου Πληροφορικής που συμμετείχαν στην αξιολόγηση της πλατφόρμας και ιδιαίτερα τον κ. Θεόδωρο Πάτκο για τις πολύτιμες συμβουλές του και τον κ. Παναγιώτη Παπαδόχο για την βοήθεια και υποστήριξη του αλλά και για την διαθεσιμότητα του καθόλη τη διάρκεια της μεταπτυχιακής μου εργασίας. Επιπλέον θα ήθελα να ευχαριστήσω θερμά το τμήμα Επιστήμης Υπολογιστών του Πανεπιστημίου Κρήτης για την υψηλής ποιότητας ακαδημαϊκή μόρφωση και υπηρεσίες που μου προσέφερε κατά τη διάρκεια των σπουδών μου.

Τέλος θα ήθελα να ευχαριστήσω τους φίλους για την στήριξή τους στις δύσκολες ώρες και την συμπαράστασή τους όλο αυτό το διάστημα και, περισσότερο από όλους θα ήθελα να ευχαριστήσω την οικογένειά μου και κυρίως την μητέρα μου Καλλιόπη που μου έδωσε την μεγαλύτερη ώθηση και συνέβαλε τα μέγιστα για να φτάσω στο σημείο που βρίσκομαι σήμερα.



στους γονείς μου



# Contents

<b>Table of Contents</b>	<b>i</b>
<b>List of Tables</b>	<b>iii</b>
<b>List of Figures</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Background &amp; Related Work</b>	<b>5</b>
2.1 Privacy Awareness . . . . .	5
2.2 Standard Privacy Policy Formats . . . . .	5
2.3 Privacy Policy Analysis . . . . .	6
2.3.1 Readability Metrics . . . . .	6
2.3.2 Extracting Privacy Policy Features . . . . .	6
2.4 Related Tools . . . . .	8
2.5 Crowdsourcing Platforms . . . . .	8
2.6 CAPrice . . . . .	9
<b>3 System Requirements &amp; Design</b>	<b>11</b>
3.1 Requirements . . . . .	11
3.2 Definitions . . . . .	11
3.2.1 Users & User Roles . . . . .	11
3.2.2 Controversial Entities . . . . .	13
3.2.2.1 Digital Products . . . . .	13
3.2.2.2 Legal Documents . . . . .	14
3.2.2.3 Annotations . . . . .	14
3.2.2.4 Entity Voting Scheme . . . . .	14
3.2.3 User Ratings & Badges . . . . .	17
3.2.4 Annotation Schema . . . . .	20
3.2.4.1 Tag Attributes & Values . . . . .	21
<b>4 System Architecture &amp; Implementation</b>	<b>23</b>
4.1 System Architecture . . . . .	23
4.1.1 Front-End . . . . .	23

4.1.2	Back-End . . . . .	23
4.2	REST API . . . . .	24
4.3	Modularity of the Annotation Schema . . . . .	26
4.4	Implementation . . . . .	26
4.4.1	Server Side vs Client Side Rendering . . . . .	26
4.4.2	Technologies Used . . . . .	27
4.4.2.1	Back-End . . . . .	27
4.4.2.2	Front-End . . . . .	28
4.4.2.3	Third Party Tools & Libraries . . . . .	28
<b>5</b>	<b>Use Cases</b>	<b>31</b>
5.1	Welcome Page . . . . .	31
5.2	Products List Page - Search Results Page . . . . .	32
5.3	Product Details Page . . . . .	33
5.4	Document View Page . . . . .	34
5.4.1	General Functionality . . . . .	34
5.4.2	Annotation Details . . . . .	36
<b>6</b>	<b>Evaluation</b>	<b>39</b>
6.1	Evaluation Setup . . . . .	39
6.2	Annotation Quality . . . . .	41
6.2.1	Process Statistics . . . . .	41
6.2.2	Relevance Metrics . . . . .	41
6.2.3	User vs Crowd Created Annotations . . . . .	43
6.3	Annotation Scores . . . . .	50
6.4	User Scores . . . . .	51
6.5	User Friendliness & General Remarks . . . . .	52
6.5.1	The Platform . . . . .	53
6.5.2	User Experience & Feelings . . . . .	54
<b>7</b>	<b>Conclusion &amp; Future Plans</b>	<b>57</b>
7.1	Performance Evaluation & Scalability . . . . .	57
7.2	UI Improvements & Extensions . . . . .	57
7.3	User Engagement Improvements . . . . .	58
7.4	More Components & Functionality . . . . .	58
7.4.1	Like/Dislike Feature . . . . .	59
7.4.2	Product External Link Resources . . . . .	59
7.4.3	Mobile Applications - Android Permissions . . . . .	59
7.4.4	Argument Web . . . . .	59
7.5	Other Tools . . . . .	60
7.5.1	Product Privacy Evaluation . . . . .	60
7.5.2	Text Difference Between Documents . . . . .	61
7.5.3	Guided/Assisted Annotation Creation . . . . .	61
7.5.3.1	Deployment of NLP & Machine Learning Algorithms	61

7.5.3.2	Ambiguous Annotation Flags . . . . .	61
7.6	Discussion . . . . .	62
<b>Appendices</b>		<b>63</b>
<b>A</b>	<b>Tags: Concerns, Attributes &amp; Values</b>	<b>65</b>
A.1	Privacy Concerns . . . . .	66
A.2	Tag Attributes . . . . .	67
A.3	Tag Attribute Values . . . . .	71
<b>B</b>	<b>Evaluation Questionnaire</b>	<b>101</b>
<b>C</b>	<b>ReST API</b>	<b>109</b>
C.1	Data API . . . . .	109
C.2	View API . . . . .	119
<b>Bibliography</b>		<b>123</b>





# List of Tables

3.1	Permission levels for each user role and actions for creating ( $\mathcal{C}$ ), deleting ( $\mathcal{D}$ ) and voting ( $\mathcal{V}$ ) controversial entities . . . . .	13
3.2	Available user role upgrades and the corresponding upgrade score thresholds . . . . .	13
3.3	Mapping between confidence level and the resulting Z value . . . .	16
3.4	The Reward Badge categories . . . . .	19
3.5	Predefined badges with the score value gain for the user score . .	22
6.1	Confusion matrix for valid & invalid annotations . . . . .	48
A.1	Privacy Concerns . . . . .	66
A.2	Attributes for <i>First Party Collection/Use</i> Privacy Concern . . . .	67
A.3	Attributes for <i>Third Party Sharing/Collection</i> Privacy Concern . .	68
A.4	Attributes for <i>User Choice/Control</i> Privacy Concern . . . . .	69
A.5	Attributes for <i>User Access, Edit and Deletion</i> Privacy Concern . .	69
A.6	Attributes for <i>Data Retention</i> Privacy Concern . . . . .	70
A.7	Attributes for <i>Data Security</i> Privacy Concern . . . . .	70
A.8	Attributes for <i>Policy Change</i> Privacy Concern . . . . .	70
A.9	Attributes for <i>Do Not Track</i> Privacy Concern . . . . .	70
A.10	Attributes for <i>International And Specific Audiences</i> Privacy Concern	71
A.11	Attributes for <i>Other</i> Privacy Concern . . . . .	71
A.12	Tag Attribute Values for <i>First Party Collection/Use</i> Privacy Concern . . . . .	78
A.13	Tag Attribute Values for <i>Third Party Sharing/Collection</i> Privacy Concern . . . . .	85
A.14	Tag Attribute Values for <i>User Choice/Control</i> Privacy Concern . .	90
A.15	Tag Attribute Values for <i>User Access, Edit and Deletion</i> Privacy Concern . . . . .	92
A.16	Tag Attribute Values for <i>Data Retention</i> Privacy Concern . . . .	95
A.17	Tag Attribute Values for <i>Data Security</i> Privacy Concern . . . . .	96
A.18	Tag Attribute Values for <i>Policy Change</i> Privacy Concern . . . . .	98
A.19	Tag Attribute Values for <i>Do Not Track</i> Privacy Concern . . . . .	98
A.20	Tag Attribute Values for <i>International and Specific Audiences</i> Privacy Concern . . . . .	99

A.21 Tag Attribute Values for <i>Other</i> Privacy Concern . . . . .	99
C.1 Data - Get Annotation Details . . . . .	109
C.2 Data - Delete Annotation . . . . .	110
C.3 Data - Create Annotation . . . . .	111
C.4 Data - Get Business Category Labels . . . . .	111
C.5 Data - Submit New Vote . . . . .	112
C.6 Data - Delete Controversial Entity . . . . .	112
C.7 Data - Submit New Document . . . . .	113
C.8 Data - Get Document Details . . . . .	113
C.9 Data - Delete Document . . . . .	114
C.10 Data - Get Document Annotations . . . . .	114
C.11 Data - Add New Product . . . . .	114
C.12 Data - Get Product Details . . . . .	115
C.13 Data - Get All Products . . . . .	115
C.14 Data - Delete Product . . . . .	115
C.15 Data - Get All Tags . . . . .	115
C.16 Data - Get Privacy Concern Tags . . . . .	116
C.17 Data - Get User Details . . . . .	116
C.18 Data - Create New User . . . . .	117
C.19 Data - User Login . . . . .	117
C.20 Data - User Logout . . . . .	118
C.21 View - Get Annotation Details . . . . .	119
C.22 View - Create Annotations . . . . .	120
C.23 View - Get User Details . . . . .	120
C.24 View - User Login . . . . .	121
C.25 View - User Logout . . . . .	121

# List of Figures

3.1	Entity clarification (valid-green, invalid-red, unresolved-gray) using the adapted Wilson score interval with 95% confidence - X-axis: number of upvotes & Y-axis number of downvotes . . . . .	17
3.2	Formula variations for different confidence levels (different z values)	18
4.1	The basic architecture of the back-end. Requests are received by the <i>Controller</i> component and then dispatched into <i>Repository</i> (GET requests) or <i>Service</i> (other HTTP request methods) components. . .	25
4.2	Sample responses from the <code>/api/{data,view}/user/login</code> requests. 25	
5.1	<i>Welcome Page</i> . . . . .	32
5.2	<i>Product List Page</i> . . . . .	33
5.3	<i>Product Details Page</i> . . . . .	34
5.4	<i>Document View Page</i> . . . . .	36
5.5	<i>Document View Page - Annotation Details Part</i> . . . . .	37
6.1	A distribution of the generated annotations per privacy concern category, which is representative of the distribution of the expert based annotations in the OPP-115 dataset . . . . .	42
6.2	Relevant and irrelevant annotations for easy and hard documents per participant . . . . .	44
6.3	Precision, recall and f-measure values per group for annotating and reviewing tasks . . . . .	46
6.4	Precision, recall and f-measure values in total for annotating and reviewing tasks . . . . .	47
6.5	Upvote/downvote ratio counts for a) total annotations, b) PC relevancy and (c) tag relevancy, as voted by the evaluation participants	49
6.6	Depiction of precision and score of grouped annotations . . . . .	50
6.7	Total created annotations and the number of relevant annotations for both relevance metrics grouped by their final score . . . . .	51
6.8	Correlation between user scores and the user precision/recall achieved during annotation creation. The lines on the diagrams correspond to line fittings . . . . .	52

6.9	The final user scores sorted for each group . . . . .	53
6.10	Detailed questionnaire responses given by the users . . . . .	55
7.1	User Interface example design of various planned features . . . . .	60

# Chapter 1

## Introduction

The huge expansion of digital products and the corresponding user generated and gathered data have raised the importance of users privacy and privacy concerns. It is widely known that most mobile applications and websites collect, process and share vast amounts of user data (i.e. personal, contact and location information) without the users being aware of it. This problem has been also noticed in various other digital products (i.e. web services) and even to hardware products (i.e. appliances and smart devices) which adopt the practice of user data collection.

Currently, businesses and organizations around the world are enforced by law (e.g. the EU General Data Protection Regulation - GDPR) to provide information about how customers' data is being treated, usually in the form of privacy policy documents. It is a common case that these documents are mostly written in a formal language with legal terms to ensure the compliance with authorities. On the other hand, the vast majority of users find these documents abstruse and difficult to understand. As a result they skip reading the corresponding privacy policy and use the web applications or services without being aware about their data privacy.

There are several works in the research community that try to address this issue along the following two basic directions:

- Formal Privacy Policy Languages (Readable by Machines)
- Annotation of Privacy Policies with Privacy Related Information

The first direction tries to create a formal privacy policy language/templates to express data management practices that can be readable from machines. Using this language, users and applications/websites can describe their privacy policies and concerns. If they do not match the user can be informed and decide if he/she will proceed in using the application/website. Although this is a legit approach, the available options and implementations that appeared in the early 00's (e.g. P3P), are now considered obsolete and were not widely adopted by users and industry, mainly due to the complexity of the software for the average user.

The second approach mainly focus to extract the basic privacy statements from these document into short sentences, tags, icons and labels for each document in order to make the privacy document easier to be understood by the users. This approach appeared to be not an easy task, since user privacy concerns can vary a lot and the language of these documents can be very ambiguous even for trained analysts and experts. However, lots of progress has been done towards this direction by deploying experts for annotating privacy policies and ML/NLP techniques. Although there is a potential in Machine Learning/NLP approaches, they currently lack the accuracy and refinement of human-based annotations. Specifically, state-of-the-art algorithms for automatic annotation of privacy policies can detect segments of text related to the basic categories of privacy statements and privacy concerns, but unfortunately can not identify the corresponding fine-grained values for the annotated privacy concern. On the other hand, the use of human experts can generate more accurate annotations. The drawbacks of this approach is that the process can be slow for the dynamic and huge environment of web services/applications and may focus on a limited set of services/applications (e.g. the popular ones). A rather small number of works focus on the feasibility of deploying users with less expertise to annotate privacy policies as this could lead to an increment of human work force. Results show that the use of plain users can be accurate enough with the help of various tools.

## Motivation

This study extends the latter direction of work, which revolves around plain users annotating privacy policy documents. To the best of our knowledge, currently there isn't any public, open and standard way for plain users to be enrolled in the process of reading and annotating privacy policies that can operate as a community, independently and outside of a controlled environment. The problem of privacy awareness is a social issue and by enrolling and empowering plain users on this task we expect to spread privacy awareness to citizens as an aftereffect.

## Problem Statement

In this thesis our primary question is to evaluate if by offering a user friendly crowdsourcing platform that gives the opportunity to plain everyday users to create, review and evaluate privacy policy annotations can result in valuable and accurate annotations of comparable quality to those offered by experts users.

## Approach

In this thesis we try to exploit the wisdom of the crowd for supporting users in reading and understanding privacy policy documents. Specifically, we have created a crowdsourcing platform named CAPPA for annotating privacy policies with

information related to privacy concerns. The platform enables users to get detailed information about any data privacy statement that is related to digital products in a direct and user friendly way. Interested users, can create annotations on privacy policy documents based on some predefined privacy concern category labels (i.e. tags) that would simplify the presentation of these documents. Users can review these annotations using a voting mechanism with the aim to filter out bad annotations whereas give prominence to more the accurate ones. We rely on successful Collaborative Awareness Platform (CAP) experiences, such as an engagement mechanism so that users have incentives to create and review annotations. The intention is to create a large, publicly open privacy policy document database reference, enriched with the corresponding annotations, which allows users to create a community around data privacy that will be able to read, discuss, analyze and evaluate the data privacy aspects for any kind of digital products and, more importantly, collectively taking steps towards improving the current problematic situation.

The implementation of such a platform will help users to focus on specific parts of privacy policy. The implemented user engagement mechanism will promote users to read more privacy policies and understand deeper specific text parts that might be in high interest for the majority of users. Moreover, mainstream citizens from various social tiers will be implicitly enforced to actually read privacy policies and receive or even spread an awareness of data privacy

We created a data schema that best represents the current user needs in order to be able to search any privacy statement for any digital product that may collect, use and manage users data. This work is not focused to optimize the various features and aspects that the platform offers. The primary goal is to create a first implementation with a core of features, ideas and mechanisms which can later be extended and evolve taking as an advantage any possible user feedback while using the platform. The main feature of the platform is that entities (i.e. products, documents and annotations) can be evaluated from users as a review mechanism using an up/down voting scheme to pin-point qualitative content and hold back the noisy/erroneous ones. As an engagement mechanism, users are evaluated by their contributions holding a score value label which represents their effort and experience level. The user score is acquired by aggregating some rewards badges that the users receive when they create content or their content is accepted from the community via the review mechanism. Another essential feature for our platform to keep the documents up-to-date. Automatic tracking of document URLs is implemented for tracking and updating the documents whenever a document revision is detected.

Another important feature/aspect of the platform is that it provides a publicly open ReST API for almost the complete set of the data that holds in the database. This feature allows primarily the research community to download and extract data for further analysis and research. Another use of the API is that it allows third party tools to be developed focused in privacy policy analysis and tagging which can further help users privacy concerns.

This work has been implemented to support the umbrella of tools that are currently developed for the CAPrice platform <sup>1</sup>, a collective awareness platform for privacy concerns and expectations.

The conducted user based evaluation shows that indeed plain users are able to provide annotations of high quality and comparable to the expert based ones, which were further refined by the reviewing process. Additionally, most users found the platform easy to use and would consider consulting the platform to get feedback for applications and services that they use, adding the relevant content if it is missing. Finally, the development of a crowd-sourcing platform can indeed improve the user awareness about data privacy.

In the next chapters we present some related work around privacy policies, design and implementation details of the platform and discuss some evaluation results and future work.

In detail, Chapter 2 presents the related work that have been done around privacy awareness, privacy policies and crowd-sourcing.

Chapter 3 unfolds the basic requirements and design decisions we undertook for the CAPPA platform, including data types and schema definitions.

Chapter 4 outlines the architecture of the system and the technologies that were used to implement the system.

Chapter 5 presents some use cases over the platform and provides a brief introduction to the UI of the platform.

Chapter 6 reports an evaluation of the platform having real users annotating privacy policies. The purpose of the test was to identify whether plain users are able to detect specific privacy statements and the level of agreement within the group first feedback test case.

Chapter 7 express some ideas on extending the platform for future work.

---

<sup>1</sup><https://www.caprce-community.net/>



## Chapter 2

# Background & Related Work

In this chapter we will focus on works that address the problem of privacy in general and the privacy awareness of the users. In addition we discuss various works related to privacy policy analysis, the deployed algorithms and the tools that have been developed.

### 2.1 Privacy Awareness

The widely adopted concept for user privacy is based on the notice and choice framework. This framework states that the primary owner of each user data is the user itself and he/she should be informed about any access and any processing to their data. Various studies have extensively examined the existing situation for users [1, 2, 3, 4] and IT systems [5], while some other propose guidelines [6]. Although various legal requirements and regulations for user privacy have been established around the world (i.e. UK Data Protection Act 1998 (DPA), European General Data Protection Regulation 2006 (GDPR), Fair Information Practice Principles (FIPs) etc.) [7, 8, 9] the compliance of business and organizations is limited [10, 11, 12].

### 2.2 Standard Privacy Policy Formats

One main research direction focused on codifying privacy policies into machine readable formats that could standardize the privacy concerns of users. Notable efforts such as the Platform for Privacy Preferences Project (P3P) [13] (an XML format of privacy policies) and the Do Not Track (DNT) flag option (an HTTP header) has been made but they have not been widely adopted since these formats are hard-to-follow and strict for the business and organization needs. As a result website operators avoid to comply with such practices.

On the other hand, the industry itself developed some standards and guidelines to raise consumers confidence when using web applications and services in the form of privacy seals. Examples of major privacy seals are TRUSTe [14], BBBOnline

and WebTrust which they ensure the privacy of the users while using the sealed services.

One study compared various other standardized formats (Privacy Finder<sup>1</sup> and layered notices) with free text privacy policies [15]. Specifically users were asked which format was more usable for them and evaluated how well users were able to understand the privacy policy by answering specific questions. The results showed that although standardized formats are faster to read, they lack the accuracy and precision of the natural language formats.

## 2.3 Privacy Policy Analysis

Free text privacy policies are currently the primary way for users to be informed about the privacy practices that a business or an organization employ regarding the user data. A major argument against the privacy policies, is that their length and complexity makes them difficult and hard to understand by the vast majority of users [16]. This problem is further amplified by the fact that privacy policies use a somewhat formal language, with a lot of legal terms, to order to comply with regulative authorities rather than inform the users [17]. There are examples in the bibliography where even experts disagree on the interpretation of the privacy statements included in privacy policies [18].

### 2.3.1 Readability Metrics

In order to measure the reading difficulty of the privacy policy documents, a number of readability metrics have been proposed in the bibliography. A number of works exploit these metrics for analyzing website privacy policy documents [19, 10, 20, 12]. Other works investigate the variance of readability levels among various business sectors and specific market categories like energy [21], healthcare [22, 23] and social networking [24]. Studies on readability levels of privacy policies for mobile applications are also popular [25, 26]. These studies showcase that the majority of privacy policies are indeed hard-to-understand for the average user.

### 2.3.2 Extracting Privacy Policy Features

To overcome the previously mentioned readability issues of privacy policies, a number of efforts try to extract basic features of privacy policy documents i.e. primary privacy statements for helping users to get a basic information regarding privacy related issues without much effort.

A major project in this direction is the [usableprivacy.org](https://usableprivacy.org/)<sup>2</sup> project [27] which aims to study progress and challenges of privacy policies [28] and develop tools, methods and frameworks that will help users to control their privacy by utilizing

---

<sup>1</sup><http://www.privacybird.org/> - Find web sites that respect your privacy (2005)

<sup>2</sup><https://usableprivacy.org/>

recent advances in machine learning (ML) and Natural Language Processing (NLP) algorithms and crowd-sourcing techniques.

Other approaches towards these directions apply crowd-sourcing methods for annotating privacy statements in privacy policies with a specific set of attributes values pairs [29] and provide an intuitive UI to present them [30]. It is notable that the set of attributes/values that should be used to represent/describe the key elements of privacy policies is of high importance. Various works have proposed some attribute/values sets and ontologies [31, 32] while others focused to find the most common sentences and terms in privacy policies [33].

Various algorithms and tools can assist users in understanding privacy policies. For example Hermes [34] can detect ambiguities in privacy policies texts by identifying the underlying semantic relations between words. Relations between segments and paragraphs among privacy policies have been studied in [35, 36]. While these methods/tools can help users to get better understandings of privacy policies they are difficult to be used to assist inexperienced users.

### **Use of Labels & Signs**

Another branch of work presents privacy policies along with labels, signs and marks. The main hypothesis is that by exploiting such auxiliary visual cues and marks they provide a user friendly way for plain users to understand the privacy statements contained in each privacy policy. For example, the work described in [37] accompanies privacy policies with a nutrition-kind label and reports that this kind of presentation of privacy policies can significantly improve the accuracy and speed of information finding. Another example, the Terms of Service; Didn't read (ToS;DR) [38] uses crowd-sourcing methods to create a collection of community labeled web applications and services by their level of privacy friendliness. Of course the process of manual labeling cannot scale up to cover the set of available applications and services if there is no strong community to support. As a result, machine learning algorithms are utilized in Privee [39] in order to extend and automate the labeling process of ToS;Dr. On the same way/logic/pattern, PrivacyGuide [40] uses ML and NLP methods to generate labels for some important aspects of privacy. Machine learning algorithms and NLP techniques have also been used to specify the completeness of privacy policy documents in [41]. Although these tools can offer summarized reports on privacy policies there isn't an easy way for the plain users to use them.

### **Crowd-sourcing Annotations**

The use of crowd workers for complex tasks has been studied in various works. For example, an approach of how a classification via clustering can occur in high-dimensional data like text is described in [42] while [43] proposes a workflow to create taxonomies with crowd workers. A common approach when crowd-sourcing a problem is to split a complex task into smaller subtasks that are easier to be

solved [44, 45] by crowd workers.

The feasibility of crowd-sourcing for privacy policy annotations has been a subject of study in [46]. In this work crowd workers were asked to highlight specific parts of privacy policies that mention specific information types. The results showed that the accuracy of their work is satisfiable when lots of crowd workers agree but they don't provide a coherent representation of the aggregated results. Another study asks crowd workers to answer specific questions on privacy practices and support their answers by highlighting the corresponding part of text that mentions the specific practice [47]. They confirm that the accuracy of the responses is high when there is an agreement within the crowd workers and they claim that the crowd can perform better with appropriate assistance from tools and algorithms.

Finally, there is a number of works with very interesting results regarding the quality of the work produced by the crowd and the expert users. For example, crowd workers were able to identify more keywords than expert annotators in [46]. Another relevant work is [18], which states that expert users can disagree in the interpretation of privacy policies.

## 2.4 Related Tools

Various text annotation tools are already available like OMTAT annotation tool [48] and GATE [49] but their features are designed for different use cases. A web based text annotation tool for crowd-sourcing has been proposed in [50] but the user inputs and options are limited.

## 2.5 Crowdsourcing Platforms

Crowdsourcing methods are often used in tasks in which humans can perform better than machines. In [51], the authors try to map the notions of human computation and classify them into a taxonomy that could reveal possible improvements for the crowd sourcing tools & platforms. The key concepts and functionalities that crowdsourcing platforms (should) consist of are discussed in [52]. A bright example of a successful crowd sourced platform, the Stack Overflow<sup>3</sup> seems to obey these concepts. Many argue that the success of Stack Overflow is due to the design of the platform while some claim that the daily enrollment of the designers, the active community and the reward system is the important factor [53].

The voting mechanism that a crowd-sourcing platform offers is an effective and efficient aggregating technique. Through highly voted or downvoted content someone can quickly find out quality content or examples where important contributions can be made. Various aggregating techniques can be designed for reviewing the

---

<sup>3</sup><https://stackoverflow.com> - In Stack Overflow users can ask and answer questions, and, through membership and active participation, to vote questions and answers up or down and edit questions and answers

output of a crowd-sourcing platform. Each one of them has its pros & cons, that someone has to take into consideration in order to choose a technique that best suits the current needs[54].

It is notable that some aspects of the platform can have a major impact to the platform's success and can motivate user participation and engagement. One study suggests that being an outlier with your engagement pattern (answer questions with low expertise density, during low peak hours etc.) can result to more reputation to the user [55].

## 2.6 CAPrice

The work described in this thesis was developed to support the umbrella of tools of the CAPrice project<sup>4</sup>. The CAPrice project is a suite of tools that facilitates community interaction and co-creation, enabling the explicit declaration of consumers' privacy expectations of the various digital products. Through a combination of socio-technical methods, such as community-generated design contractualism, crowd sourcing and a knowledge commons approach to privacy policy, the outcome will be a new innovation model that will allow consumers to collectively express their concerns and developers to adopt more privacy-friendly practices and respond to the needs of consumers with novel products and services.

---

<sup>4</sup><https://www.caprice-community.net/>



## Chapter 3

# System Requirements & Design

In this section we will discuss in detail the requirements and the basic concepts of the system, review various aspects and argue over the design decisions that have been taken.

### 3.1 Requirements

The CAPPA platform should follow some commonly accepted guidelines and have a common base/concept with other popular platforms. The task is to create a crowd-sourced platform that any user can easily interact with. Most importantly though we expect our platform to provide a straight-forward way for any user to grasp privacy related information about the various digital products that he/she might be interested in. Borrowing ideas and concepts from other popular crowd-sourced platforms (e.g., stack overflow), our platform evolves around the following three dimensions:

- User Management and Evaluation
- Entity Management and Evaluation
- A Complete Annotation Schema

### 3.2 Definitions

This section describes the basic concepts and the primary components and entities of the CAPPA platform.

#### 3.2.1 Users & User Roles

The aim of the platform is to be publicly open to anyone that wants to search, retrieve, or input information related to privacy policies and their privacy statements. On the other hand, we should provide different access levels to the data,

since a crowd-sourced platform can easily be targeted for various kinds of attacks, including malicious or unwanted input. For that reason, it is required that users should be registered in case they want to contribute to the platform. Registered users should be able to create new content to the platform (i.e., annotations) and/or review existing ones through voting. Of course, the separation of novice and expert users is a must requirement. Newcomers should have a guided introduction with limited privileges, so that they can gradually learn the platform and the offered functionality, and get feedback from other users of the platform.

The idea is that initially, all users should have restricted access to the platforms functionality, until they provide the appropriate cues that they have gained the need experience and are able to contribute valuable content. As users become more experienced, according to the rest users (i.e., we measure experience through a score formula), they will gain access to actions and functionality of crucial importance and impact to the platforms data. On the other hand users might lose grants if it seems that they produce noisy content. Based on this approach, we have defined a user role hierarchy, where each role corresponds to different permission grants, which are described in Table 3.1. In detail:

- **Simple User:** The entry role that each new user has (after registration). This role has permissions for creating annotations and deleting annotations created by the user himself/herself.
- **Elevated User:** This role gives more functionality and permissions. Specifically, elevated users are able to review, evaluate, and vote annotations that already exist in the platform.
- **Content Editor:** This role is the most advanced one a user can have. Content editors are able to do anything that an elevated user can do, plus the ability to create, edit, delete and review entities more crucial than annotations, like privacy policies, documents and products. Of course, each user can edit/delete only entities that are owned by him/her.
- **Administrator:** Administrator is a special role, not available to the normal registered users. An administrator can add, edit, delete any entity that exists in the platform. Administrators are manually set by setting the corresponding values in the CAPPA platform's database.

As we mentioned previously, users can upgrade their role (and their access levels) by gaining experience. In our context experience is gained by creating content (i.e. annotations) and receiving positive feedback (i.e. votes) from other users for these annotations, while the opposite (i.e. down-voting of annotations) removes experience. User experience is measured by the user score which occurs by specific scoring formula. Details on the user score and ratings are provided in the section 3.2.3. In order to upgrade a user role, the user score has to reach a specific score threshold. These thresholds are shown in table 3.2.



User Roles	Annotation			Legal Document			Digital Product		
	$\mathcal{C}$	$\mathcal{D}$	$\mathcal{V}$	$\mathcal{C}$	$\mathcal{D}$	$\mathcal{V}$	$\mathcal{C}$	$\mathcal{D}$	$\mathcal{V}$
Simple User	✓	Owned	-	-	-	-	-	-	-
Elevated User	✓	Owned	✓	-	-	✓	-	-	✓
Content Editor	✓	Owned	✓	✓	Owned	✓	✓	Owned	✓
Administrator	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 3.1: Permission levels for each user role and actions for creating ( $\mathcal{C}$ ), deleting ( $\mathcal{D}$ ) and voting ( $\mathcal{V}$ ) controversial entities

User Role	Next Role	Score To Upgrade
Simple User	Elevated User	10
Elevated User	Content Editor	100
Content Editor	-	-
Administrator	N/A	N/A

Table 3.2: Available user role upgrades and the corresponding upgrade score thresholds

By default, the **CAPPA** platform contains two registered users: the ADMIN with an Administrator role, and CAPRICE as Content Editor. These are predefined users, responsible for handling and managing the initial content of the platform.

### 3.2.2 Controversial Entities

A controversial entity is defined as any user generated entity that can receive votes. This definition allows to maintain a data schema where user input can be evaluated by other users and can be clarified as malicious or legit. In the following subsections we discuss in detail these entities and details about the aggregation of the user voting scheme.

#### 3.2.2.1 Digital Products

Digital Products is a controversial entity type which represents any product that has the ability to collect or manage user data. These kind of products can be any software (i.e mobile or web applications) or hardware product (i.e. appliance, smart devices and objects). The collection of these entities are the entry point for which any user can search details and annotations about the products that he/she is interest in. Each digital product is also tagged with a set of business category labels that this digital product is related to. These labels are exploited for facilitating the retrieval of products of a specific category/label.

### 3.2.2.2 Legal Documents

A legal document is the basic entity that represents any kind of document that a digital product publishes about data privacy related issues to the public. Each legal document has a document type attribute that is used to recognize the distinctive type of the published document. Although privacy statements can be included in many types of documents (e.g., license agreement, terms of service or terms and conditions etc), the current version of our platform only supports 2 distinctive types of documents; the Privacy Policy and the Privacy Notice. The difference between these two documents is that Privacy Policy documents thoroughly describe the actual privacy policies an organization or a business has about its users' data privacy, while the privacy notice is a document explicitly written to report the privacy practices to the users, like a summary of the privacy policy. Privacy policies are so detailed that are also used by developers/employees of the service provider as a guideline for the implementation and the compliance of the respective privacy related functionality and services.

A major issue regarding privacy documents is that they are dynamic and can change any time, sometimes without any prior notification. It's a common case that an organization or a business may change/update the privacy practices that follow the users' data from time to time. When this happens, it is up to the organization or business whether the users will be notified. To address this issue we included in our design a process to detect and update document changes on regular basis in order to keep the platform's documents up to date. Under an interval, the platform compares the document published online with the snapshot available in the database. Whenever an alteration of a privacy policy is found the platform updates the available document by creating a new version of the document in the database. Currently we check the document difference using the last-modified HTTP header date and the text difference in the HTML response.

### 3.2.2.3 Annotations

An annotation is a primitive entity created by the users, that holds privacy related information in the form of tags that users have extracted from the privacy statements described in a document (e.g. privacy policy). It is composed by some (usually one) highlighted parts of text, the labeled tags and an optional comment. Annotations are independent entities that belong to a specific document and only registered users have permissions to generate these kind of controversial entities. Each annotation is related to a specific privacy concern category which restricts the available tag set (attribute/value pairs) that is available to the user for creating the annotation.

### 3.2.2.4 Entity Voting Scheme

Since the CAPPA platform is a crowd-sourced platform and the quality of user input can not be guaranteed, a mechanism to filter out malicious user input and give

prominence to the legit one is necessary. Any user can contribute to the platform by creating new entities (i.e. annotations, documents or products) but there is no guarantee that these entities are legit and valid. The solution to this problem is to allow the community to review any entity and decide as a whole (wisdom of the crowd) whether something is valid or not. Following one of the most popular crowd-sourced platform’s pattern, the stack overflow voting scheme, the **CAPPA** platform adopts the vote up/down mechanism for reviewing controversial entities. Each entity can be up-voted or down-voted by elevated or more experienced users and maintain a score value which results from the difference between the up-votes and the down-votes. This way, users can identify high valued entities and perceive them as legit and important. On the other hand, low or negative valued entities will not be considered by the users as legit but rather as noise.

Although the controversial entity score value resulting by this voting scheme is straightforward and understandable by the majority of users it has some disadvantages. Computing the difference of the aggregation of the up-votes and down-votes completely loses any information about whether there is agreement between the users that have voted. As an example, consider an entity (a) that has received 7 upvotes and 1 downvote and another entity (b) that has received 27-21 votes up/down respectively. In both cases the score value for the entity will result to 6. Based on our intuition though, we consider the entity (a) with score  $7-1 = 6$  more valid than the entity (b) with score  $27-21 = 6$  (b), since users’ agreement in (a) is much higher than in (b), where users opinion seems to diverge.

Another issue that arises when we exploit only the difference of the upvotes and downvotes sums, is that such a metric implicitly depends on the popularity of the entity. For example lets consider what will happen between a popular product (i.e. facebook, google) and a less popular one. Although both products can be equally valid, the score of the popular product is expected to be much higher since more user will search for it and potentially upvote it. Of course this applies for all entity types (i.e. annotations, documents). In all these cases, the assumption that an entity is more valid than another is not accurate.

To resolve these issues, each controversial entity was extended to maintain a separate clarification value, beyond the score of the upvotes and downvotes difference. This clarification value is a string that marks each entity as valid, invalid or unresolved by aggregating the entity votes into a more sophisticated formula rather than just taking the difference of the sums. Clearly, such a formula should consider the total votes that each entity received.

We adapted the Wilson score interval in order to resolve controversial entities. The Wilson score interval is an interval estimate of a success probability  $p$  when only the number of experiments  $n$  and the number of successes  $n_S$  are known. In other words, it tries to estimate the probability of a success for the next trial. This formula takes into consideration the number of trials and results to a possible interval for which the actual probability  $p$  is most likely to be with a given confidence level. The method is based on the central limit theorem with the assumption that the probability of a success (and a fail accordingly) in any trial on

the same experiment is constant. The interval of the actual probability is given by the following equations:

$$threshold_{min} = \frac{n_s + \frac{z^2}{2}}{n + z^2} - \frac{z}{n + z^2} * \sqrt{\frac{n_s n_f}{n} + \frac{z^2}{4}} \quad (3.1)$$

$$threshold_{max} = \frac{n_s + \frac{z^2}{2}}{n + z^2} + \frac{z}{n + z^2} * \sqrt{\frac{n_s n_f}{n} + \frac{z^2}{4}} \quad (3.2)$$

where  $n_f$  is the number of failures (i.e., downvotes in our case) and  $z$  is used to set the desired confidence level. Table 3.3 shows the corresponding  $z$  values that should be used for different confidence levels.

Confidence Level	Z value
70%	1.036
80%	1.281
90%	1.645
95%	1.960
98%	2.326
99%	2.576

Table 3.3: Mapping between confidence level and the resulting Z value

In order to exploit the above formula to match our needs, we model our problem as follows: At first, we view an upvote as a successful trial for our entity whereas a downvote as a failed trial. We deployed the Wilson Score Interval in order to estimate the probability of an upvote (success) in the next trial.

Since every user of the platform can give only one vote per entity, we make the following hypothesis: in the case that all users of the platform vote for an entity, the proportion of upvotes to the total votes is equivalent to the probability  $p$  that the Wilson score interval tries to find. Using this model we get an interval estimate of the actual upvotes that the entity will have if all users vote for it.

We consider an entity as valid iff we are sure at a given confidence level that more than half of the users consider the entity as valid. Under the same logic, we consider an entity as invalid iff we are sure at a given confidence level that more than half of the users consider the entity as invalid. Based on the output range of the Wilson score interval and the previous definitions, an entity is valid iff the min threshold for  $p$  is over 0.5 and invalid iff max threshold for  $p$  is under 0.5 for a given confidence level.

Figure 3.1 shows the landscape between up-votes and down-votes and the resulting entity resolution under the Wilson Score Interval for 95% confidence level ( $z = 1.960$ ) which is a commonly used confidence level in the bibliography and the default confidence level in the CAPPA platform. The X axis represents the number

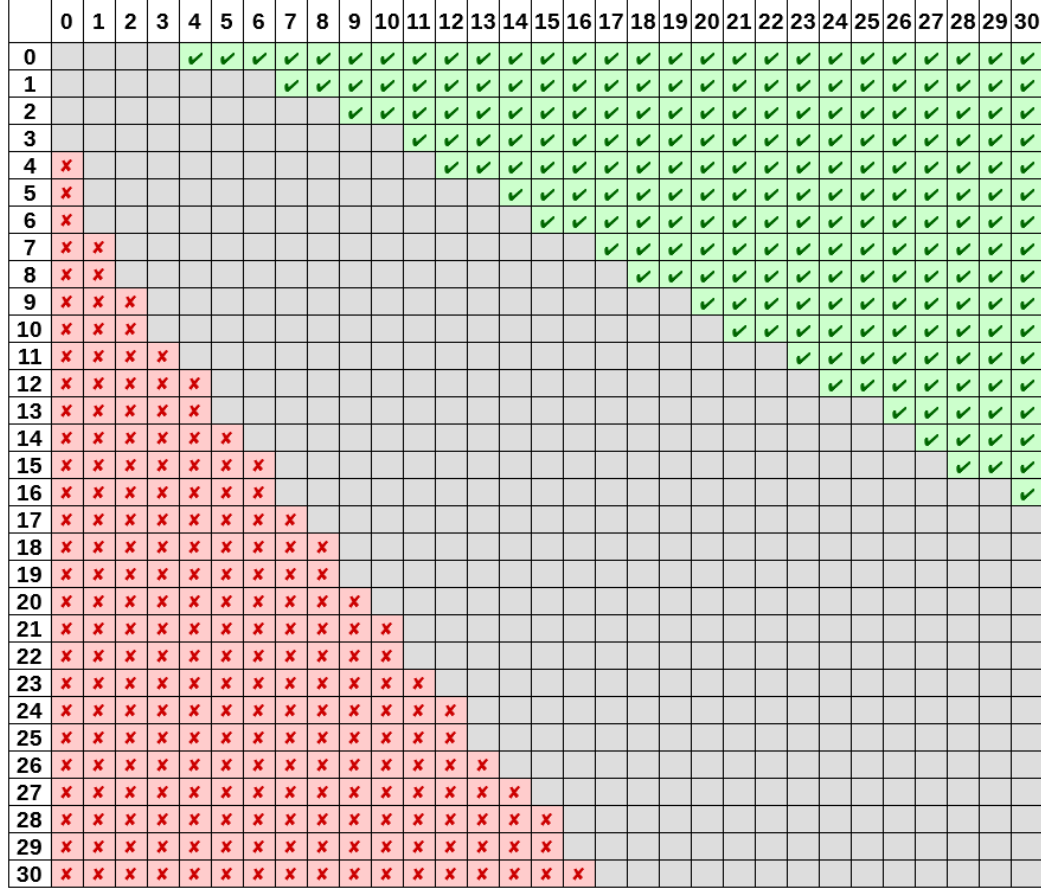


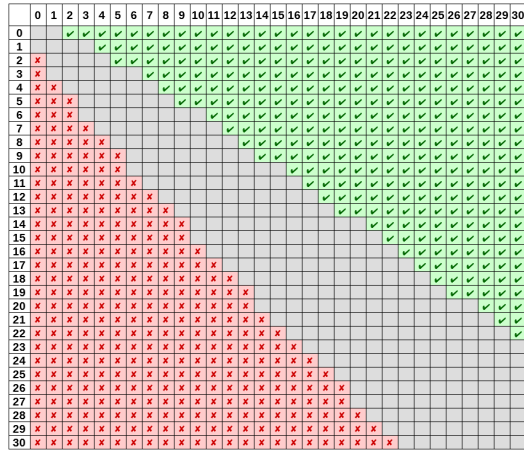
Figure 3.1: Entity clarification (valid-green, invalid-red, unresolved-gray) using the adapted Wilson score interval with 95% confidence - X-axis: number of upvotes & Y-axis number of downvotes

of upvotes, while the Y axis represents the number of downvotes. The green area maps the entity as valid whereas the red area clarifies the entity as invalid. The gray area results to unresolved entities. Figure 3.2 illustrates the various entity resolution maps for various confidence levels.

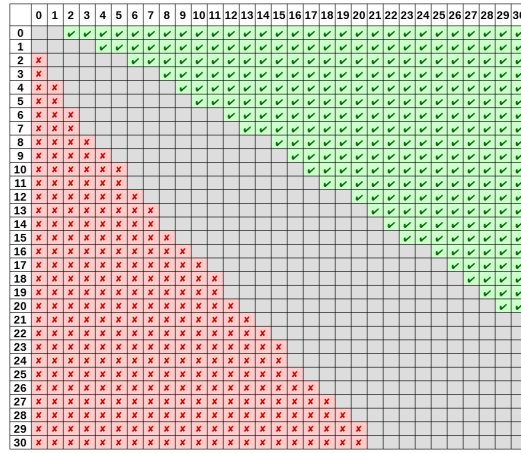
Finally, in order to provide more intuitive representations of the confidence thresholds to better match downvotes as a negative aspect, we mapped the thresholds from interval  $[0,1]$  to interval  $[-1,1]$ .

### 3.2.3 User Ratings & Badges

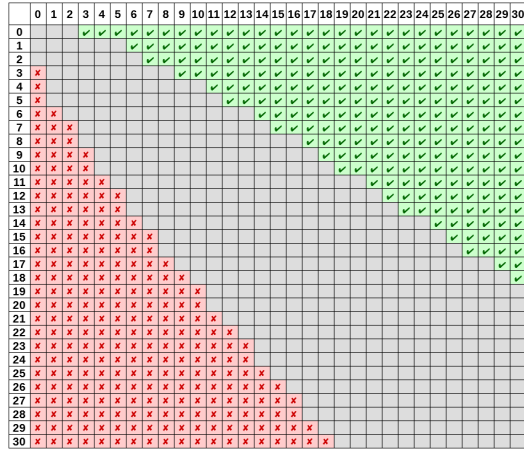
In order to promote user engagement, keeping them enrolled and active with the CAPPA platform in the long term, we need engagement mechanisms and policies that promote users' interest and reward their effort. The users should be able to evolve and progress in a steady pace, feeling that they have something to win and



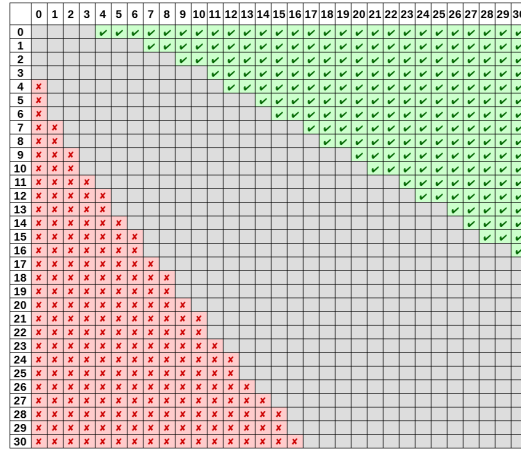
(a) 70% confidence level



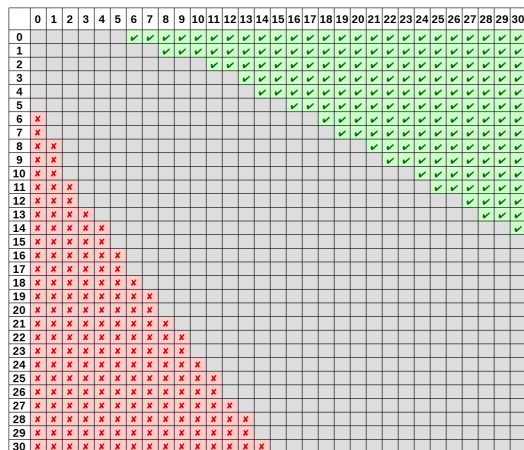
(b) 80% confidence level



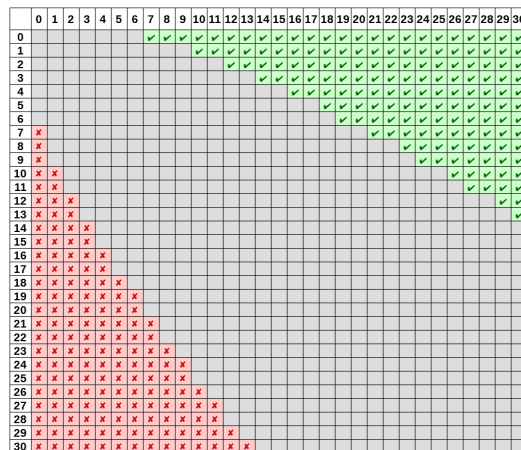
(c) 90% confidence level



(d) 95% confidence level



(e) 98% confidence level



(f) 99% confidence level

Figure 3.2: Formula variations for different confidence levels (different z values)

that their work has an impact to other users or the society in general.

The most widely used way to keep users enrolled long term is by giving them some kind of rewards that represents their effort. We designed and implemented a collection of reward badges that each user can earn to improve his score status within the platform. A reward badge is assigned when a specific task has been achieved. The score of each user is altered whenever a user contributes to the platform or other users recognize his effort. Users can earn various badges of various importance and value.

We have grouped the reward badges into three different categories, based on the different kind of user contribution. In detail, badges are classified into the following three categories - Simple, Special and Goal badges (see Table 3.4).

Badge Category	Description
Simple Badges	Characterizes the most primitive contribution that a user can commit i.e. creating/adding entities, reviewing existing
Special Badges	Express rewards for essential user contributions i.e. a valid entity
Goal Badges	Specify milestones for the user i.e. own 10 valid annotations

Table 3.4: The Reward Badge categories

*Simple badges*, characterize the most primitive contribution that a user can commit. Badges of this type include rewards for simple actions like creating or adding new entities, reviewing the existing ones etc. These reward badges have a minimum impact on user score and focus on keeping users active and engaged in the short-term.

Another type of badges is the *Special badges*. These kind of badges reward users for essential and important contributions to the knowledge created by CAPPa platform’s community. Such a badge is given when an entity resolves into valid or invalid. The value of these rewards is usually higher than the value of *Simple badges*, since the contribution of the user is recognized by other users as valuable/non-valuable and has a great influence.

The third type of badges are called *Goal badges*. They should be considered as milestones that a user can complete in order to get his score boosted. Their purpose is to promote the user interest in long-term.

Using these badge concepts/types we implemented various reward badges for each type. The complete list of reward badges is given in table 3.5.

The final user score is acquired by aggregating the score values from all rewards that the user has received as formula 3.3 describes. The user score turns out to be considered as an experience level and a value of fame by the community. There are many examples of crowd-sourced platforms (i.e. stack overflow) where the user score is considered an important factor that denotes the experience of the

user on a specific domain (e.g., programming). Since a high score can also result to various real-life benefits like getting a job offer, it is an important factor for the engagement of the user.

$$score_{user} = \sum_{i \in badges} value(i) * owned_{user}(i) \quad (3.3)$$

### 3.2.4 Annotation Schema

One of the most important aspects of the CAPPA platform is the way that users create annotations. The selection of a rich annotation schema that is able to describe usable and expressive annotations that fit the user needs and are easy to understand is of a critical importance.

Major data privacy regulation authorities have defined the aspects and the basic principles of the user privacy concerns, so that companies and service providers can produce privacy policies that can address any potential privacy concern of their users. These specifications can constitute a basic guide that we can follow to define an expressive annotation schema for the platform needs.

In our case, instead of defining a new annotation schema from scratch, we used the annotation schema introduced in [29] which has been defined following these principles. This annotation schema has been created by experts and seems to describe the main user privacy concerns. This annotation schema decomposes the users' privacy concerns into 10 categories.

- First Party Collection/Use
- Third Party Sharing/Collection
- User Choice/Control
- Data Retention
- Data Security
- User Access, Edit and Deletion
- Other
- International and Specific Audiences
- Policy Change
- Do Not Track

Every annotation refers to only one privacy concern category like in [29]. Each category defines a specific set of attributes (some of them are mandatory), and each attribute has a specific set of values (attributes can be either single valued



or multi-valued). Details on the attributes/values will be explained in the next subsection.

Due to its length the complete table with the defined privacy concern categories and their attributes/values is given in the appendix (Appendix A).

#### 3.2.4.1 Tag Attributes & Values

The primary aim of the privacy concern categories is to make clear of the different aspects of privacy related information, so that users can easily recognize and classify their concern to any of these categories. This annotation schema is using single attribute-value pairs for each privacy category to describe any highlighted text with privacy statements using these values. Every privacy category stipulates a specific set of attributes-values that tend to answer a basic set of question about the privacy category.

For example, the privacy category "First Party Collection/Use" describes any action that the first party does to collect and use users data. Some defined attributes under this privacy category are the action it applies (i.e. collection or use), the type of information data that is applied to the previous set action (i.e. personal data, computer IP, contacts), the purpose for which this action is done etc. This annotation schema defines some mandatory and some optional attribute-value pairs for each privacy category to further refine the annotations.

In order to comply with this approach, we defined a set of tags under each privacy category where each tag represents an attribute-value pair as defined in the annotation schema. The defined tags are distinguished into mandatory and optional in alignment with the annotation schema. In order to permit the creation and storage of an annotation in the platform, the users must complete all mandatory attributes.

Moreover we relaxed the annotation schema to permit some attribute-value pairs to be multivalued (i.e., allowing multiple tags for some attributes per privacy concern). This way we might possibly reduce the number of annotations needed to describe the privacy statement, since the needed attribute-value pairs (tags) can be added in a single annotation.

A complete list of the available attributes/values (tags) that are available in each privacy concern category can be found in the Appendix A.

Type	Badge	Value	Description
Simple	UpVote	1	Received when the user contributes by giving an upvote to an entity
Simple	DownVote	1	Received when the user contributes by giving an downvote to an entity
Simple	Annot. Created	0	Received when the user contributes by creating a new annotation
Simple	Doc. Added	0	Received when the user contributes by adding a new document
Simple	Prod. Added	0	Received when the user contributes by adding a new product
Simple	Annot. UpVoted	2	Received when someone votes up an annotation created by the user
Simple	Doc. UpVoted	2	Received when someone votes up a document added by the user
Simple	Prod. UpVoted	2	Received when someone votes up a product added by the user
Simple	Annot. DownVoted	-1	Received when someone votes down an annotation created by the user
Simple	Doc. DownVoted	-3	Received when someone votes down a document added by the user
Simple	Prod. DownVoted	-5	Received when someone votes down a product added by the user
Special	Annot. Validated	2	Received when an annotation created by the user resolves into valid
Special	Doc. Validated	5	Received when a document added by the user resolves into valid
Special	Prod. Validated	10	Received when a product added by the user resolves into valid
Special	Annot. Invalidated	-2	Received when an annotation created by the user resolves into invalid
Special	Doc. Invalidated	-5	Received when a document added by the user resolves into invalid
Special	Prod. Invalidated	-10	Received when a product added by the user resolves into in valid
Goal	#Valid Annot. 1	1	Received when the user reaches to 1 valid annotation in total
Goal	#Valid Annot. 10	10	Received when the user reaches to 10 valid annotations in total
Goal	#Valid Annot. 50	50	Received when the user reaches to 50 valid annotations in total
Goal	#Valid Annot. 100	100	Received when the user reaches to 100 valid annotations in total
Goal	#Valid Annot. 500	500	Received when the user reaches to 500 valid annotations in total
Goal	#Valid Annot. 1000	1000	Received when the user reaches to 1000 valid annotations in total

Table 3.5: Predefined badges with the score value gain for the user score

## Chapter 4

# System Architecture & Implementation

### 4.1 System Architecture

The platform's maintenance is heavily based on the code organization & structure. A clean architecture of the core code structure is a vital requirement, so that other developers can contribute and extend the platform and its features<sup>1</sup>. The platform's architecture has been organized into different layers and components, in order to decouple code with different roles.

#### 4.1.1 Front-End

The platform's front end is organized into two main packages - the fragments and the pages. The pages package contains the landing page (root) HTML files. These files are based into smaller pieces of UI elements in order to be rendered and they serve as containers for more primitive UI fragments. It is expected that a page file represents a specific URL pattern.

The fragments represent UI components that the pages include, usually as static files, in order to build the final page. There are cases that some fragments should dynamically be rendered in order to be returned as responses of the view (`/api/view`) REST sub-API. These fragments are implemented mostly in the form of custom tags.

#### 4.1.2 Back-End

The back-end has been organized into three main components:

- *Controller*, responsible for receiving the requests
- *Service*, contains all the business logic

---

<sup>1</sup>The code of the platform is available as open-source

- *Repository*, responsible for all database I/O

The *Controller* package, contains all the classes that are responsible to receive the incoming HTTP requests and apply a basic error checking on requests' parameter data including form validations and resource checking. If there are not errors, the request is forwarded to the appropriate component for processing, either the *Repository* component in the case of HTTP GET requests or the service component for other HTTP methods. Specifically, when the method of a request is GET, which means that it's a read-only request, the request is directly dispatched into the *Repository* component which handles all database access, since the execution of the request is straight-forward. For create, update or delete requests (i.e. POST, PUT, DELETE methods) or when the request is more complex and demands some kind of processing, the request is firstly guided through the *Service* component, which contains the needed implementation of the business logic. In addition, the service package is responsible for a second layer of error checking like permission checking and entity resource resolutions. As figure 4.1 illustrates, all database I/O access for the *Service* component are (as in the *Controller* component) managed by the *Repository* component. Last but not least, as already mentioned, the *Repository* package is responsible for all database I/O. It is used by the *Service* and the *Controller* components to perform the needed data retrieval and updates.

## 4.2 REST API

As discussed previously, a major aspect of the platform is the REST API it offers. Anyone can use HTTP requests to the platform to receive data contained in the database in the form of JSON format. All API requests are directed under /api using intuitive URL patterns for various entities.

The API is divided into 2 sub-APIs - the data and the view API. In detail:

- **data:** The **data API** is located under /api/data/. The corresponding responses are in JSON.
- **view:** The **view API** can be found under /api/view/. Each response includes an HTML part than can be directly embedded into the clients page. The server responses are also in JSON format having the rendered HTML data in a predefined field.

Response samples of these two sub-APIs can be found in figure 4.2. In the case of **view** API format the resource is included in the form of a rendered HTML form whereas in the case of the **data** API format the resource is included as JSON data. The basic resources entities are included under both sub-APIs.

- **User:** Functionality that refers to users can be found under /api/{data,view}/user
- **Controversial Entity:** Requests about voting and deleting entities can be found under /api/{data,view}/controversialEntity

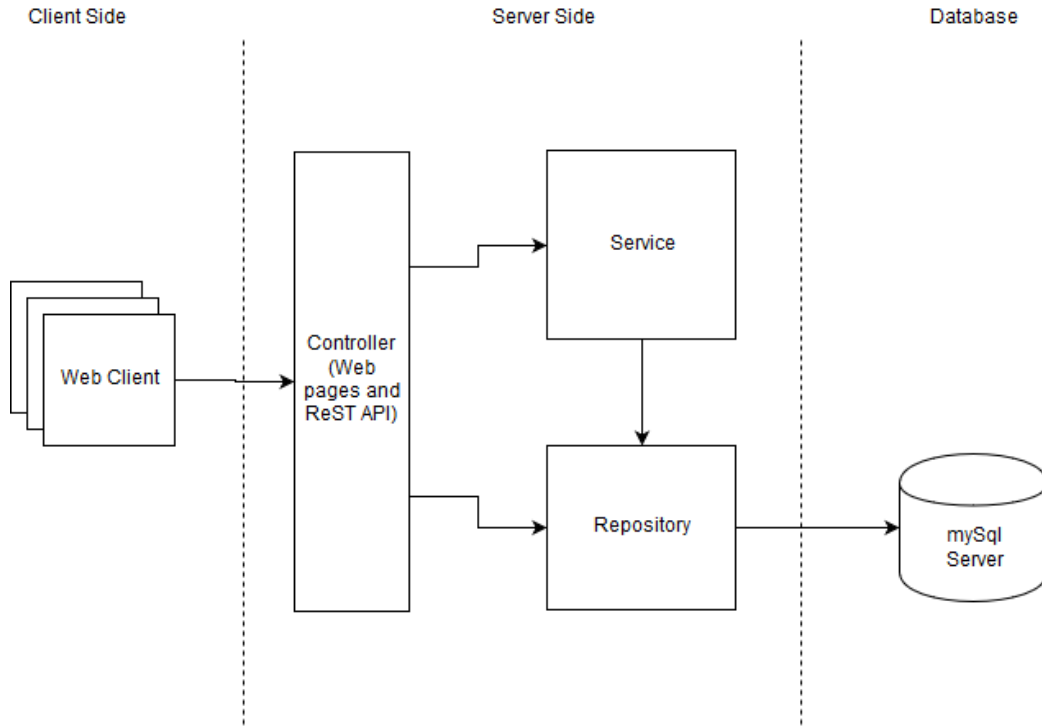


Figure 4.1: The basic architecture of the back-end. Requests are received by the *Controller* component and then dispatched into *Repository* (GET requests) or *Service* (other HTTP request methods) components.

```

{
  "status": "success",
  "message": "Login Successful",
  "html": "\n\n\n<script>\n    function log
    /logout',\n        method: 'POST'
    jqXHR){\n        if( response
    .loginInfo').html(response.html);\n
    2000);\n        }\n
    <h1 class=\"d-inline\"><i class=\"fa
    \n    <div class=\"d-inline-block ali
    -light v-middle dropdown-toggle\" hre
    -toggle=\"dropdown\" aria-haspopup=\"
    &nbsp;(0)&nbsp;</span>\n    </a>
    =\"dropdownMenuLink\">\n    <
  }

```

(a) Sample response from a /api/view/user/login request

```

{
  "status": "success",
  "message": "Login Successful",
  "data": {
    "id": 2,
    "age": 45,
    "createTime": 1512654456000,
    "details": null,
    "educationLevel": "Phd",
    "email": "test@test2",
    "profession": "Developer",
    "username": "user1",
    "score": 0,
    "rewards": [],
    ...
  }
}

```

(b) Sample response from a /api/data/user/login request

Figure 4.2: Sample responses from the `/api/{data,view}/user/login` requests.

- **Product:** Retrieval and management for products can be found under `/api/{data,view}/product`
- **Document:** Retrieval and management for documents can be found under `/api/{data,view}/document`
- **Annotation:** Retrieval and management for annotations can be found under `/api/{data,view}/annotation`

Various other URL paths are defined in order to support other functionality that was needed for the platform implementation. The complete documentation of the provided API can be found in the appendix (Appendix C).

The API has been designed so that it can be easily extended to support more functionalities and features in the future. It provides great flexibility since with the use of the API anyone can build third-party tools, analyze the database or even offer a different user interface (front-end).

### 4.3 Modularity of the Annotation Schema

Since this is a research platform, it is expected that the data values and schema will evolve rather fast. The platform's modularity level meets this criteria since all the data type definitions with their possible values (i.e. privacy concerns, tag attributes/values, user roles, etc.) are defined as values in the database. More specifically we can completely redefine the annotation schema (i.e. tags used to create new annotations) by simply changing values in the database.

### 4.4 Implementation

The implementation of an adaptive and scalable system relies heavily on the underlying technologies used for building it. The selection of technologies plays a crucial role in keeping the platform easily expandable and adaptable to design changes while new features are being deployed. Moreover, in order to keep the platform alive, to provide support in the near future (i.e. bugfixes and implementation of new features) and reduce the learning curve for other developers to contribute, we have built the CAPPa platform using popular and mature technologies. Developers with experience in popular technologies are more common than for technologies with poor support and usage. Later on in this section we review and compare existing popular technologies and justify our decisions for choosing each.

#### 4.4.1 Server Side vs Client Side Rendering

Since the birth of the Internet, the common way to get the rendered view show up in browser was to request it from the server. This method has started to fade out in the last years having client side rendering taking its place. The reason

for this change is that applications evolved into highly reactive pages demanding components and various parts of the page to change their content in real time. Although building a client side rendering provides direct and rich interaction with the user, it has a major disadvantage. The Search Engine Optimizations (SEO) capabilities of client side applications are very low since the search engines are failing to get the actual content of the page in most of the cases. On the other side, server side rendering with an appropriate REST route design can make the platform's documents and products easy to find by search engines. Currently our platform's functionality is closer to a DBMS system rather than an interactive application. We chose to adapt the server side rendering pattern for most of the pages meeting the requirements of a REST design whereas we implemented a hybrid design for the document view page since the user has a variety of actions there (i.e. create, delete or vote an annotation, filter and reorder annotations, etc.) getting out the most between the two architectures. In case the requirements change in the future, the modular design of the platform allows to alter this pattern in future versions without too much effort.

#### 4.4.2 Technologies Used

The development of large platforms rely on various frameworks in order to manage the complexity and the different needs of the various components and their functionality. In the case of the CAPPA platform, the development of the platform was based on a number of popular frameworks, both for the front-end and the back-end.

##### 4.4.2.1 Back-End

The platform relies heavily on the back-end, since most of the features pertain to data management. We implemented the back-end of the CAPPA platform based on the use of Java EE<sup>2</sup> coupled with the popular Spring framework<sup>3</sup>. Spring offers some high standard technologies and it is widely used for enterprise and business applications. The Spring Boot<sup>4</sup> module of the Spring framework provides some fast server setup and deployment tools which provides a boost on the development side. Spring Data JPA<sup>5</sup> which uses Hibernate<sup>6</sup> tools allows any relational data schema to be mapped to Java objects. Spring, Spring Boot, and Spring Data JPA are all under the Apache License 2.0, while Hibernate is under the LGPL 2.1 license.

---

<sup>2</sup><http://www.oracle.com/technetwork/java/javaee/overview/index.html>

<sup>3</sup><https://spring.io/>

<sup>4</sup><https://projects.spring.io/spring-boot/>

<sup>5</sup><https://projects.spring.io/spring-data-jpa/>

<sup>6</sup><http://hibernate.org/>

#### 4.4.2.2 Front-End

For the front-end we used some of the most popular frameworks for creating web user interfaces. Specifically, we used Bootstrap 4.0<sup>7</sup> that offers a variety of ready-to-use UI components alongside with jQuery 3.2<sup>8</sup>, which exponentially reduces the JavaScript code needed for the web client side. Pairing these two technologies is a popular/mainstream option for front-end development and the basis for extending it with more advanced libraries and frameworks. The SaSS project<sup>9</sup>, which is a CSS extension language has also been exploited for describing the presentation of the web pages.

#### 4.4.2.3 Third Party Tools & Libraries

We used also some other third party libraries (mostly in the front-end) that are worth mentioning.

- **bootstrap-notify**<sup>10</sup> is a JavaScript library that we used to create UI pop-up information messages for the platform user. The library is distributed under MIT license.
- **selectize**<sup>11</sup> is a JavaScript library which offers an advanced and customized user input field. It was used in various form input fields of the platform where tag-labeled input from the user was needed (i.e. product or document types, annotation tags, etc.). This library is distributed under Apache v2.0 license.
- **flag-icons**<sup>12</sup> is an easy-to-use country flag icon set. It was used to display the language of each document. This library is distributed under MIT license.
- **font-awesome**<sup>13</sup> is one of the most popular vector icon set. This library is distributed under MIT & OFL 1.1 license.
- **text-highlighter**<sup>14</sup> is a JavaScript component used for text highlighting of documents. This library is distributed under MIT license.
- **pagemap**<sup>15</sup> is a JavaScript UI minimap generator component. It was used for showcasing the highlighted text of each annotation in a small minimap. This library is distributed under MIT license.

---

<sup>7</sup><https://getbootstrap.com/docs/4.0/getting-started/download/>

<sup>8</sup><http://jquery.com/download/>

<sup>9</sup><https://sass-lang.com/guide>

<sup>10</sup><http://bootstrap-notify.remabledesigns.com/>

<sup>11</sup><https://selectize.github.io/selectize.js/>

<sup>12</sup><http://flag-icon-css.lip.is/>

<sup>13</sup><https://fontawesome.com/v4.7.0/license/>

<sup>14</sup><https://github.com/mir3z/texthighlighter>

<sup>15</sup><https://larsjung.de/pagemap/>



- **ipeirotis/readability metrics**<sup>16</sup> is a Java project that contains implementations of various document readability metrics. This library was used in the platform's back-end, to compute the readability metrics of each document (i.e., privacy policy). This library is distributed under Apache v2.0 license.

**NOTE:** All mentioned libraries and projects were used as on top level dependencies, but their contributions are offered by their dependencies also.

---

<sup>16</sup><https://github.com/ipeirotis/ReadabilityMetrics>



## Chapter 5

# Use Cases

This chapter demonstrates in detail typical system use cases, showing how a user can search for a digital product, read the available privacy policies, create, check or review the available annotations, etc. In parallel we provide some screenshots along with the corresponding descriptions of the various parts of the user interface (UI).

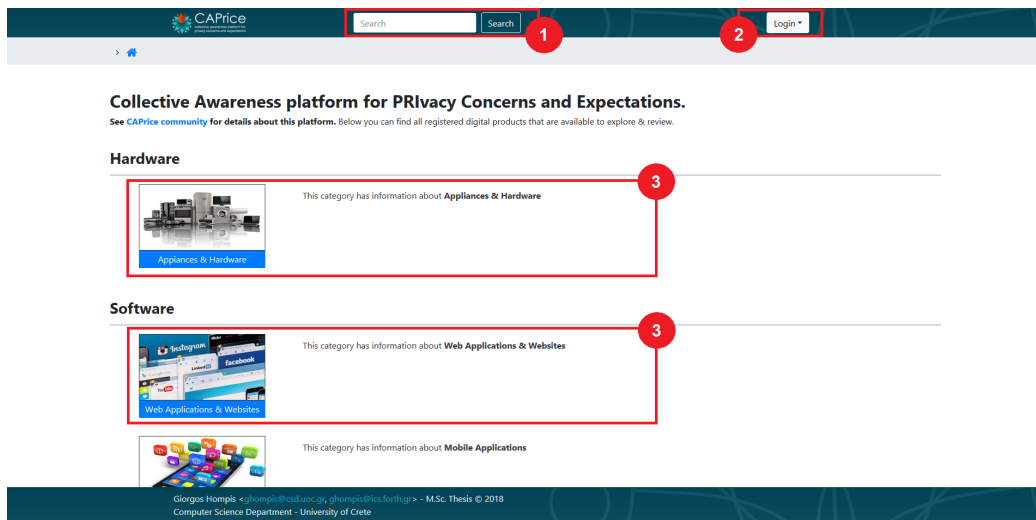
A user-friendly, straight-forward, self-descriptive and efficient<sup>1</sup> UI is a significant requirement for the success of the system, since the long term interaction and use the platform by the users and their engagement is done through various UI components. Efforts has been made to get an efficient and user-friendly UI. Based on the evaluation though(see 6), there is still space for improvements.

### 5.1 Welcome Page

Figure 5.1 depicts the *Welcome Page* of the **CAPPA** platform. The part labeled with (1) allows the user to search for any product that is available in our platform by its registered name or URL. The part labeled with (2) enables the user to view his login status, login to the platform or create a new account to the platform. Both (1) and (2) parts are located in the top navigation bar of the platform. This functionality is available to the user in all other pages. The landing page introduces the user to the platform with a short description and offers the 3 different categories of digital products that are available in our platform (the predefined product categories are labeled with (3)). This design allows any user to search for the desired digital product or navigate through all the available products in the platform.

---

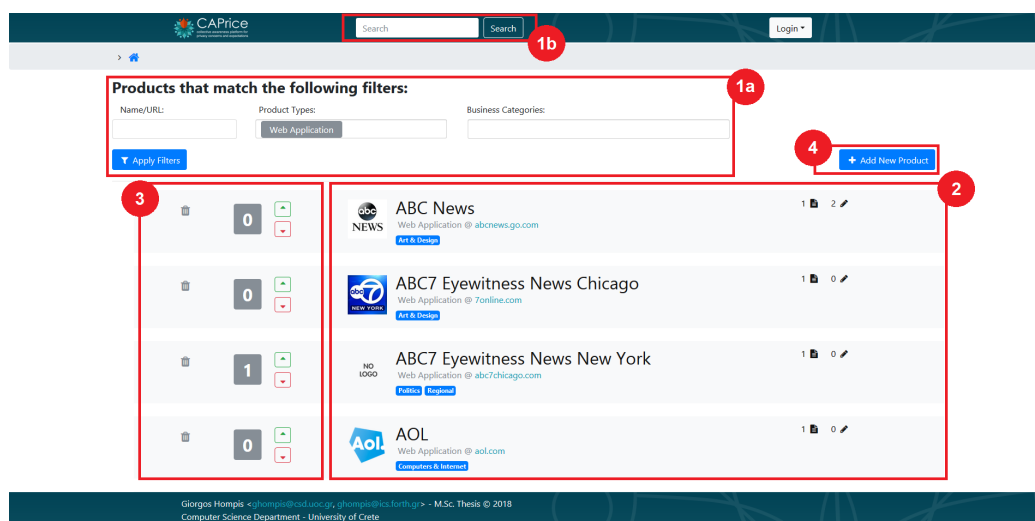
<sup>1</sup>Efficient UIs require a minimum number of clicks, mouse movements, idleness, etc., to complete some basic tasks

Figure 5.1: *Welcome Page*

## 5.2 Products List Page - Search Results Page

The list of available products can be considered as one page that contains the results of a search. By selecting one product category the user will land on the corresponding *Search Results Page* that contains all the available products of that specific category. Of course, a user can refine the search criteria by using the tools provided in the area marked as (1a) in Figure 5.2. The user can search a product by providing a string in the corresponding form field (either its name or its URL), and can exploit business category labels to filter out the results. Moreover the user can always use the search field in the header bar (1b) to quickly navigate into the search results. The header field input is used as a name or URL search string and it is equivalent with the name/url field that exists in (1a). The user can navigate through the results list and view the available list of products (2). Each product item contains information like the name, logo, product type, URL, and business category labels, that helps him/her to identify product he is looking for. Moreover, each product includes information generated by the platform and its users, like the entity score, the number of available documents (e.g., different privacy policy versions) and the total annotations this specific document contains. When an entity has been reviewed by the users of the platform and there is a consensus about its validity/invalidity clarified, the score label in box (3) will turn from gray to green or red respectively.

The search results list also offers some functionality for logged-in users only. If the user wishes (and has the necessary permissions), he/she can review (i.e. vote) or delete the product from the platform by clicking the available buttons showed in (3). In case of a missing product, the user can insert it to the platform by clicking the button displayed in (4). The user has to insert the relevant product

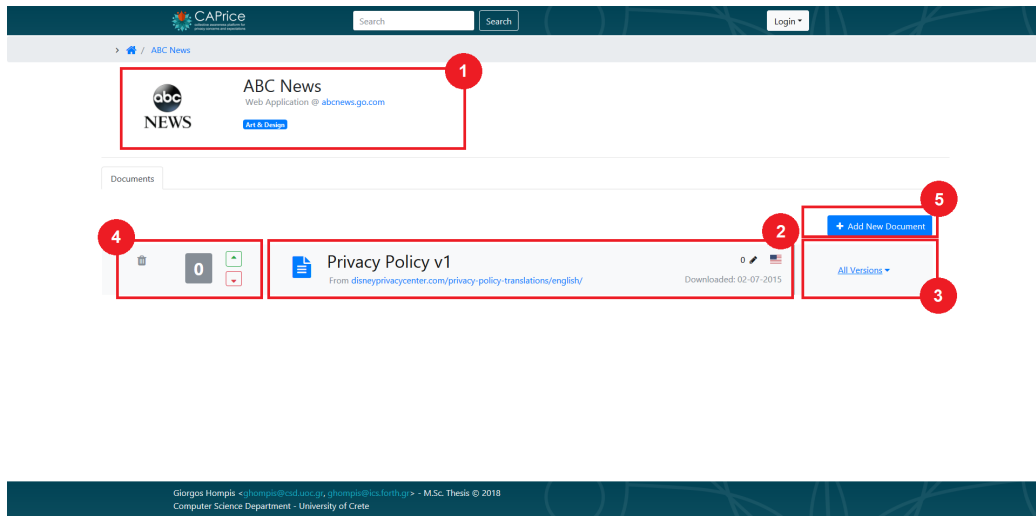
Figure 5.2: *Product List Page*

information to the corresponding inputs of a modal form. By submitting the form, if there are no errors, the product will be added to the platform.

### 5.3 Product Details Page

In the *Product Details Page* the user can among others view the available documents for a specific product. The document list consists of various document types related to the product at hand, that may contain privacy statements about what, how and why user data are being collected and used. These documents, as previously described can be either Privacy Policies or Privacy Notices. All documents are identified by their URL. The platform tracks the URLs and periodically checks and updates each document creating a new version, if they are different, in order to maintain different document versions. By clicking the link showed in area (3) of Figure 5.3 all versions of a document will be displayed (in descending order based on their downloaded time stamp), so that the user can select the document version he wants to read and review. Specifically, in area (1) detailed information about the product is displayed. The user can also navigate to the link of the product item. The area labeled with (2) presents detailed information about each document, like the document type, the internal version number, the document URL from where the document was downloaded, the date the document was downloaded, the document language and an indication of the total annotations that the document contains. This area is click-able and navigates to the actual document view.

Analogously to the *Products List Page*, some actions are only provided to the logged-in user. When permissions are granted, the user can review each document entity (i.e. vote) or delete it using the buttons illustrated in area (4). Furthermore,

Figure 5.3: *Product Details Page*

in case of a valid/invalid document entity, the score label in area (4) will turn from gray to green for a valid entity and red for an invalid.

Last but not least, if the user wants to add a new document (for example when the document URL changes) he can click the button shown in area (5). By clicking the button, a form will show up requesting some input fields to be completed. Submitting the form with no error, will add the document to the platform under the specified product.

## 5.4 Document View Page

The most important page in the system is the page that displays the actual document. It allows to read and explore the actual document content and provides functionality to create and review annotations. By utilizing a color palette for different annotations types (as shown in area (3) of Figure 5.4, the process of reviewing the actual document is rather easy. Our designed UI for the document view page was influenced by the UI offered in [30].

### 5.4.1 General Functionality

As depicted in Figure 5.4, the main part of the document view page displays some document information, 3 mode buttons and three panels parts. In the area labeled with (1), the user can find details about the current document. These details include the logo of the product, and the type, the version and the URL of the document, along with the date that the document was downloaded. In addition, in the right of area (1), there is an indication about the document readability level for various readability metrics. The primary metric that was deployed for

the readability level of the document is the Flesh-Kincaid Grade Level (FKG), but other readability metrics are also supported by clicking the corresponding link. These metrics include:

- Flesche Readability Ease Score
- Simple Measure of Gobbledygook
- Coleman-Liau Index
- Flesh-Kincaid Grade Level
- Automated Readability Index
- Gunning Fog Index

For simplicity purposes, each metric score is mapped into a predefined 7-scale readability level - *Very Easy, Easy, Somewhat Easy, Standard, Somewhat Hard, Hard, and Very Hard*. These metrics can be quite useful for inexperienced users since they can get some kind of indication about the document readability difficulty and proceed accordingly into reading the text and create new annotations or review annotations that are already there.

An important aspect of the platform is the three mode buttons available in area (2). These buttons represent different view/functionality modes (*View, Review, and Annotate*) for the document, that the user to focus on a specific task.

For example if the user is interested to just read the document, it is helpful to avoid noisy annotations that might distract him from this task. Specifically, the *View* mode allows only valid annotations to be visible over the document. Highlighted parts of text from unresolved or invalid annotations will not be shown when this mode is selected.

By default, the *Review* mode is selected. The *Review* mode does not filter out any annotations. The intention is to let users explore and review the complete set of annotations that are available in the document.

In case a user wants to create a new annotation, the *Annotate* mode should be selected. In this mode, the document is cleared from any highlighted parts. The user can select and markup the relevant parts of text that should be included in the annotation that he/she is currently creating. By cleaning up any possible noise included by other text highlights helps the process of annotation creation.

The panel appeared in area (3) can be considered as an index for the annotations that appear in the document. Someone can use this index to investigate the number of annotations for each privacy concern category. Each privacy concern category can be expanded showing all the available tag attributes and values that have been used. By clicking any of these tags, the document (area (4)) scrolls to the first available annotation highlight, making the navigation from tags to the corresponding annotations rather simple. A second click to the same tag, brings into focus the next relevant annotation, using a round-robin based formula. Additionally, each privacy concern is color coded, to help users identify quickly the

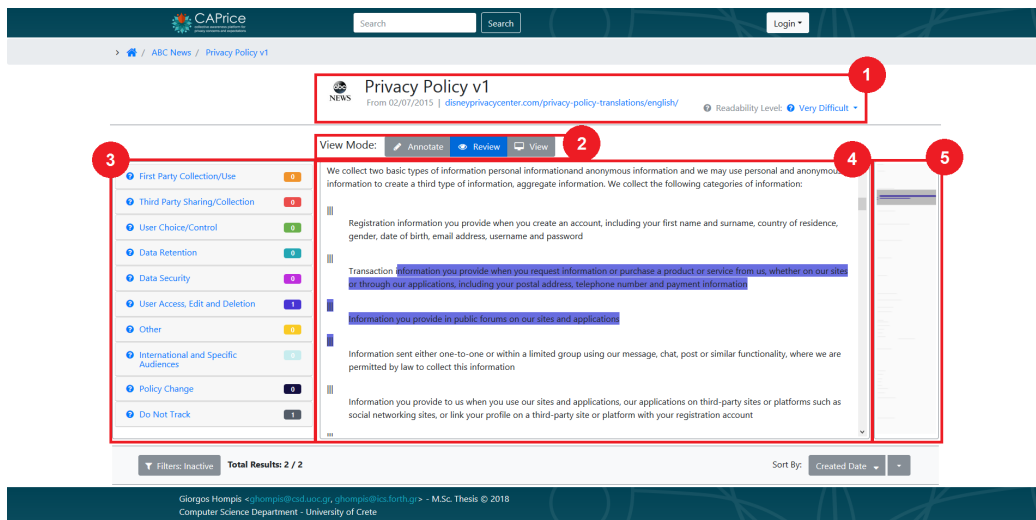


Figure 5.4: Document View Page

privacy concern category of an annotation. The same color is used also for the highlighted text and in the minimap.

The main document text is displayed in area (4) bounded by a fixed-sized scrollable div. This way the document can be scrolled (while the user reads the document) leaving a steady window with all the available accessories in a fixed place that will not confuse the user.

Finally, a document minimap was deployed on the right part of the document (area (5)) which gives insights to the user about the document state and the placement of the annotations in it. Each annotation is displayed in the minimap using its predefined representation color which makes rather simple for the user to identify the type of the annotation. The minimap can also be used to scroll and focus on specific parts of the document since its click-able and scrollable.

### 5.4.2 Annotation Details

On the bottom part of the page, someone can find the list of annotations along with their details, displayed as colored text highlights. Figure 5.5 shows the relevant part in the *Document View Page*.

The number of annotations in each document can be high since anyone can create annotations. A high number of annotations in the document will result in many (may crossed) colored parts of text in the document which can frustrate the user as a side-effect. For that reason, some annotation filters has been implemented as shown in area (1). This button offers to the user some input fields which can be used to refine the set of annotations that are visible. As a result the user can focus on a specific subset of annotations.

Another way to control the list of annotations, is to sort them using some



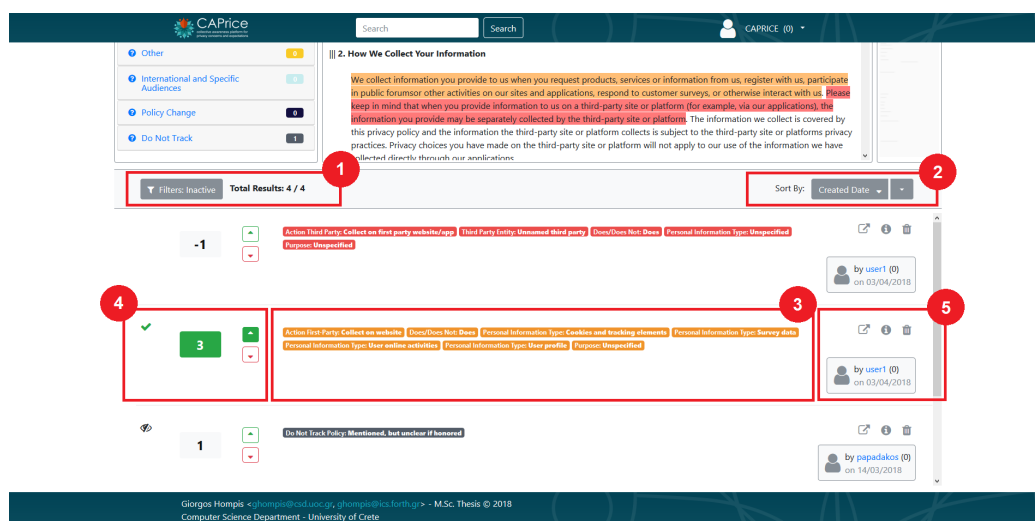


Figure 5.5: Document View Page - Annotation Details Part

criteria (area 2). By default, the annotation list is ordered based on the creation date of the annotation. Other options like ordering by entity or user score are also available.

As described in the previous chapters, an annotation entity consists of the highlighted text, some tags and an optional text comment. These details are displayed on area (3) for each annotation in the list.

In correspondence with other entities, annotations entity score is shown in area (4) together with the corresponding vote buttons that allow users to review each annotation (if they have the corresponding permissions).

In the area (5), someone can find some annotation metadata (i.e. when the annotation was created and by whom) along with some button icons. The first two button help the user to find the highlighted text part, either by focusing the document on the highlighted text of the annotation (the first button) or by showing a popup with the highlighted text (the second button). The functionality of the last button is to delete the annotation (when permissions are granted).



## Chapter 6

# Evaluation

The previous chapters discussed in detail the design of the platform and its features and showcased the final user interface. In order to measure the platform effectiveness and various other features that it supports we conducted a user based evaluation. The primary aim of the evaluation was to measure the quality of annotations that plain users can generate and whether the crowd (as a community of users) can give prominence to the most valuable ones. Other aspects of the platform (i.e. the entity scores and the user scores) were also evaluated/measured in order to get feedback to be able to make targeted improvements in the future.

Although the evaluation was conducted in a controlled and a small scale environment, it can be considered as a first effort to get some very useful feedback around various aspects of the platform and user impressions. In our case, the results appeared to be very encouraging, showcasing that non-expert users are able to provide high quality annotations, while noisy or non-refined annotations can be pin-pointed and downvoted through the reviewing/voting process. As a result the recall and precision of the valid annotations (as voted by the crowd), were rather satisfactory compared to the annotations offered by the experts, reaching in some tasks the ideal score. On the other hand the participants provided valuable feedback regarding the user-friendliness of the platform and suggested a number of improvements as future work.

In the next section we will describe the setup of the evaluation, some details about the users that participated, the process and the environment of the evaluation. Later on, we analyze the results and discuss our findings.

### 6.1 Evaluation Setup

The evaluation took place in a controlled environment where 12 highly educated users (6 with a MSc degree and 6 with a PhD degree) were asked to create annotations and then review the annotations created by other users. At first the participants were given a short tutorial about the platform and a simple example of how to create annotation. Then they were separated into 2 groups (i.e. group

A and group B). Each group was asked to annotate 2 privacy policies, an easy and a hard one as indicated by the Flesh-Kincaid Grade (FKG) readability metric. All privacy policies were taken from the OPP-115 collection, for which privacy experts have provided annotations (the ground truth in our case). At first each group annotated the corresponding easy privacy policy, and then each group reviewed the annotations of the other group. Then they were asked to do the same process also for the hard privacy policy as indicated by the FKG readability score (annotate and then review the annotations of the other group). Finally they were asked to complete a questionnaire about different aspects of the platform. The whole process was split into 4 main phases (i.e. A, B, C and D) and lasted about 2 hours.

In detail, the first phase (A) of the evaluation task was to introduce/showcase the platform to the users. Specifically, we provided a short tutorial about the platform usage and an introduction to the available privacy concern categories and tags. Participants were asked to inspect the available categories and tags and try to create an annotation based on a simple privacy policy excerpt. The duration of this phase was around 20 minutes.

In the second phase (B), we split the *IronHorseVineyards.com* privacy policy, a long but easy document according to the FKG readability metric, into 2 sets of sections/paragraphs. Initially, the first group (group A) was asked to create annotations for the first set of sections/segments, whereas the second group (group B) was asked to create annotations for the second set of sections. Each group had about 20 minutes to create the corresponding annotations and both groups worked in parallel. Then, each group was asked to review (vote up/down) the annotations that had been created by the other group, again with a time limit of 20 minutes. The whole second phase lasted about 40 minutes and resulted to a user-created set of 69 annotations. The review process provided a set of 12 valid annotations and 6 invalid.

A similar process was followed for the third phase (C) of the evaluation. This time each group was given one hard document (*Mohegan Sun* and *Restaurant-News.com* privacy policies respectively), and were firstly asked to create the respective annotations (within 20 minutes) and then review the annotations created by the other group (also within 20 minutes). The two groups created 33 and 37 annotations respectively, while the review/voting process produced 7 valid and 6 invalid annotations for group A and 5 valid and 3 invalid for group B respectively.

For the forth and final phase (D) of the evaluation task, the participants were asked to complete a questionnaire in order to give feedback about the platform design/features, the user-friendliness of the platform, the difficulty of the requested tasks, the expressiveness of the privacy concerns, attributes and values, and their final impressions and comments.

In the next sections we will try to analyze and discuss a) the user created annotations, b) the valid/invalid annotations provided by the voting processes and if the correlation between the annotations' scores and their quality/value. Finally, we will discuss the results of the questionnaire about the aspects of the platform

mentioned earlier (.e.g., user-friendliness, etc.).

## 6.2 Annotation Quality

The quality of the annotations that are created and reviewed in the CAPPA platform may dependent on a number of factors, and since it is a crowd-sourcing platform it depends on the wisdom of the crowds, which is manifested by aggregating the efforts of a lot of users. Our evaluation process was not a direct representative of the real-life usage of the platform, due to the limited time constraints, the small number of users and documents, and the in parallel creation of annotations. Despite the above limitations, the conducted evaluation provides an estimate about the quality of the annotations that can be generated from users that are not experts in privacy related issues. The provided results should not be seen as proof of concept but rather as evidence that the platform and its various aspects can be effective, and that the content (annotations and reviews) created by plain users can be of good quality and value.

### 6.2.1 Process Statistics

After a quick summary of the user-generated annotations we can say that the annotation creation is rather a process that needs time. Users created 139 annotations in total which results to 11.5 annotations per user within a 40 minutes time span or 1 annotation per 3.4 minutes on average (notice that we did not observe a big difference in the total number of annotations given for the easy document and the hard document, i.e. 69 vs 70 annotations). At this point we have to mention that the selected privacy policy documents that the users were asked to annotate were not lengthy ones (less than 450 words). A distribution of the user created annotations per each privacy concern category is depicted in Figure 6.1. It is obvious that the vast majority of the privacy policies mostly contain statements that match the first two categories (i.e. what kind of data are being collected/used, with whom are being shared) whereas there is a limited description on the user choices/control and data security measures (i.e. categories 3,5). Only a fraction of privacy policies contain statements that match the rest categories (i.e. categories 4,6,8,9,10). The same distribution of privacy statements per privacy category is also observed in the average privacy policy [29].

Regarding the review/voting processes, participants submitted 441 votes in total (262 up-votes and 179 down-votes). Each participant received 36.75 votes on average (on average a vote every 1.09 minutes).

### 6.2.2 Relevance Metrics

To measure the quality of the crowd-sourced annotations we deployed standard Information Retrieval (IR) metrics i.e. Precision (equation (6.1)), Recall (equation

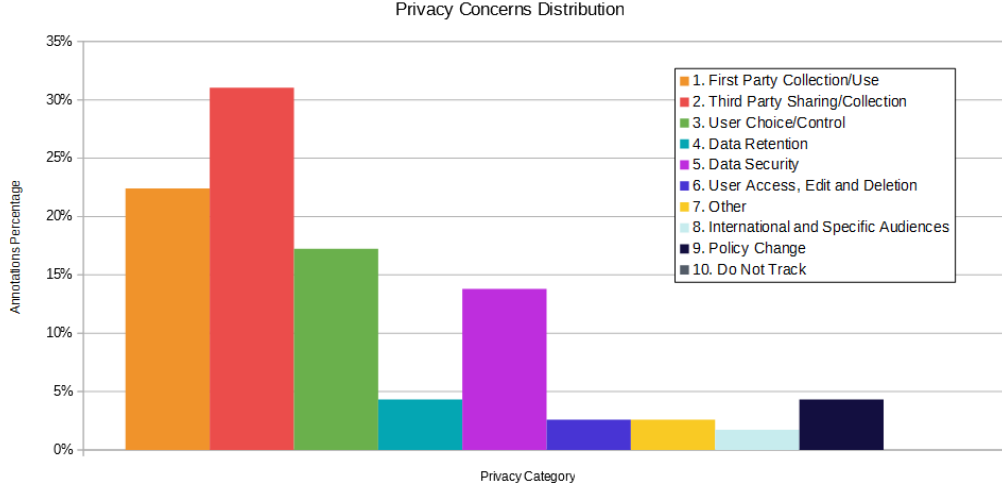


Figure 6.1: A distribution of the generated annotations per privacy concern category, which is representative of the distribution of the expert based annotations in the OPP-115 dataset

(6.2)) and F-measure (equation (6.3)). Those metrics were computed by considering the expert based annotations given in the OPP-115 dataset, which consisted the ground truth in our case.

$$Precision = \frac{|\{UserAnnotations\} \cap \{ExpertAnnotations\}|}{|\{UserAnnotations\}|} \quad (6.1)$$

$$Recall = \frac{|\{UserAnnotations\} \cap \{ExpertAnnotations\}|}{|\{ExpertAnnotations\}|} \quad (6.2)$$

$$F - Measure = \frac{2 * Precision * Recall}{Precision + Recall} \quad (6.3)$$

We computed the above relevance metrics across two different aspects: a) a coarse-grain one regarding the relevance of privacy concern categories and b) a fine-grain one regarding the annotated tags relevance.

- **Relevant by Privacy Concern category (PC):** A user-created annotation is considered relevant when the privacy concern category of the annotation matches the category of any annotation created by experts on the same segment of the privacy policy.
- **Relevant by Tags:** A user-created annotation is considered relevant when the mandatory attribute-values (i.e. tags in the platform's case) of the annotation match the corresponding mandatory attribute-values of an annotation created by experts on the same segment of text. This kind of relevance presupposes relevance by privacy concern (i.e. the coarse-grain relevance).

The interpretation of the second definition of the annotation relevance shows that it is a stricter relevance metric. When the tags between two annotations match each other the privacy concern categories will also match since each tag belongs to a specific privacy concern category. So, the set of annotations that occurs with the relevant by tags metric should be a subset of the set of the annotations that occurs using the relevant by PC metric.

It is a surprise that the difference between these two annotations sets is very small. More specifically, the total relevant annotations by PC were 70 (out of 139 annotations in total) whereas the relevant annotations by tags were 65. This means that when the participants were able to identify the correct privacy concern category, they were able to identify 93% (65/70) the mandatory tags that should be placed.

Since the difference between the two metrics is rather small, for reasons of simplicity we will report results for the relevant by PC metric<sup>1</sup>. By using the PC relevancy metric the provided numbers can be directly compared with the corresponding metrics of the current state-of-the-art algorithms which can automatically detect the privacy concern category per segment [29].

### 6.2.3 User vs Crowd Created Annotations

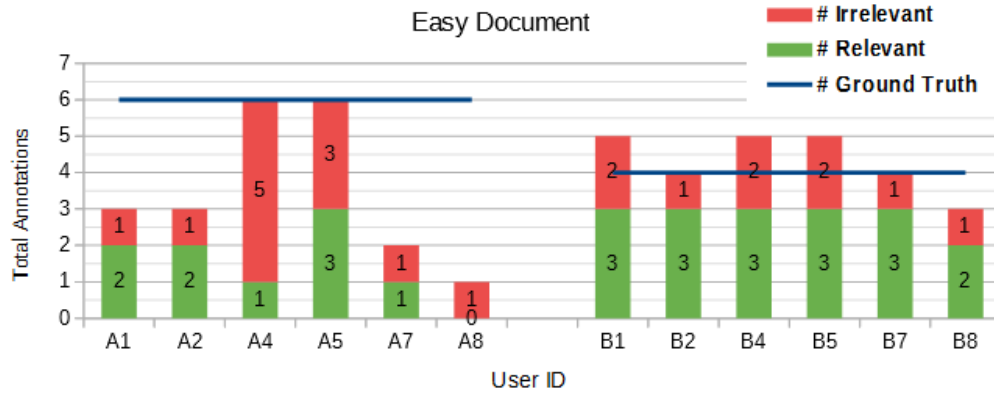
As mentioned earlier, the evaluation process was designed in such a way that we could measure the quality of: a) the created annotations and b) the annotations produced by the reviewing/voting process.

Recall, that annotations were created in the first step of phases B and C respectively. Figure 6.2 provides a detailed description of the annotations generated by the participants. As previously discussed the provided numbers are for the relevance by PC metric, since there are no big difference between relevance by PC and tags relevance. The figure shows the total annotations per participant created for the easy and the hard document, and the number of relevant/irrelevant annotations. The blue line represents the number of the relevant (by privacy concern category) annotations identified by the experts in each task. They denote the actual user performance in the annotation creation tasks.

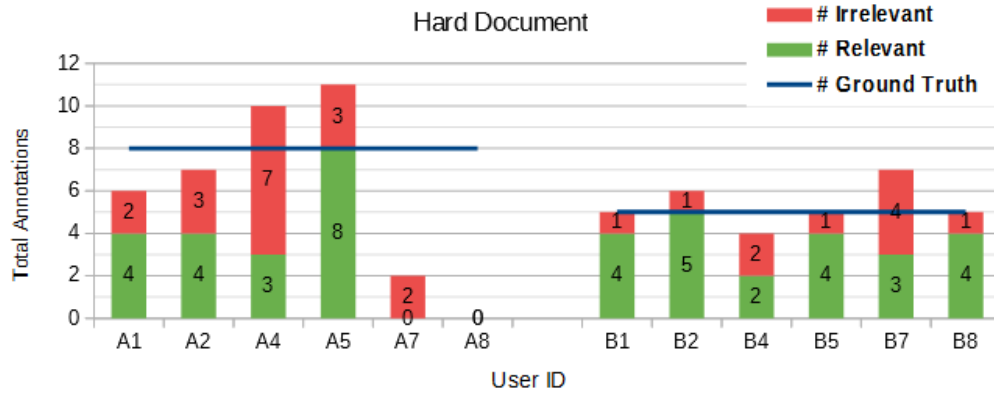
Notice that none of the participants managed to produce the exact annotations that were given by the experts. Some notable exceptions are users A5 and B2, who were able for TaskB to provide all the expert based annotations, plus some irrelevant ones. On the other hand most participants managed to produce a bigger number of relevant annotations than irrelevant ones, except from users A4 and A8. The first one produced consistently irrelevant annotations for both tasks, while the latter one had almost zero output of annotations. The output of groupB is more consistent with less variance and better overall results, while the opposite holds for groupA. Notice that although both groups had to read the same amount of text per privacy policy, the number of expert annotations for groupA was bigger

---

<sup>1</sup>when not stated otherwise



(a) Generated annotations per user on the easy document



(b) Generated annotations per participant on the hard documents

Figure 6.2: Relevant and irrelevant annotations for easy and hard documents per participant

for both tasks (although of lesser complexity), i.e. groupA participants had to produce more annotations per task (6 instead of 4 for taskA and 8 instead of 5 for taskB). But at least for taskA, groupB produced more annotations than groupA (consider though the two outlier participants of groupA we mentioned earlier).

After the first step of creating annotations, the participants were asked to review/vote the annotations created by the participants of the other group. The task was to vote up/down annotations so that the wisdom of the crowd (i.e., aggregation of votes) could pinpoint those annotations that it assumes as correct or more precise and mark as invalid those that are erroneous or imprecise. The validity of each annotation was based on the Wilson Score Interval (as described in Chapter 3) with a 90% confidence level.



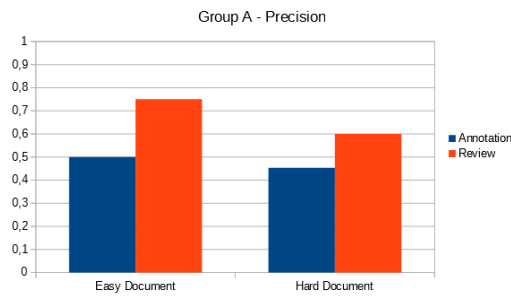
The valid annotations that occurred after the review/voting process are considered to be the community created annotations, since an annotation is marked as valid only when there is a strong user agreement on up-votes. Valid annotations were filtered to keep only distinct annotations since users were working in parallel and the number of duplicates would interfere the results (currently there is no way to mark duplicate annotations).

A comparison between the precision/recall/f1 values per group can be found in Figure 6.3, while Figure 6.4 provides the precision/recall/f1 values in total. The relevance of each annotation is defined as the privacy concern category match between the user and the expert in the document segment it belongs. The blue bars illustrate the average precision/recall/f1 measures for the annotations created by the users (first step of each task), whereas the red bars illustrate the scores of the annotation set that occurred by the community (second step that marks valid annotations after the review/voting phase).

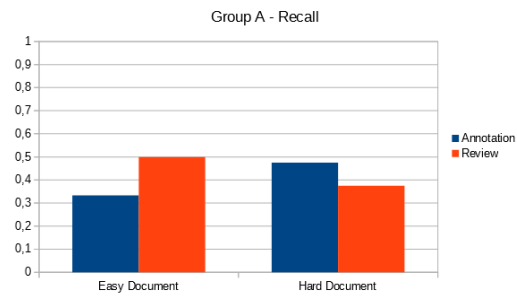
We can observe a large increase in the precision of the annotations after the review step for all groups and documents. This means that the valid annotations that the crowd/community was able to feature was of higher quality and more relevant than the annotations produced by the each participant. On the other hand, while there was also an increase in the recall value on the easy document after the review phase, there was a decrease of the recall value on the hard documents for both groups. One speculation on this effect could be that the review process of the crowd/community is somewhat slower than the task of creating annotations, since participants had to compare all available annotations, some of them with minor variations, that were produced by the other group. Despite the fact of smaller recall levels in our evaluation, we believe that this inefficiency could be addressed if annotations were reviewed by more participants with more available time. Since a lot of users have to agree for making an annotation valid and some annotations might be ambiguous, confusing or even not so precise, minor differences in the approaches of how each participant upvote, downvote, or does not provide any vote at all, affect the final result of the review process. For example some participants did not upvote/downvote an annotation if they found a more refined one (which they upvoted), while other users downvoted coarse-grained annotations and upvoted only the one they thought was the most precise. Notice that we did not propose a specific approach on how to vote/review the annotations, in order to not affect the different approaches for reviewing annotations.

Ambiguous annotations can also help policy makers to detect statements that confuse plain user and may be used as warning flags for updates/clarifications on the text. Notice that many annotations remained with 0 votes when the process was completed (they were neither upvoted or downvoted). It seems possible that most of the participants were not willing to spend effort reviewing similar annotations.

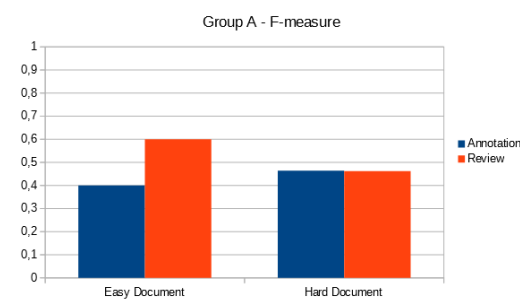
The number of annotations for each upvote/downvote ratio voted by the crowd are given in Figure 6.5 (a). This figure also depicts the annotation resolution



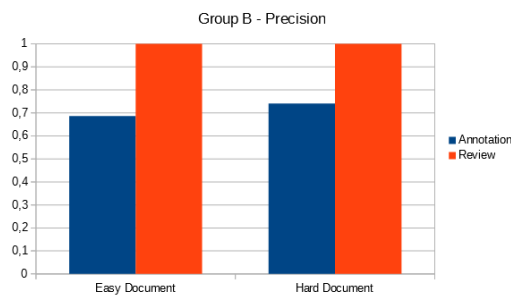
(a) Precision for group A



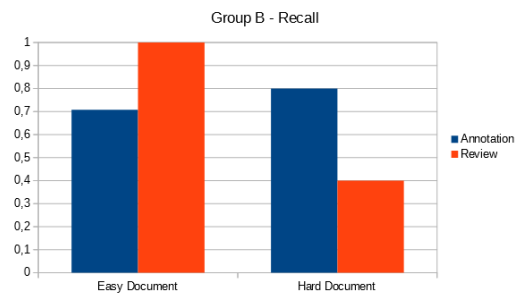
(b) Recall for group A



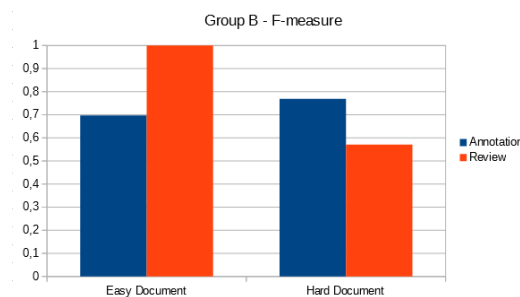
(c) F-Measure for group A



(d) Precision for group B



(e) Recall for group B



(f) F-Measure for group B

Figure 6.3: Precision, recall and f-measure values per group for annotating and reviewing tasks

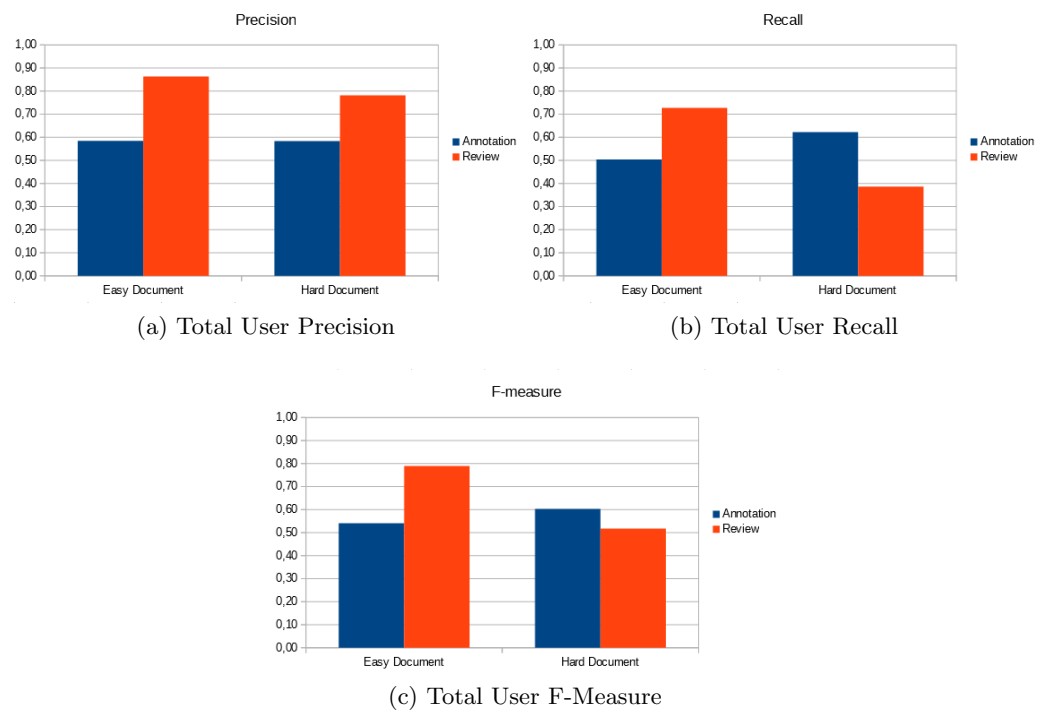


Figure 6.4: Precision, recall and f-measure values in total for annotating and reviewing tasks

within the relevant annotation sets of both relevance metrics. Figure a) demonstrates the total annotations that the users created/reviewed based on the total up-votes/down-votes each annotation received. Figures b and c, illustrates the relevant annotations on both defined metrics based on the received votes. Given that each annotation could receive at most 6 votes (each user of the group could submit one vote per entity), cells that represents more than 6 votes were left empty. The green/red areas denote that an annotation resolves into valid/invalid using the Wilson Score Interval with 90% confidence level. Note that the difference (cell numbers are shown in bold) between the 2 relevance metrics (b and c) appeared on annotations that tend to resolve into invalid since their aggregated score is less or equal to zero. A number of interesting observations can be extracted from this figure.

Firstly, we can output the confusion matrix and measure the error rate on the valid/invalid annotations of the crowd. As Table 6.1 shows, the crowd/community was able to correctly resolve 26 annotations out of 33 which results to an error rate of 21%. A rather interesting fact is that 3 (out of 4) valid but irrelevant annotations according to the expert based annotations were, upvoted by all users. After a careful inspection of the corresponding annotations, we believe that the crowd is correct and that those specific annotations were missed by the experts! On the other hand some relevant annotations were downvoted and were found invalids by the crowd.

	Irrelevant	Relevant
Invalid	9	3
Valid	4	17

Table 6.1: Confusion matrix for valid & invalid annotations

Secondly, based on our results it seems that the difference between the 2 relevance metrics (i.e. privacy category relevance and tag relevance) as shown in Figure 6.5 b) and c) appeared on annotations that tend to resolve into invalid since their aggregated score is less or equal to zero (cell numbers are shown in bold).

Last but not least, there are lots of relevant annotations that confused the crowd/community having similar or equal number of up-votes and down-votes as we can notice from the corresponding cells in the figure 6.5 (b) and (c).

Concluding, based on our evaluation results, plain users are able to produce most of the expert-based annotations. Further, there is an impressive increase in the quality of the accepted as valid annotations by the crowd after the review/voting phase. This quality improvement mainly focus in the precision of the valid annotations, while there is a small drop in the recall, which we believe can be addressed with more participants and more time. Another interesting fact is that the crowd can generate annotations that might be missed by the limited set of expert users.

		Up Votes						
Down Votes		0	1	2	3	4	5	6
	0	2	0	5	8	7	3	3
	1	5	3	8	14	8	2	
	2	3	7	5	6	2		
	3	4	2	2	2			
	4	5	5	2				
	5	3	0					
	6	0						

(a) Total annotations based on the received up-votes/downvotes

		Up Votes						
Down Votes		0	1	2	3	4	5	6
	0	1	0	3	7	7	2	1
	1	4	1	4	10	7	2	
	2	0	3	5	5	1		
	3	0	2	0	0			
	4	0	2	0				
	5	3	0					
	6	0						

(b) Total **relevant by privacy concern** annotations based on the received up-votes/down-votes

		Up Votes						
Down Votes		0	1	2	3	4	5	6
	0	1	0	3	7	7	2	1
	1	4	1	4	10	7	2	
	2	0	2	3	5	1		
	3	0	1	0	0			
	4	0	2	0				
	5	2	0					
	6	0						

(c) Total **relevant by tags** annotations based on the received up-votes/down-votes

Figure 6.5: Upvote/downvote ratio counts for a) total annotations, b) PC relevancy and (c) tag relevancy, as voted by the evaluation participants

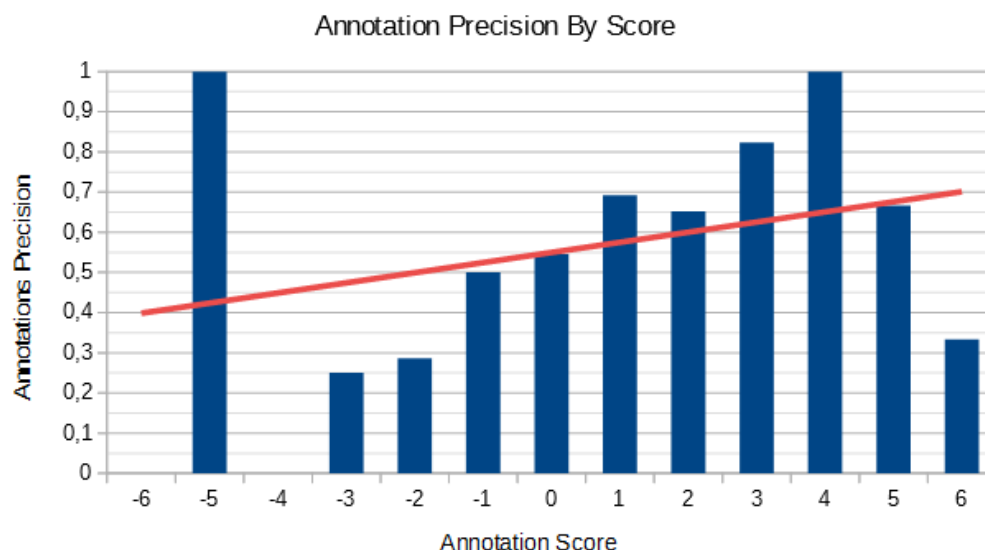


Figure 6.6: Depiction of precision and score of grouped annotations

### 6.3 Annotation Scores

Another important aspect worth measuring is how representative is the annotation score compared to the annotation quality/value. Each annotation maintains a score value based on the aggregation of the up/down votes which is a straightforward way for users to get an indication about the validity and/or popularity of the annotation. Of course, since we are using the Wilson Score Interval, a relatively high difference between the number of upvotes and downvotes does not necessarily result to a valid annotation, since the level of agreement between the users does matter. On the other hand, it is expected that the majority of the annotations with a high aggregated score will be more probable to be valid and accurate whereas annotations with a low/negative score scores will probably be invalid and noisy.

Figure 6.6 provides some clues to this specific question. This figure illustrates the annotations precision (y-axis) in comparison to the annotation score (x-axis), when we grouped annotations by their score. If we except some edge cases, we can see that the annotation precision increases when the annotation score increases with a peak on 100% for the 9 annotations with score equals to 4. While this provides a strong evidence of the annotation quality, as we can see from Figure 6.5 there are 2 annotations (cell 1,5) which maintain a score of 4 but the received negative vote prevents them to resolve into valid. Also there is some noise (e.g., downvoted relevant annotations with score -5). The red line illustrates the best line fit which denotes the trend.

A detailed image of the annotation relevance and the corresponding upvote

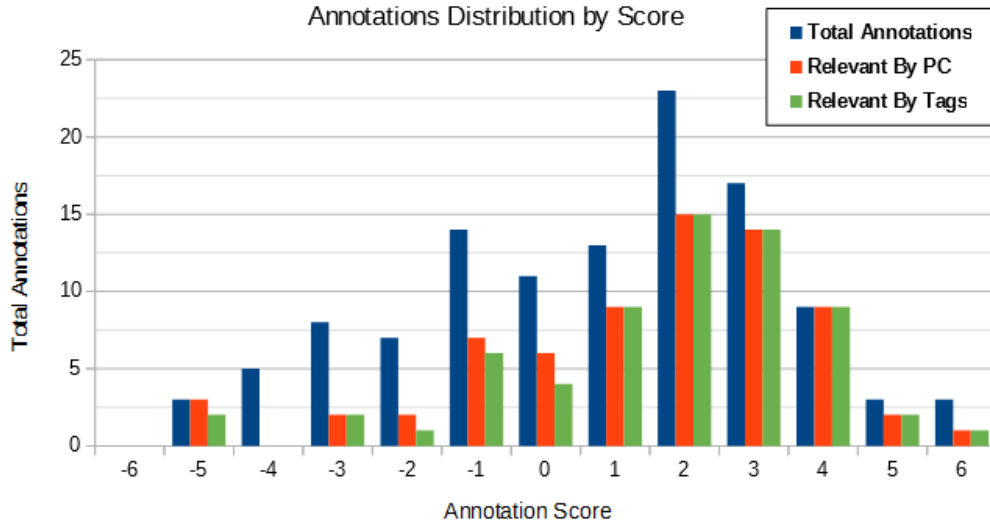


Figure 6.7: Total created annotations and the number of relevant annotations for both relevance metrics grouped by their final score

and downvote difference is given in Figure 6.7. This image depicts the number of annotations grouped by their score difference and the number of the relevant annotations per each relevance metric. Notice that a difference between the red (relevant by PC) and green (relevant by tags) bars occurs only for the annotations with score equal and less to zero. This means that the annotations for which our users were able to correctly detect the privacy concern category but not the correct tags, tend to resolve into invalid since their score is mostly negative (something that was also mentioned earlier). Another interesting fact is that edge cases do not contain enough annotations to make safe assumptions. For this reason, the precision of the annotations grouped by score as illustrated in Figure 6.6 can be considered somewhat noisy.

## 6.4 User Scores

Another aspect of the platform that was evaluated is the reward system deployed by the platform. Each user maintains a score represented by an integer value that results from the aggregation of the reward badges that he/she received. By design, the user score should represent the effort and ability of a user to create or pin-point valuable content (i.e., annotations), as well as the level and the effort of his/her contributions.

In our evaluation process, all participants were new users of the platform, who had to register to the platform and as a result started with a score of 0. Due to this fact, users that were able to create valuable and relevant annotations should

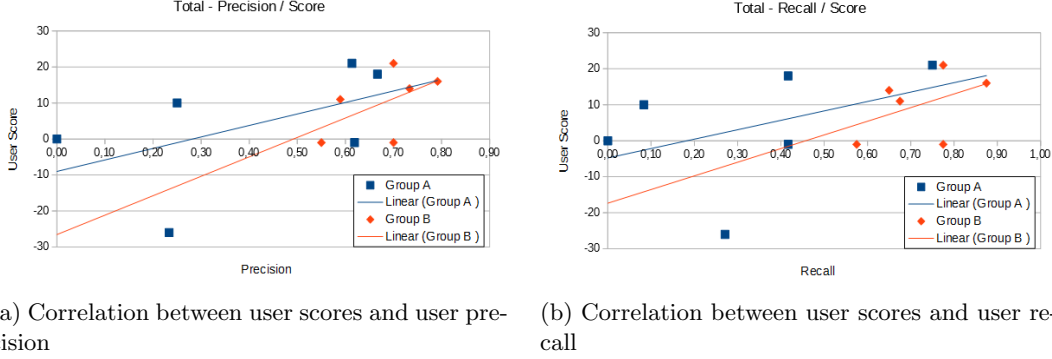


Figure 6.8: Correlation between user scores and the user precision/recall achieved during annotation creation. The lines on the diagrams correspond to line fittings

gather a higher score than the less efficient and effective ones. Figure 6.8 partially confirm this hypothesis. There seems to be some correlation (Pearson corr = 0.511) between the user score and the user precision on the generated annotation. Although less obvious, the same pattern is observed with the user recall level of each user in relation with his score. Higher recall level is found to users with higher scores (Pearson corr = 0.413). Notice the outlier user in group A with the rather low score (-27), which is the user that consistently provided non-relevant annotations.

The not so big correlation coefficients between the users' scores and their precision/recall levels could possibly be ascribed to the small number of users and small number of generated annotations. We expect that with a bigger number of annotations and users we would get a higher correlation of users' scores and recall/precision levels.

Figure 6.9 depicts the final user scores, sorted per group for completeness reasons. As it is obvious the aggregated scores show that participants in group B gathered a higher score than participants of group A (mainly due to the low scored user and the inactive user that gathered non score points).

## 6.5 User Friendliness & General Remarks

The main goal of the evaluation was to measure the effectiveness of our platform in creating and reviewing valuable annotations on privacy policies. Since the platform is based on plain everyday users for this task, user friendliness and user experience while interacting with the platform is of a rather crucial factor. The user friendliness of the platform's interface and the feelings of the users while interacting with the platform can have a major impact in user engagement which affects the usage of the platform and its success and effectiveness. In order to get some insights on the aforementioned issues, we asked users to fill in a questionnaire, in order



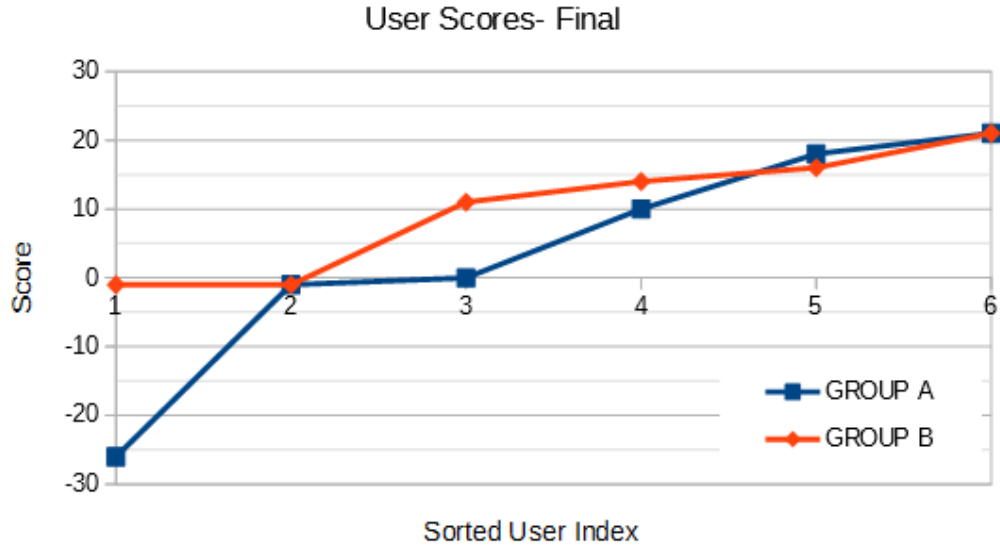


Figure 6.9: The final user scores sorted for each group

to get feedback about their experience with the platform during the evaluation task. The questionnaire basically consists of questions with answers organized in a Likert scale with 5 points. Participants were also asked to write their opinion and comments in free text for some open questions. The given questionnaire is provided in Appendix B.

### 6.5.1 The Platform

The user experience of the platform was one of the main parts that the questionnaire examined, where the participants were asked about the effectiveness and UI friendliness of the platform. In a nutshell, 11/12 participants agreed that the platform's user interface and functionality helps to create and review annotations on privacy policies, while regarding user-friendliness, 9/12 users consider the platform as user friendly, whereas the other 3 users were neutral.

Regarding the deployed annotation schema, participants were asked about the expressiveness and the comprehension of the privacy concerns categories and attribute/values. The results show that 10/12 users consider that the privacy concern categories are expressive enough to describe the privacy policy contents whereas 11/12 users do not disagree that the tags under each category are expressive enough.

Some users suggested that more tags should be added in order to make the schema even more expressive while others mentioned that the number of tags from which they have to choose is too high. The trade-off between the number of tags and their expressiveness is something that should be investigated further for the

annotation schema. Having lots of tags can confuse users and may be a deterrent factor to not use the platform for inexperienced users, while more numerous and expressive tags can be exploited by more experienced and advanced users.

As described in the next chapter about future work, the annotation schema could be divided into a basic and extended set of tags (based on the obligatory attributes) so that newcomers and inexperienced users will have a smoother learning curve. A drawback of this policy is that since it leads to coarse-grain annotations, the newcomers' annotations might get down-voted by expert users. As a result newcomers might get a negative score right at the beginning of their interaction and as a result might disengage from the CAPPa platform. A modification of the scoring formula, so that each participant's annotations are reviewed based on his/her allowed expressiveness for his/her current experience, along with some kind of annotation refinement of newcomers' annotations by experienced ones, could possibly address this problem.

Regarding the questions related to any difficulties in the comprehension of the privacy concerns and their attribute/values, only 1/12 believes that the privacy concern categories are hard to understand, while 2/12 were a bit confused about the supported tags.

### 6.5.2 User Experience & Feelings

In order to get feedback about the user experience and feelings regarding the CAPPa platform, at first we asked participants some questions to see what is their opinion regarding privacy related issues and if they are experienced in reading privacy policies. The gathered responses provide evidence that everyday users are concerned about the privacy of their data. The vast majority though has only partially read some privacy policies. This result showcases the need for tools like the CAPPa platform, that can help everyday users understand the important privacy statements of privacy policies. Another interesting fact is that although most participants did not have any prior experience in reading privacy policies, 10 out of 12 users found the reading of privacy policies relatively easy, with the difficult document being a little harder to understand. This result combined with the quality of the created annotations as previously discussed shows the viability of the proposed crowd-sourced approach.

Participants were also asked to comment on the tasks at hand. Almost all participants (11/12) believe that the time was enough for the requested tasks and 10 out of 12 users were satisfied with their overall performance. There were 4 participants that found the task of creating annotations difficult and frustrating, while only 3 participants had the same opinion about the review task. All participants agree though that the tasks were getting easier, as they were getting familiar with the platform.

Finally, participants found the platform rather useful. Based on their responses, they would consider as an option for getting privacy information and statements about the products/services they use. Further, 75% of the participants

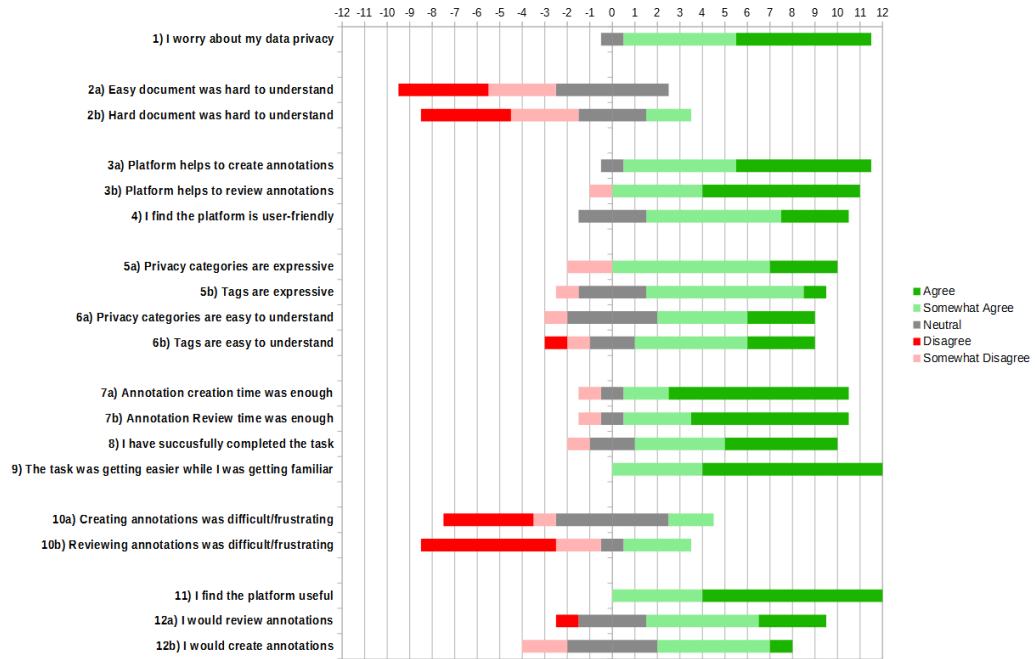


Figure 6.10: Detailed questionnaire responses given by the users

would also contribute to the CAPPA platform by reviewing annotations, if valid annotations are missing for the privacy policy they are interested in, while 66.6% of the participants are willing to read privacy policies of interest to them, and create new annotations (in case of missing annotations). These are rather promising results for the sustainability, impact and effectiveness of the CAPPA platform.



## Chapter 7

# Conclusion & Future Plans

As already mentioned, the aim of this thesis was to design and implement a basic core for a platform to enable users to initiate an active community revolving around awareness on data privacy and the annotation of privacy policies. Some aspects of the platform were not considered as part of the addressed problem though<sup>1</sup> In this chapter we describe possible next steps for improving and extending the CAPPA platform, along with directions for future work and some final thoughts.

### 7.1 Performance Evaluation & Scalability

A dimension of the platform that was not extensively examined is the performance of the back-end and the optimization of the request/response sequence and times. Although, during the evaluation phase we did not observed any significant performance issues, some HTTP responses were not optimized in this thesis and may require improvements (i.e. pagination of the results page). This will help to smoothly scale up, as the platform gains more popularity, active users and content (i.e. annotations).

### 7.2 UI Improvements & Extensions

The evaluation of the platform appeared to be a strong test for the platform's success. Although the users were satisfied from the platform's usability level, some users' comments showcased needed improvements. The most important feedback for the platform was that the set of the available tags for creating annotations on some privacy concern categories were too many. Despite the fact that some users asked for more tags to be supported (those that were able to provide rather refined annotations), it is admitted that for some inexperienced users the quantity of the displayed tag values can seem a bit confusing, keep them away from using the platform, either for the retrieval or the creation of content. One solution for this

---

<sup>1</sup>Some of the discussed missing features are currently under development

issue could be to adapt the number of tag values shown to the user based on the user experience and his/her user score, excluding for example non required tag values. Another option is to let him/her choose the tag level and expressiveness of tag values he/she finds most useful. Of course, a hybrid solution by deploying both suggestions (i.e. adapt to user experience by default and also give the option to the user) is also possible. In both solutions, tags should be grouped by their importance level so that the user can choose between a basic or a more extended set of tags. A drawback of the above approach is that it leads to coarse-grain annotations. As a result, the annotations of the newcomers might get down-voted by expert users, and newcomers might get a negative score right at the beginning of their interaction, disengaging them from the CAPPA platform. A modification of the scoring formula, so that each participant's annotations are reviewed based on his/her allowed expressiveness according to his/her current experience, along with some kind of annotation refinement of non-experienced users' annotations by experienced ones, could possibly solve this problem.

### 7.3 User Engagement Improvements

The current implementation of the user engagement mechanism provides a very basic level. There are a bunch of ideas and extensions that can result to improvements over the user engagement on the platform. One idea is to design more customized badges that focus on specific tasks that will be important for the platforms entities. Specifically, since the annotation completeness of documents is we can create badges that will reward users when they contribute on specific privacy policy in order to complete them i.e. when he/she creates annotations from various privacy concern categories for the same document, or when reaches to a specific number of valid annotations. In addition to these badges, we can use notifications to request users specific actions that they can complete in order to receive some extra reward badges. This will result to crowd users to be more targeted/focused on specific tasks that might be more important than others.

### 7.4 More Components & Functionality

Although the evaluation of the platform's functionality was positive, there is always space for extensions and improvements. This work mainly focused on providing a solid core for creating and voting crowd-source annotations. But more tools can be developed to support the aforementioned task and to further extend the platform for more features. In this section, we elaborate on some ideas that were mentioned during the design and implementation chapters.

### 7.4.1 Like/Dislike Feature

The 'like' button has gained its place in social networks and it is rather important feature of available collaborative platforms. It's been a standard way to get a popularity measure or an acceptance level for any content on the web and along its extensions (e.g. dislike or other emotive icons) can offer valuable feedback. We can apply this feature on some entities of the platform to get information about the crowd acceptance of the privacy policy, the available source and any security issues. Some work has been done towards this direction (user interface and database table design) but it was not completed during this thesis. A mock-up example of this feature is depicted in figure 7.1 (2).

### 7.4.2 Product External Link Resources

Another step for the completeness of the privacy concerns for the user is to support a way for user to post/publish some external links (i.e. news articles, forum and blog posts) as resources for some products that will keep up-to-date the users for privacy related stuff regarding some product. These external documents may report privacy leaks, data thefts or on the other way that happen from time to time. A collection of these links (or the absence of them) could have a major impact for each product usage since it can also hold as a proof of concept for the privacy statements they report. Some work has been done towards this direction (partial user interface and database table design) but it was not completed during this thesis. A mock-up example of this feature is depicted in figure 7.1 (4).

### 7.4.3 Mobile Applications - Android Permissions

Mobile applications are a major issue for data privacy. The vast majority of mobile applications collect huge amounts of data, with most of the users not being aware of it. Specifically, mobile applications ask for permissions that are not necessary in order to offer the service. Since most users are not willing to stop using an application in order to gain their privacy back, it could be a valuable feature of the platform if the crowd could review the permissions asked by the application for the functionality offered. They could further compare them with the privacy statements in the corresponding privacy policies and check their agreement and validity. Some work has been also done towards this direction (partial user interface and database table design) but it was not completed during this thesis. A mock-up example of this feature is depicted in figure 7.1 (3).

### 7.4.4 Argument Web

An interesting addition to the platform would be the deployment of argumentation tools and approaches developed for the Argument Web. The ecosystem of argument web has been expanded with a large number of inter-operable and cross compatible tools for the analysis, navigation and evaluation of arguments across

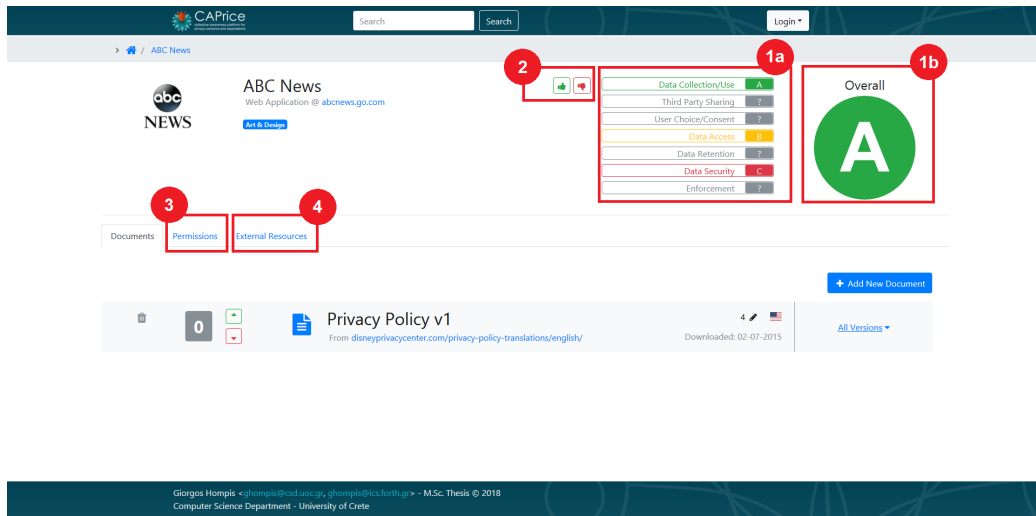


Figure 7.1: User Interface example design of various planned features

a broad range of domains. By definition, each controversial entity of our platform can be considered as an argument raised by a user. This view enables each user to create a new thread for discussion to express it's opinion about how much valid, relevant or complete a controversial entity is. There is a high chance that parts of text in privacy policies can have ambiguous meanings and can result in disagreements between users. Moreover, various annotations can have text overlaps, some can be more refined than others, their tags are superset/subset of others, or they may be duplicates of others. The application of the argument web can equip the platform with a new set of algorithms and tools to resolve all of the above issues which. Some work has been also done towards this direction (database table design) but it was not completed during this thesis.

## 7.5 Other Tools

### 7.5.1 Product Privacy Evaluation

In an abstract perspective, different levels of data privacy are offered by various products. Users do not have a direct way to realize these privacy level in a straightforward manner. A concise and summarized indication of the privacy friendliness level of the product for each privacy aspect (i.e. privacy concern category) is estimated to be a very useful feature for the platform. Users can rely on labels to get a sign on how safe regarding privacy each product is. In case of bad signs, users may be triggered to look for details in privacy statements which implicitly lead them to read and review the documents. The generation of privacy friendliness labels may exploit the valid annotations that each product contains in its privacy policy documentation. Another approach is that users could explicitly label each



product with sign labels and then offer the final score as an aggregation of the available user feedback. Both cases have the potential to work rather well. Some primary UI design effort has been done on this thesis for this feature but we mark it as future work since it's completeness level is too low. A mock-up example of this feature is depicted in figure 7.1 (1).

### 7.5.2 Text Difference Between Documents

It is known, that privacy policies are dynamic and can change, sometimes without any notification. Especially for major applications this can happen very often. It's on the organization's competence (some times described in the privacy policy also) if the user should be notified somehow or not when a change happens. For some users, it would be rather useful to know that a privacy related change has occurred and have a way to review the changes (i.e. diff) of the text in an easy way. Since the platform already tracks the uploaded policy documents, it would really helpful to provide a tool for highlighting text differences between document versions, helping them to update current annotations or provide new ones. The implementation of such a tool could also be helpful for the implementation of an annotation transfer process. Such a process could help transfer annotations from one version of privacy policy into another, for example when the parts of text that the annotation refers to has remained untouched. Currently no much effort has been committed towards this feature.

### 7.5.3 Guided/Assisted Annotation Creation

#### 7.5.3.1 Deployment of NLP & Machine Learning Algorithms

As discussed in previous chapters, another primary direction of the research community for analyzing privacy policy documents is to deploy NLP and ML techniques. Despite the fact that currently these techniques do not perform well on annotating privacy policies, it could be a helpful tool to assist users creating annotations by suggesting text parts that might match with specific privacy statement & concerns. The current state of the art algorithms can perform well on finding the privacy concern category in a coarse level. This means that the algorithms are not able to precisely annotate the corresponding part of text (i.e. they annotate only segments), reason on this annotation and provide refined annotations. In the future we plan to exploit the valid annotations produced by the crowd in CAPPA in order to train ML algorithms on a bigger dataset than the currently used OPP-115 collection.

#### 7.5.3.2 Ambiguous Annotation Flags

Beyond the guided annotation generation and intelligent assistance implementations much simpler guidance would also be rather useful. For example the user could be assisted by the use of marks, icons and messages that would notify him

for possible duplication or ambiguity with an existing annotation before creating a new (duplicate) one. The annotation duplication or ambiguity can be directly checked by comparing the tag sets (i.e. subsets and supersets) of the annotation, and/or the highlighted text strings. As a result the chances of noisy and duplicated annotations could be eliminated, and the user voting effort would not spread over duplicate annotations. In this way we could maintain clean versions of the privacy policies and their annotations, and keep the **CAPPA** system fast and responsive.

## 7.6 Discussion

During the last months, the importance of data privacy has attracted considerable attention for the social masses due to the application deadline of the GDPR (General Data Privacy Regulation) and the recent major data leak of user's data (i.e. facebook's scandal [56]). A movement on data privacy has been initiated with a portion of users being more aware about.

In this thesis, we developed the **CAPPA** system, an online open crowdsourcing platform, that allows the insertion and tracking of privacy policies and supports the annotation of privacy policies with privacy statements about various privacy concerns. Users are able to review annotations through an upvote/downvote process in order to pinpoint qualitative annotations and remove noisy and erroneous ones.

We conducted a user based evaluation, where the crowd-sourced annotations were compared with expert based ones. One important result is that indeed plain users are able to provide annotations of high quality. Although not all users are able to consistently add high quality content, participants of the platform are able to highlight qualitative annotations through the upvoting/downvoting process. Further, the crowd strongly upvoted some annotations that were not given by the experts, which indeed seem rather relevant and somehow skipped the eyes of the experts. Users that participated in the evaluation process reported that reading privacy policies is not as hard as they initially thought and that they become more tolerant on reading privacy policies. Another crucial result of the evaluation is that most users would use the platform to get feedback for applications and services that may use or be interested in the future, and would consider reading privacy policies and reviewing annotations in case the corresponding privacy policies were not annotated.

# Appendices



## Appendix A

### Tags: Concerns, Attributes & Values

The complete list of currently available privacy concerns, tag attributes and tag values as defined in [usableprivacy.org](http://usableprivacy.org)

## A.1 Privacy Concerns

Privacy Concern	Description
First Party Collection/Use	Privacy practice describing data collection or data use by the company/organization owning the website or mobile app.
Third Party Sharing/Collection	Privacy practice describing data sharing with third parties or data collection by third parties. A third party is a company/organization other than the first party company/organization that owns the website or mobile app.
User Choice/Control	Practice that describes general choices and control options available to users.
User Access, Edit and Deletion	Privacy practice that allows users to access, edit or delete the data that the company/organization has about them.
Data Retention	Privacy practice specifying the retention period for collected user information.
Data Security	Practice that describes how users' information is secured and protected, e.g., from confidentiality, integrity, or availability breaches. Common practices include the encryption of stored data and online communications.
Policy Change	The company/organization's practices concerning if and how users will be informed of changes to its privacy policy, including any choices offered to users.
Do Not Track	Practices that explain if and how Do Not Track signals (DNT) for online tracking and advertising are honored.
International and Specific Audiences	Specific audiences mentioned in the company/organization's privacy policy, such as children or international users, for which the company/organization may provide special provisions.
Other	Another aspect not covered in the other categories is discussed in the text segment.

Table A.1: Privacy Concerns

## A.2 Tag Attributes

First Party Collection/Use			
Tag Attribute	Mandatory	Multivalue	Description
Does/Does Not	false	false	Use this optional attribute to denote if the policy explicitly states that something is NOT done. Defaults to <i>Does</i> .
Collection Mode	false	false	Use this optional attribute to denote if the data collection performed by the first party is implicit (e.g., company collects information without user's explicit awareness) or explicit (e.g., user provides information). Defaults to <i>Not selected</i> .
Action First-Party	true	false	How does the first party collect, track, or obtain user information?
Identifiability	false	false	Use this optional attribute if it is explicitly stated whether the information or data practice is linked to the user's identity or if it is anonymous. Defaults to <i>Not selected</i> .
Personal Information Type	true	true	What category of information is collected or tracked by the company/organization?
Purpose	true	false	What is the purpose of collecting or using user information?
User Type	false	false	Use this optional attribute if this practice applies specifically to users with an account or users without an account.
Choice Type	false	false	Use this optional attribute if user choices are explicitly offered for this practice. Defaults to <i>Not selected</i> .
Choice Scope	false	false	Use this optional attribute to indicate the scope of user choices. In some cases, even if user choices are not clear or specific, this attribute can be selected. Defaults to <i>Not selected</i> .

Table A.2: Attributes for *First Party Collection/Use* Privacy Concern

Third Party Sharing/Collection			
Tag Attribute	Mandatory	Multivalue	Description
Third Party Entity	true	false	The third-party involved in the data practice.
Does/Does Not	false	false	Use this optional attribute to denote if the policy explicitly states that something is NOT done. Defaults to <i>Does</i> .
Action Third Party	true	false	How does the third-party receive, collect, track, or see user information.
Identifiability	false	false	Use this optional attribute if it is explicitly stated whether the information or data practice is linked to the user's identity or if it is anonymous. Defaults to <i>Not selected</i> .
Personal Information Type	true	true	What category of information is shared with, collected by or otherwise obtained by the third-party.
Purpose	true	false	What is the purpose of a third party receiving or collecting user information?
User Type	false	false	Use this optional attribute if this practice applies specifically to users with an account or users without an account.
Choice Type	false	false	Use this optional attribute if user choices are explicitly offered for this practice. Defaults to <i>Not selected</i> .
Choice Scope	false	false	Use this optional attribute to indicate the scope of user choices. In some cases, even if user choices are not clear or specific, this attribute can be selected. Defaults to <i>not selected</i> .

Table A.3: Attributes for *Third Party Sharing/Collection* Privacy Concern



User Choice/Control			
Tag Attribute	Mandatory	Multivalue	Description
Choice Type	true	false	The type of user choice or privacy control options available to users.
Choice Scope	true	false	What scope does the user choice or control apply to, i.e., first party collection/use or third party collection/use. Note that sometimes use of information can be limited, but the information is still collected from users.
Personal Information Type	true	true	What category of information does the user choice apply to?
Purpose	true	false	What purpose/use of information does the user choice apply to?
User Type	false	false	Use this optional attribute if this practice applies specifically to users with or without an account.

Table A.4: Attributes for *User Choice/Control* Privacy Concern

User Access, Edit and Deletion			
Tag Attribute	Mandatory	Multivalue	Description
Access Type	true	false	Options offered for users to access, edit, delete information that the company/organization has about them.
Access Scope	true	false	If access is offered, what data does it apply to.
User Type	false	false	Use this optional attribute if this practice applies specifically to users with or without an account.

Table A.5: Attributes for *User Access, Edit and Deletion* Privacy Concern

Data Retention			
Tag Attribute	Mandatory	Multivalue	Description
Retention Period	true	false	Description of the retention period, i.e., how long data is stored.
Retention Purpose	true	false	The purpose to which the retention practice applies (may be <i>unspecified</i> ).
Personal Information Type	true	true	The information type for which the retention period is specified (may be <i>unspecified</i> ).

Table A.6: Attributes for *Data Retention* Privacy Concern

Data Security			
Tag Attribute	Mandatory	Multivalue	Description
Security Measure	true	false	Policy statements that describe the type of security that the website/app implements to protect users' information.

Table A.7: Attributes for *Data Security* Privacy Concern

Policy Change			
Tag Attribute	Mandatory	Multivalue	Description
Change Type	true	false	For what type of changes to the website/app's policy are users notified.
Notification Type	true	false	How is the user notified when the privacy policy changes.
User Choice	true	false	What choices/options are offered to the user when the policy changes.

Table A.8: Attributes for *Policy Change* Privacy Concern

Do Not Track			
Tag Attribute	Mandatory	Multivalue	Description
Do Not Track policy	true	false	If and how Do-Not-Track signals (DNT) are honored.

Table A.9: Attributes for *Do Not Track* Privacy Concern

International and Specific Audiences			
Tag Attribute	Mandatory	Multivalue	Description
Audience Type	true	false	Select which audience the policy segment refers to

Table A.10: Attributes for *International And Specific Audiences* Privacy Concern

Other			
Tag Attribute	Mandatory	Multivalue	Description
Other Type	true	false	What other aspect not covered in the other categories is discussed in the text segment?

Table A.11: Attributes for *Other* Privacy Concern

### A.3 Tag Attribute Values

First Party Collection/Use		
Tag Attribute	Value	Description
Does/Does Not	Does	The first party does engage in the described practice.
Does/Does Not	Does Not	The first party does not engage in the described practice.
Collection Mode	Explicit	The company/organization collects or uses information that the user explicitly provides, e.g., the user enters information in a web form.
Collection Mode	Implicit	The company/organization collects or uses information that the user does not explicitly provide, e.g., data is collected or transferred automatically in the background. The user may or may not have given consent to such implicit collection/use.
Collection Mode	Unspecified	It is not specified or unclear whether the information is collected explicitly or implicitly.
Action First-Party	Collect on website	The company/organization collects user information directly on the website.

Action First-Party	Collect in mobile app	The company/organization has a mobile app and it collects user information through that platform.
Action First-Party	Collect on mobile website	The company/organization has a mobile version of its website through which it collects user information. This value is only needed if the policy explicitly distinguishes between its normal and mobile websites.
Action First-Party	Track user on other websites	The company/organization tracks its users' activities when they visit other websites, typically without the user being aware of it.
Action First-Party	Collect from user on other websites	This company/organization (the first party) has a widget or element on other websites, in which the user can explicitly provide data to the first party while being on the other website. For example, Facebook users can comment on news or other content on different websites and Facebook collects that data.
Action First-Party	Receive from other parts of company/affiliates	The company/organization is part of a family of companies/organizations or has subsidiaries (e.g., physical store, other websites that belong to same company). It receives user information from those other units.
Action First-Party	Receive from other service/third-party (unnamed)	The company/organization receives user information from an unnamed third-party (e.g., the policy just speaks of "data brokers" or "partners" in the abstract).

Action First-Party	Receive from other service/third-party (named)	The company/organization acquires user information from a third party that is explicitly named (e.g., Facebook when the user signs in using Facebook account; or a specific "partner").
Action First-Party	Other	The specified type of collection is not covered by the options above.
Action First-Party	Unspecified	The type of collection is not specified or unclear, e.g., "We collect your personal information" without further specification whether the collection occurs on a website, an app, offline, etc.
Identifiability	Identifiable	It is explicitly stated that the information/data practice is linked to the user's identity.
Identifiability	Aggregated or anonymized	The collected data is anonymized (e.g., link to user's identity is removed) or aggregated (e.g., merged with other users' information so that it is not possible to uniquely identify a single user).
Identifiability	Other	The practice makes an explicit statement about identifiability that is not covered by the options above.
Identifiability	Unspecified	It's not explicitly stated or unclear if the information/data practice is linked to the user's identity.
Personal Information Type	Financial	Financial information, such as credit/debit card data, other payment information, credit scores, etc.
Personal Information Type	Health	Health Information, such as information about health conditions, prescriptions, medication, as well as health monitoring data, e.g., heart rate, step count, activity level, etc.

Personal Information Type	Contact	Contact Information, such as name, email address, phone number, street address, etc.
Personal Information Type	Location	Geo-location information (e.g., user's current location) regardless of granularity, i.e., could be exact location, ZIP code, city-level.
Personal Information Type	Demographic	Demographic Information, e.g., gender, age, occupation, education, etc.
Personal Information Type	Personal identifier	Identifiers that uniquely identify a person, e.g., SSN, driver's license number, etc.
Personal Information Type	User online activities	The user's online activities on the first party website/app or other websites/apps, e.g., pages visited, time spent on pages, general user behavior online, etc.
Personal Information Type	User profile	The user's profile on the first-party website/app and its contents, e.g., data in user profile, data that user uploaded to website, user comments, user profile preferences, etc. This is common for websites/apps where users can create an account or profile, e.g., on twitter, youtube, Facebook, Amazon, etc.
Personal Information Type	Social media data	User profile and data from a social media website/app or other third party service to which the user gave the first party access, e.g., by connecting with Facebook, twitter, or other services. Exchanged data may include user profile, photos, comments, friends, etc.

Personal Information Type	IP address and device IDs	Permanent (e.g., device IDs, MAC address) or temporary (e.g., IP address) identifiers needed to establish a connection for the current browsing session.
Personal Information Type	Cookies and tracking elements	Identifiers locally stored on user's device by the company/organization or third-parties including cookies, beacons, or similar that are commonly used to uniquely identify users, but that are not essential to establish a connection with the user's device or to provide a service.
Personal Information Type	Computer information	The type of operating system (OS) or web browser that the user uses, or similar computer or device information.
Personal Information Type	Survey data	Any data that is collected through surveys
Personal Information Type	Generic personal information	No specific type of information is mentioned, but the policy talks about "personal information" or "personal identifiable information" in general.
Personal Information Type	Other	A specific type of information not covered by the above categories.
Personal Information Type	Unspecified	The type of information is not explicitly stated or unclear (e.g., refers to "information" very generically).
Purpose	Basic service/feature	Provide a service that the user explicitly requests and that is part of the website/app's basic service or functionality. Examples are watching a video, reading an article, making a purchase, creating an account, contacting the company, etc.

Purpose	Additional service/feature	Provide a service that the user explicitly requests but that is not a necessary part of the website/app's basic service. Additional services/features may enhance user experience or add convenience but require additional data, e.g., social media integration, comments, blog participation, a store finder that needs location information, etc.
Purpose	Advertising	To show ads that are either targeted to the specific user or not targeted.
Purpose	Marketing	To contact the user to offer products, services, or other promotions (e.g., send marketing emails, calling or texting user with marketing messages). Marketing typically requires the use of contact information.
Purpose	Analytics/Research	For understanding the website/app's audience, improving the website/app, inform company strategy, or general research.
Purpose	Personalization & Customization	For providing user with a personalized experience, e.g., by allowing to arrange how the website/app looks, based on the user's preferences or language, etc.
Purpose	Service Operation and Security	For website/app operation and security, enforcement of terms of service, fraud prevention, protecting users and property, etc.
Purpose	Legal requirement	For compliance with legal obligations, e.g., regulations, government data requests, government retention requests, law enforcement requests in general, etc.



Purpose	Merger/Acquisition	If company/organization merges or is acquired it transfers users' information to another company/organization.
Purpose	Other	Other specific purpose not covered above.
Purpose	Unspecified	The purpose is not explicitly stated or is unclear.
User Type	User without account	This data practice specifically applies to users that do not have an account or are not registered with the website or mobile app.
User Type	User with account	This data practice specifically applies to users with an account or who are registered with the website or mobile app.
User Type	Other	This data practice applies to a specific user type not covered by the options above.
User Type	Unspecified	It is not specified whether this practice applies to users with or without account.
Choice Type	Don't use service/feature	Only option is not to use the feature or service. Only select this if explicitly stated in policy (i.e., don't interpret silence as "Don't use website or feature").
Choice Type	Opt-in	User must consent before data can be collected or used by first party.
Choice Type	Opt-out link	Link provided in privacy policy, on website, in mobile app, or in email, etc.
Choice Type	Opt-out via contacting company	Must contact company via email, phone, or postal mail to opt-out.
Choice Type	First-party privacy controls	Website/app provides user settings for privacy configuration.
Choice Type	Third-party privacy controls	Choices provided by a third party (e.g., privacy settings on social media site) or industry (e.g., AdChoices Opt-out).

Choice Type	Browser/device privacy controls	Policy suggests the use of browser or mobile device's privacy settings, e.g., to block trackers or cookies, activate Do-Not-Track, disable location sharing, clear history, etc.
Choice Type	Other	Other specific user choice or control option not captured above.
Choice Type	Unspecified	No user choices mentioned for this practice.
Choice Scope	Collection	Choices apply to collection only.
Choice Scope	Use	Choices apply to use only.
Choice Scope	Both	Choices apply to both collection and use.
Choice Scope	Unspecified	No specific scope of choices is mentioned.

Table A.12: Tag Attribute Values for First Party Collection/Use Privacy Concern

Third Party Sharing/Collection		
Tag Attribute	Value	Description
Third Party Entity	Unnamed third party	The third party is not explicitly named, i.e., it is just generically referred to as "third-party," "partner," or similar.
Third Party Entity	Named third party	The third party is explicitly named (e.g. Facebook) or at least characterized (e.g., "advertising partner" or "data broker").
Third Party Entity	Other part of company/affiliate	Data is made available to other parts of the company/organization, e.g., it is shared with other services, apps, or websites operated by the company, could also be data exchange between online and offline company units (e.g. physical stores).
Third Party Entity	Other users	The third-party involved are other users of the first party website or mobile app.

Third Party Entity	Public	User information is made public or can be obtained from public sources.
Third Party Entity	Other	Other specific third-party entity not covered above.
Third Party Entity	Unspecified	The third-party entity is not specified. This is uncommon, "unnamed third-party" will likely be the right value instead.
Does/Does Not	Does	The third party does engage in the described practice.
Does/Does Not	Does Not	The third party does not engage in the described practice.
Action Third Party	Receive/Shared with	The third party receives information from the first party. (i.e., the first party explicitly shares data with third-party)
Action Third Party	Collect on first party website/app	The third party explicitly collects data from users on the first party website/app, e.g., by functionality on the website/app that allows users to directly provide information to the third party, such as social media sharing buttons or commenting forms.
Action Third Party	Track on first party website/app	The third party implicitly collects data about users directly on the first party website/app, typically without the user being aware of it, e.g., by tracking users with cookies, beacons, third party ad libraries, or other functionality.

Action Third Party	See	Third-party can see user information that is publicly available either on the website/app or somewhere else. Remember, a third-party can be another user of the website/app. The difference to "receive" is that the information is available but not explicitly given to a specific third party.
Action Third Party	Other	How the third party collects or receives user information is specified in the policy but it is not covered by the options above.
Action Third Party	Unspecified	The type of collection is not specified or unclear, e.g., "Our outside partners collect your personal information" without further specification whether the collection occurs on a website, app, offline, etc.
Identifiability	Identifiable	It is explicitly stated that the information/data practice is linked to the user's identity.
Identifiability	Aggregated or anonymized	The collected data is anonymized (e.g., link to user's identity is removed) or aggregated (e.g., merged with other users' information so that it is not possible to uniquely identify a single user).
Identifiability	Other	The policy makes an explicit statement about identifiability that is not covered by the options above.
Identifiability	Unspecified	It is not explicitly stated or unclear if the information/data practice is linked to the user's identity.

Personal Information Type	Financial	Financial information, e.g., credit/debit card data, other payment information, credit scores, etc.
Personal Information Type	Health	Health Information, such as information about health conditions, prescriptions, medication, as well as health monitoring data, e.g., heart rate, step count, activity level, etc.
Personal Information Type	Contact	Contact Information, e.g., name, email address, phone number, street address, etc.
Personal Information Type	Location	Geo-location information (e.g., user's current location) regardless of granularity, i.e., could be exact location, ZIP code, city-level.
Personal Information Type	Demographic	Demographic Information, e.g., gender, age, occupation, education, etc.
Personal Information Type	Personal identifier	Identifiers that uniquely identify a person, e.g., SSN, driver's license number, etc.
Personal Information Type	User online activities	The user's online activities on the first party website/app or other websites/apps, e.g., pages visited, time spent on pages, general user behavior online, etc.
Personal Information Type	User Profile	The user's profile on the first-party website/app and its contents, e.g., data in user profile, data that user uploaded to website/app, user comments, user profile preferences, etc. This is common for websites/apps where users can create an account or profile, e.g., on twitter, youtube, Facebook, Amazon, etc.

Personal Information Type	IP address and device IDs	Permanent (e.g., device IDs, MAC address) or temporary (e.g., IP address) identifiers of the user device (e.g., computer, mobile device, etc.) needed to establish a connection for the current browsing session.
Personal Information Type	Cookies and tracking elements	Identifiers locally stored on user's device by company/organization or third-parties including cookies, beacons, or similar that are commonly used to uniquely identify users, but that are not essential to establish a connection with the user's device or to provide a service.
Personal Information Type	Computer information	The type of operating system (OS) or web browser that the user uses, or similar computer or device information.
Personal Information Type	Survey data	Any data that is collected through surveys.
Personal Information Type	Generic personal information	No specific type of information is mentioned, but the policy talks about "personal information" or "personal identifiable information" in general.
Personal Information Type	Other	A specific type of information not covered by the above categories.
Personal Information Type	Unspecified	The type of information is not explicitly stated or unclear (e.g., refers to "information" very generically).
Purpose	Basic service/feature	Provide a service that user explicitly requests and that is part of website/app's basic service or functionality. Examples are watching a video, reading an article, making a purchase, creating an account, contacting the company, etc.

Purpose	Additional service/feature	Provide a service that the user explicitly requests but that is not a necessary part of the website/app's basic service. Additional services/features may enhance user experience or add convenience but require additional data or sharing with third parties, e.g., social media integration, comments, blog participation, a store finder that needs location information, etc.
Purpose	Advertising	To show ads that are either targeted to the specific user or not targeted.
Purpose	Marketing	To contact the user to offer products, services, or other promotions (e.g., send marketing emails, calling or texting user with marketing messages). Marketing typically requires the use of contact information.
Purpose	Analytics / Research	For understanding the website/app's audience, improving the website/app, inform company strategy, or general research.
Purpose	Personalization & Customization	For providing user with a personalized experience, e.g., by allowing to arrange how the website/app looks, based on the user's preferences or language, etc.
Purpose	Service operation and security	For website/app operation and security, enforcement of terms of service, fraud prevention, protecting users and property, etc.
Purpose	Legal requirement	For compliance with legal obligations, e.g., regulations, government data requests, government retention requests, law enforcement requests in general, etc.

Purpose	Merger/Acquisition	If company/organization merges or is acquired it transfers users' information to another company/organization.
Purpose	Other	Other specific purpose not covered above.
Purpose	Unspecified	The purpose is not explicitly stated or is unclear.
User Type	User without account	This data practice specifically applies to users that do not have an account or are not registered with the website or mobile app.
User Type	User with account	This data practice specifically applies to users with an account or who are registered with the website or mobile app.
User Type	Other	This data practice applies to a specific user type not covered by the options above.
User Type	Unspecified	It is not specified whether this practice applies to users with or without account.
Choice Type	Don't use service/feature	Only option is not to use the feature or service. Only select this if explicitly stated in policy (i.e., don't interpret silence as "Don't use website or feature").
Choice Type	Opt-in	User must consent before data can be shared with or collected/used by third party.
Choice Type	Opt-out link	Link provided in privacy policy, on website, in mobile app, or in email, etc.
Choice Type	Opt-out via contacting company	Must contact company/organization via email, phone, postal mail to opt-out.
Choice Type	First-party privacy controls	Website/app provides user settings for privacy configuration.
Choice Type	Third-party privacy controls	Choices provided by a third party (e.g., privacy settings on social media site) or industry (e.g., AdChoices Opt-out).



Choice Type	Browser/device privacy controls	Policy suggests the use of browser's or mobile device's privacy settings, e.g., to block trackers or cookies, activate Do-Not-Track, disable location sharing, clear history, etc.
Choice Type	Other	Other specific user choice or control option not captured above.
Choice Type	Unspecified	No user choices mentioned for this practice.
Choice Scope	Collection	Choices apply to collection by or sharing with third party only.
Choice Scope	Use	Choices apply to use by third party only.
Choice Scope	Both	Choices apply to both collection/sharing and use.
Choice Scope	Unspecified	No specific scope of choices is mentioned.

Table A.13: Tag Attribute Values for *Third Party Sharing/Collection* Privacy Concern

User Choice/Control		
Tag Attribute	Value	Description
Choice Type	Don't use service/feature	Only option is not to use the feature or service. Only select this if explicitly stated in policy (i.e., don't interpret silence as "Don't use website or feature").
Choice Type	Opt-in	User must consent before data can be shared with or collected/used by third party.
Choice Type	Opt-out link	Link provided in privacy policy, on website, in mobile app, or in email, etc.
Choice Type	Opt-out via contacting company	Must contact company/organization via email, phone, postal mail to opt-out.
Choice Type	First-party privacy controls	Website/app provides user settings for privacy configuration.

Choice Type	Third-party privacy controls	Choices provided by a third party (e.g., privacy settings on social media site) or industry (e.g., AdChoices Opt-out).
Choice Type	Browser/device privacy controls	Policy suggests the use of browser's or mobile device's privacy settings, e.g., to block trackers or cookies, activate Do-Not-Track, disable location sharing, clear history, etc.
Choice Type	Other	Other specific user choice or control option not captured above.
Choice Type	Unspecified	No user choices mentioned for this practice.
Choice Scope	First party collection	Choices apply to data collection by first party.
Choice Scope	First party use	Choices apply to the use of information by first party.
Choice Scope	Third party sharing/collection	Choices apply to data sharing with / collection by third party.
Choice Scope	Third party use	Choices apply to the use of information by third party.
Choice Scope	Unspecified	No specific scope of choices is mentioned.
Personal Information Type	Financial	Financial information, such as credit/debit card data, other payment information, credit scores, etc.
Personal Information Type	Health	Health Information, such as information about health conditions, prescriptions, medication, as well as health monitoring data, e.g., heart rate, step count, activity level, etc.
Personal Information Type	Contact	Contact Information, such as name, email address, phone number, street address, etc.
Personal Information Type	Location	Geo-location information (e.g., user's current location) regardless of granularity, i.e., could be exact location, ZIP code, city-level.

Personal Information Type	Demographic	Demographic Information, e.g., gender, age, occupation, education, etc.
Personal Information Type	Personal identifier	Identifiers that uniquely identify a person, e.g., SSN, driver's license number, etc.
Personal Information Type	User online activities	The user's online activities on the first party website/app or other websites/apps, e.g., pages visited, time spent on pages, general user behavior online, etc.
Personal Information Type	User profile	The user's profile on the first-party website/app and its contents, e.g., data in user profile, data that user uploaded to website, user comments, user profile preferences, etc. This is common for websites/apps where users can create an account or profile, e.g., on twitter, youtube, Facebook, Amazon, etc.
Personal Information Type	Social media data	User profile and data from a social media website/app or other third party service to which the user gave the first party access, e.g., by connecting with Facebook, twitter, or other services. Exchanged data may include user profile, photos, comments, friends, etc.
Personal Information Type	IP address and device IDs	Permanent (e.g., device IDs, MAC address) or temporary (e.g., IP address) identifiers needed to establish a connection for the current browsing session.

Personal Information Type	Cookies and tracking elements	Identifiers locally stored on user's device by company/organization or third-parties including cookies, beacons, or similar that are commonly used to uniquely identify users, but that are not essential to establish a connection with the user's device or to provide a service.
Personal Information Type	Computer information	The type of operating system (OS) or web browser that the user uses, or similar computer or device information.
Personal Information Type	Survey data	Any data that is collected through surveys
Personal Information Type	Generic personal information	No specific type of information is mentioned, but the policy talks about "personal information" or "personal identifiable information" in general.
Personal Information Type	Other	A specific type of information not covered by the above categories.
Personal Information Type	Unspecified	The type of information is not explicitly stated or unclear (e.g., refers to "information" very generically).
Purpose	Basic service / feature	Provide a service that the user explicitly requests and that is part of the website/app's basic service or functionality. Examples are watching a video, reading an article, making a purchase, creating an account, contacting the company/organization, etc.

Purpose	Additional service / feature	Provide a service that the user explicitly requests but that is not a necessary part of the website/app's basic service. Additional services/features may enhance user experience or add convenience but require additional data, e.g., social media integration, comments, blog participation, a store finder that needs location information, etc.
Purpose	Advertising	To show ads that are either targeted to the specific user or not targeted.
Purpose	Marketing	To contact the user to offer products, services, or other promotions (e.g., send marketing emails, calling or texting user with marketing messages). Marketing typically requires the use of contact information.
Purpose	Analytics / Research	For understanding the website/app's audience, improving the website/app, inform company strategy, or general research.
Purpose	Personalization & Customization	For providing user with a personalized experience, e.g., by allowing to arrange how the website/app looks, based on the user's preferences or language, etc.
Purpose	Service Operation and Security	For website/app operation and security, enforcement of terms of service, fraud prevention, protecting users and property, etc.
Purpose	Legal requirement	For compliance with legal obligations, e.g., regulations, government data requests, government retention requests, law enforcement requests in general, etc.

Purpose	Merger/Acquisition	If company/organization merges or is acquired it transfers users' information to another company/organization.
Purpose	Other	Other specific purpose not covered above.
Purpose	Unspecified	The purpose is not explicitly stated or is unclear.
User Type	User without account	This data practice specifically applies to users that do not have an account or are not registered with the website or mobile app.
User Type	User with account	This data practice specifically applies to users with an account or who are registered with the website or mobile app.
User Type	Other	This data practice applies to a specific user type not covered by the options above.
User Type	Unspecified	It is not specified whether this practice applies to users with or without account.

Table A.14: Tag Attribute Values for *User Choice/Control* Privacy Concern

User Access, Edit and Deletion		
Tag Attribute	Value	Description
Access Type	None	Users cannot access, edit, or delete data. Only select this if explicitly stated that users don't have access, otherwise select "Unspecified."
Access Type	View	Users can access their information, but not edit or delete it
Access Type	Export	Users can export their information to other services or download it to own computer.
Access Type	Edit information	User can modify or delete specific information

Access Type	Deactivate account	User can deactivate account so that the user's information is not visible to other users anymore, but the company/organization keeps all the data.
Access Type	Delete account (partial)	User can delete account, but the company/organization may continue to keep some of the user's data.
Access Type	Delete account (full)	User can delete account and all of the user's information is removed from company/organization's servers/databases.
Access Type	Other	An access, edit, or delete option not covered above.
Access Type	Unspecified	Access options are not mentioned or unclear.
Access Scope	User account data	Information explicitly provided by the user, such as contact, demographic, and any other explicitly provided information that is part of the user's profile/account. This includes user preferences and settings.
Access Scope	Transactional data	Purchases made, online activity, products watched, comments or questions submitted.
Access Scope	Profile data	Information that the company has learned about user, even if the user did not explicitly provide it
Access Scope	Other data about user	Other information that the company/organization has learned about the user, e.g., inferred preferences, data from other third parties, etc.
Access Scope	Other	A specific access scope is described that is not covered above.
Access Scope	Unspecified	Access scope is not mentioned or unclear.

User Type	User without account	This data practice specifically applies to users that do not have an account or are not registered with the website or mobile app.
User Type	User with account	This data practice specifically applies to users with an account or who are registered with the website or mobile app.
User Type	Other	This data practice applies to a specific user type not covered by the options above.
User Type	Unspecified	It is not specified whether this practice applies to users with or without account.

Table A.15: Tag Attribute Values for *User Access, Edit and Deletion* Privacy Concern

Data Retention		
Tag Attribute	Value	Description
Retention Period	Indefinitely	Collected user information is retained indefinitely.
Retention Period	Limited	Data is deleted, anonymized, or aggregated at some point, but no specific retention period is stated, e.g., "only stored as long as needed to perform requested service" or "as required by legal obligations".
Retention Period	Stated Period	Collected user information is deleted, anonymized or aggregated after a specific time period, e.g., "activity data is anonymized after 30 days".
Retention Period	Other	A specific retention type not covered above.
Retention Period	Unspecified	Retention period is not stated or unclear.
Retention Purpose	Perform service	Collected user information is only stored as long as it is needed to perform the requested service.



Retention Purpose	Legal requirement	Collected data is only stored as long as required for legal or law enforcement purposes.
Retention Purpose	Analytics/Research	For understanding the website/app's audience, improving the website/app, inform company strategy, or general research.
Retention Purpose	Service operation and security	For website/app operation and security, enforcement of terms of service, fraud prevention, protecting users and property, etc.
Retention Purpose	Advertising	To show ads that are either targeted to the specific user or not targeted.
Retention Purpose	Marketing	To contact the user to offer products, services, or other promotions (e.g., send marketing emails, calling or texting user with marketing messages). Marketing typically requires the use of contact information.
Retention Purpose	Other	Other specific retention purpose not covered above.
Retention Purpose	Unspecified	The retention purpose is not explicitly stated or is unclear.
Personal Information Type	Financial	Financial information, such as credit/debit card data, other payment information, credit scores, etc.
Personal Information Type	Health	Health Information, such as information about health conditions, prescriptions, medication, as well as health monitoring data, e.g., heart rate, step count, activity level, etc.
Personal Information Type	Contact	Contact Information, such as name, email address, phone number, street address, etc.

Personal Information Type	Location	Geo-location information (e.g., user's current location) regardless of granularity, i.e., could be exact location, ZIP code, city-level.
Personal Information Type	Demographic	Demographic Information, e.g., gender, age, occupation, education, etc.
Personal Information Type	Personal identifier	Identifiers that uniquely identify a person, e.g., SSN, driver's license number, etc.
Personal Information Type	User online activities	The user's online activities on the first party website/app or other websites/apps, e.g., pages visited, time spent on pages, general user behavior online, etc.
Personal Information Type	User profile	The user's profile on the first-party website/app and its contents, e.g., data in user profile, data that user uploaded to website, user comments, user profile preferences, etc. This is common for websites/apps where users can create an account or profile, e.g., on twitter, youtube, Facebook, Amazon, etc.
Personal Information Type	Social media data	User profile and data from a social media website/app or other third party service to which the user gave the first party access, e.g., by connecting with Facebook, twitter, or other services. Exchanged data may include user profile, photos, comments, friends, etc.
Personal Information Type	IP address and device IDs	Permanent (e.g., device IDs, MAC address) or temporary (e.g., IP address) identifiers needed to establish a connection for the current browsing session.

Personal Information Type	Cookies and tracking elements	Identifiers locally stored on user's device by website or third-parties including cookies, beacons, or similar that are commonly used to uniquely identify users, but that are not essential to establish a connection with the user's device or to provide a service.
Personal Information Type	Computer information	The type of operating system (OS) or web browser that the user uses, or similar computer or device information.
Personal Information Type	Survey data	Any data that is collected through surveys
Personal Information Type	Generic personal information	No specific type of information is mentioned, but the policy talks about "personal information" or "personal identifiable information" in general.
Personal Information Type	Other	A specific type of information not covered by the above categories.
Personal Information Type	Unspecified	The type of information is not explicitly stated or unclear (e.g., refers to "information" very generically).

Table A.16: Tag Attribute Values for *Data Retention* Privacy Concern

Data Security		
Tag Attribute	Value	Description
Security Measure	Secure data transfer	Data transfer between user and website/app is encrypted, e.g., SSL, TLS, HTTPS.
Security Measure	Secure user authentication	User authentication, e.g., login to a user account, is encrypted/secured.
Security Measure	Secure data storage	Data is stored securely, e.g. in an encrypted format or database.
Security Measure	Data access limitation	Data is accessible to employees/third parties on a need-to-know basis.

Security Measure	Privacy training	The company/organization trains its employees/third parties in applicable privacy and security practices to protect user data.
Security Measure	Privacy review/audit	Privacy practices and security measures of the first party or third party are reviewed/audited by internal or external reviewers/auditors.
Security Measure	Privacy/Security program	The company/organization has a privacy or security program/organization in place addressing, for example, how to protect data against unauthorized access or privacy training for employees.
Security Measure	Generic	The policy makes generic security statements, e.g., "we protect your data" or "we use technology/encryption to protect your data".
Security Measure	Other	A specific security measure not covered above.
Security Measure	Unspecified	Security measures are not mentioned or unclear.

Table A.17: Tag Attribute Values for *Data Security* Privacy Concern

Policy Change		
Tag Attribute	Value	Description
Change Type	Non-privacy relevant change	Minor change to the privacy policy that does not significantly affect data practices.
Change Type	Privacy relevant change	A change to the privacy policy significantly impacting current data practices, e.g., use, collection, sharing, retention, etc.
Change Type	In case of merger or acquisition	Users are notified if the policy changes as the result of a merger or acquisition.

Change Type	Other	Other specific policy change type not covered above.
Change Type	Unspecified	It is not mentioned or unclear for what kind of policy changes users are notified.
Notification Type	No notification	User is not notified of changes to the privacy policy.
Notification Type	General notice in privacy policy	The policy date is updated or information about the change is posted as part of the privacy policy site.
Notification Type	General notice on website	Users will be notified when visiting the main website, i.e., not only when looking at the privacy policy.
Notification Type	Personal notice	Users will be personally informed about a privacy policy change, e.g., via email, text message or when logging into their account.
Notification Type	Other	Users are notified in another specific way not covered above.
Notification Type	Unspecified	How users are notified about policy changes is not mentioned or unclear.
User Choice	None	The user has no options when the policy changes.
User Choice	Opt-out	Users can decline the new policy within a certain time period (e.g., 30 days), e.g., by canceling their account, opting-out of new practices, etc.
User Choice	Opt-in	User must agree before their data is collected/used/shared according to the new privacy policy.
User Choice	User participation	Users can decide or influence policy change (e.g., the company/organization proposes a change and asks for users' opinions).
User Choice	Other	Other specific user choice not covered above.

User Choice	Unspecified	Choices regarding policy changes are not mentioned or unclear.
-------------	-------------	--

Table A.18: Tag Attribute Values for *Policy Change* Privacy Concern

Do Not Track		
Tag Attribute	Value	Description
Do Not Track policy	Not mentioned	There is no statement concerning Do Not Track
Do Not Track policy	Honored	The website/app reads and adheres to the user's DNT preference
Do Not Track policy	Not honored	The website/app ignores DNT headers and the user's DNT preference.
Do Not Track policy	Mentioned, but unclear if honored	DNT headers are mentioned but it is unclear if the company/organization adheres to the user's DNT preference.
Do Not Track policy	Other	The website/app handles DNT headers in a different way not covered above.

Table A.19: Tag Attribute Values for *Do Not Track* Privacy Concern

International and Specific Audiences		
Tag Attribute	Value	Description
Audience Type	Californians	How data from Californian users is treated, e.g., California privacy rights.
Audience Type	Europeans	How data from European users is treated, e.g., Safe Harbor provisions.
Audience Type	Citizens from other countries	Specific provisions for international audiences or citizens from countries other than US or Europe, e.g., international data transfer.
Audience Type	Children	How data from children is treated.
Audience Type	Other	Other specific audience group not mentioned above.

Table A.20: Tag Attribute Values for *International and Specific Audiences* Privacy Concern

Other		
Tag Attribute	Value	Description
Other Type	Introductory / Generic	It's a paragraph that introduces the policy, a section, or a group of practices, but does not mention a specific privacy/data practice. The paragraph makes generic statements, but does not describe specific privacy/data practices
Other Type	Practice not covered	The paragraph describes a specific data practice, which is not covered by any of the other data practice categories.
Other Type	Privacy contact information	The paragraph describes how to contact the company with questions, concerns, or complaints about the privacy policy.
Other Type	Other	The paragraph does not fit any of the values above.

Table A.21: Tag Attribute Values for *Other* Privacy Concern





## Appendix B

# Evaluation Questionnaire

## 1/7. Demographics

\* Required

1. **Username used for Evaluation \***

---

2. **Education Level \***

*Mark only one oval.*

- ☐ High School or Lower
- ☐ Bachelor Degree (or similar)
- ☐ Master Degree (or similar)
- ☐ PhD Degree
- ☐ Other: 

---

3. **Are you worried about your data privacy? \***

*Mark only one oval.*

	1	2	3	4	5	
Never	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Always

## 2/7. Privacy Policies

4. **Have you read any Privacy Policy in the past? \***

*Mark only one oval.*

- ☐ Never
- ☐ I have taken a brief look at least once
- ☐ I have read specific parts of privacy policy at least once
- ☐ I have completely read a privacy policy at least once

5. **I found the document A hard to understand \***

*Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

6. I found the document B hard to understand \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

### 3/7. Platform Evaluation

7. The platform helps to create annotations \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

8. The platform helps to review annotations (vote up/down) \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

9. I find the platform user friendly \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

10. Which features did you find intuitive and easy to use? \*

---



---



---



---



---

11. Which features did you find confusing and difficult? Do you have suggestions for improvement? \*

---



---



---



---



---

### 4/7. Privacy Concerns

## 12. I find the Privacy Concerns categories easy to understand \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

## 13. I find the Privacy Concerns categories expressive enough \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

## 14. I find the Privacy Concerns attribute/values expressive enough \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

## 15. I find the Privacy Concerns attribute/values easy to understand \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

## 16. Any comments? \*

---



---



---



---



---

## 5/7. User Engagement

## 17. The user score kepted me engaged with the tasks \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

18. I believe that my final user score represents my effort \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

19. Do you have any suggestions to improve the engagement of the users? \*

---



---



---



---



---

## 6/7. Task Success

20. The time for the annotating tasks was enough \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

21. The time for the reviewing tasks was enough \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

22. I feel that I have successfully completed the task \*

Mark only one oval.

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

23. If not why? \*

---



---



---



---



---

## 7/7. General

24. **The task of creating annotations was difficult and frustrating \****Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

25. **The task of reviewing annotations was difficult and frustrating \****Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

26. **The task was getting easier while I was getting familiar \****Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

27. **I find the platform useful \****Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

28. **In case I want to use a new service, I would consider the valid annotations provided by the platform \****Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

29. **In case I want to use a new service, if there are not that many valid annotations in the corresponding PP I would review the available annotations \****Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

30. In case I want to use a new service, if there are not that many annotations in the corresponding PP I would consider reading the privacy policy and create new annotations \*

*Mark only one oval.*

	1	2	3	4	5	
Disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Agree

31. Do you have any suggestions to improve the platform? \*

---

---

---

---

---

32. Comments \*

---

---

---

---

---

---





## Appendix C

# ReST API

The REST api is included under the path /api of the platform It is seperated into data and view api under /api/data and /api/view/ respectively.

### C.1 Data API

<b>Title</b>	Get Annotation Details
<b>Description</b>	Returns the annodation details of the given annotation
<b>URL Path</b>	/api/data/product/{productId}/ /document/{documentId}/annotation/{id}
<b>Method</b>	GET
<b>Path params</b>	<ul style="list-style-type: none"><li>• {productId}: (integer) the id of the corren- sponding product where the annotation belongs</li><li>• {documentId}: (integer) the id of the corren- sponding document where the annotation be- longs</li><li>• {id}: (integer) the id of requested annotation</li></ul>
<b>URL params</b>	-
<b>Data params</b>	-

Table C.1: Data - Get Annotation Details

<b>Title</b>	Delete Annotation
<b>Description</b>	Deletes a specific annotation
<b>URL Path</b>	/api/data/product/{productId}/ /document/{documentId}/annotation/{id}
<b>Method</b>	DELETE
<b>Path params</b>	<ul style="list-style-type: none"> <li>• {productId}: (integer) the id of the corresponding product where the annotation belongs</li> <li>• {documentId}: (integer) the id of the corresponding document where the annotation belongs</li> <li>• {id}: (integer) the id of requested annotation</li> </ul>
<b>URL params</b>	-
<b>Data params</b>	-

Table C.2: Data - Delete Annotation

<b>Title</b>	Create Annotation
<b>Description</b>	Creates a new Annotation
<b>URL Path</b>	/api/data/product/{productId}/ /document/{documentId}/annotation/new
<b>Method</b>	POST
<b>Path params</b>	<ul style="list-style-type: none"> <li>• <b>{productId}</b>: (integer) the id of the corresponding product where the annotation belongs</li> <li>• <b>{documentId}</b>: (integer) the id of the corresponding document where the annotation belongs</li> </ul>
<b>URL params</b>	-
<b>Data params</b>	<ul style="list-style-type: none"> <li>• <b>textParts</b>: (List of Json objects) A list of json data that describes the highlighted parts of text. These data are provided by the TextHighlighter js library used by the platform</li> <li>• <b>privacyConcernId</b>: (Integer) The internal id that refers to the privacy concern defined in the platform</li> <li>• <b>tagIds</b>: (List of Integers) A list of integer ids that refer to the tag values under the same privacy concern defined by the platform</li> <li>• <b>comment</b>: (String, optional) A text string with an optional comment to be included in the annotation</li> </ul>

Table C.3: Data - Create Annotation

<b>Title</b>	Get Business Categories
<b>Description</b>	Returns all available business categories labels that are available to the platform
<b>URL Path</b>	/api/data/businessCategories/all
<b>Method</b>	GET
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.4: Data - Get Business Category Labels

<b>Title</b>	Submit New Vote
<b>Description</b>	Submits a new vote value for a given controversial entity <b>NOTE:</b> Resubmission of the same vote for the same controversial entity results to undo the vote that has been previously submitted
<b>URL Path</b>	/api/data/controversialEntity/vote/new
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	<ul style="list-style-type: none"> <li>• <b>controversialEntityId:</b> (integer) the id of the corresponding entity that the vote refers to</li> <li>• <b>value:</b> (integer) a positive value (typical +1) for vote up or negative integer (typical -1) to submit a downvote</li> </ul>

Table C.5: Data - Submit New Vote

<b>Title</b>	Delete Controversial Entity
<b>Description</b>	Deletes any controversial entity with the given id
<b>URL Path</b>	/api/data/controversialEntity/{id}
<b>Method</b>	DELETE
<b>Path params</b>	{id}: (integer) The id of the controversial entity to be deleted
<b>URL params</b>	-
<b>Data params</b>	-

Table C.6: Data - Delete Controversial Entity

<b>Title</b>	Submit New Document
<b>Description</b>	Submits a new document
<b>URL Path</b>	/api/data/document/new
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	<ul style="list-style-type: none"> <li>• <b>productId:</b> (integer) The product id for which the new document will be added</li> <li>• <b>language:</b> (integer) The internal id that refers to the language that the document is written</li> <li>• <b>documentType:</b> (integer) The internal id of the document type that the document belongs</li> <li>• <b>url:</b> (string) The url from where the document was/ can be downloaded</li> <li>• <b>htmlText:</b> (string, optional) The initial text of the document. If missing the document will be downloaded from the provided url</li> <li>• <b>lastModified:</b> (Date, optional) The date that corresponds to the last modified header of the downloaded document. Should be used only when htmlText param is provided.</li> <li>• <b>collectedAt:</b> (Date, optional) The date for which the document provided with htmlText had been downloaded. Should be used only when htmlText param is provided.</li> </ul>

Table C.7: Data - Submit New Document

<b>Title</b>	Get Document details
<b>Description</b>	Returns details for a specified document
<b>URL Path</b>	/api/data/document/{id}
<b>Method</b>	GET
<b>Path params</b>	{id}: The id of the document to be retrieved
<b>URL params</b>	-
<b>Data params</b>	-

Table C.8: Data - Get Document Details

<b>Title</b>	Delete Document
<b>Description</b>	Deletes a specified document
<b>URL Path</b>	/api/data/document/{id}
<b>Method</b>	DELETE
<b>Path params</b>	{id}: The id of the document to be deleted
<b>URL params</b>	-
<b>Data params</b>	-

Table C.9: Data - Delete Document

<b>Title</b>	Get Document Annotations
<b>Description</b>	Retrieves all annotations that have been created for the given document
<b>URL Path</b>	/api/data/document/{id}/annotations
<b>Method</b>	GET
<b>Path params</b>	{id}: The id of the document to be deleted
<b>URL params</b>	-
<b>Data params</b>	-

Table C.10: Data - Get Document Annotations

<b>Title</b>	Add Product
<b>Description</b>	Add a new product to the platform
<b>URL Path</b>	/api/data/product/new
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	<ul style="list-style-type: none"> <li>• <b>name:</b> (string) The name of product</li> <li>• <b>url:</b> (string) The webpage for which the product refers to</li> <li>• <b>imageUrl:</b> (string) A logo image URL for the created product. The logo image will be downloaded into the platform</li> <li>• <b>productType:</b> (integer) The internal id that corresponds to the product type that will be added</li> <li>• <b>businessCategories:</b> (list of integer) A comma seperated list of integer ids of the business category labels that the new product belongs to</li> </ul>

Table C.11: Data - Add New Product

<b>Title</b>	Get Digital Product Details
<b>Description</b>	Returns the details of the given product
<b>URL Path</b>	/api/data/product/{id}
<b>Method</b>	GET
<b>Path params</b>	{id}: The id of the product to be retrieved
<b>URL params</b>	-
<b>Data params</b>	-

Table C.12: Data - Get Product Details

<b>Title</b>	Get All Digital Products
<b>Description</b>	Returns all products that the platform contains
<b>URL Path</b>	/api/data/product/all
<b>Method</b>	GET
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.13: Data - Get All Products

<b>Title</b>	Delete Digital Product
<b>Description</b>	Deletes the given digital product from the platform
<b>URL Path</b>	/api/data/product/{id}
<b>Method</b>	DELETE
<b>Path params</b>	{id}: The id of the product to be deleted
<b>URL params</b>	-
<b>Data params</b>	-

Table C.14: Data - Delete Product

<b>Title</b>	Get All Tags
<b>Description</b>	Returns the list of available tags defined in the platform
<b>URL Path</b>	/api/data/tag/all
<b>Method</b>	GET
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.15: Data - Get All Tags

<b>Title</b>	Get Privacy Concern Tags
<b>Description</b>	Returns all tags defined in the platform that belong to a specific privacy concern category
<b>URL Path</b>	/api/data/tag/
<b>Method</b>	GET
<b>Path params</b>	-
<b>URL params</b>	<ul style="list-style-type: none"> <li>• <b>privacyConcern:</b> (integer) The internal id that refers to the privacy concern defined in the platform</li> </ul>
<b>Data params</b>	-

Table C.16: Data - Get Privacy Concern Tags

<b>Title</b>	Get User Details
<b>Description</b>	
<b>URL Path</b>	/api/data/user/{id}
<b>Method</b>	GET
<b>Path params</b>	{id}: the user id that identifies a specific user on the system
<b>URL params</b>	-
<b>Data params</b>	

Table C.17: Data - Get User Details



<b>Title</b>	Create New User
<b>Description</b>	
<b>URL Path</b>	/api/data/user/signup
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	<ul style="list-style-type: none"> <li>• <b>email:</b> (string) An email of the user</li> <li>• <b>username:</b> (string) The username of the user</li> <li>• <b>password:</b> (string) The password field of user credentials</li> <li>• <b>retypePassword:</b> (string) A string that matches the given password of the user</li> <li>• <b>age:</b> (integer, optional) A number indicating the age of the user</li> <li>• <b>educationLevel:</b> (string, optional) A string the describes the education level of the user</li> <li>• <b>profession:</b> (string, optional)</li> <li>• <b>details:</b> (string, optional) A string that adds some details about the user</li> <li>• <b>acceptTerms:</b> (boolean) A boolean that indicates that the user agree with the term and conditions while using the platform. Must be true in order the user account to be created</li> </ul>

Table C.18: Data - Create New User

<b>Title</b>	User Login
<b>Description</b>	Logins the user in platform. If the login is successful a session id cookie will be returned in the response.
<b>URL Path</b>	/api/data/user/login
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.19: Data - User Login

<b>Title</b>	User Logout
<b>Description</b>	Logout the current logged user from the platform and invalidates the corresponding cookie value. The cookie should be included to the request
<b>URL Path</b>	/api/data/user/logout
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.20: Data - User Logout

## C.2 View API

<b>Title</b>	Get Annotation Details
<b>Description</b>	Returns the rendered annotation details of the given annotation
<b>URL Path</b>	/api/view/product/{productId}/ /document/{documentId}/annotation/{id}
<b>Method</b>	GET
<b>Path params</b>	<ul style="list-style-type: none"> <li>• {productId}: (integer) the id of the corresponding product where the annotation belongs</li> <li>• {documentId}: (integer) the id of the corresponding document where the annotation belongs</li> <li>• {id}: (integer) the id of requested annotation</li> </ul>
<b>URL params</b>	-
<b>Data params</b>	-

Table C.21: View - Get Annotation Details

<b>Title</b>	Create Annotation
<b>Description</b>	Creates a new Annotation and returns the rendered HTML part of the annotation details of the created annotation
<b>URL Path</b>	/api/view/product/{productId}/ /document/{documentId}/annotation/new
<b>Method</b>	POST
<b>Path params</b>	<ul style="list-style-type: none"> <li>• <b>{productId}</b>: (integer) the id of the corresponding product where the annotation belongs</li> <li>• <b>{documentId}</b>: (integer) the id of the corresponding document where the annotation belongs</li> </ul>
<b>URL params</b>	-
<b>Data params</b>	<ul style="list-style-type: none"> <li>• <b>textParts</b>: (List of Json objects) A list of json data that describes the highlighted parts of text. These data are provided by the TextHighlighter js library used by the platform</li> <li>• <b>privacyConcernId</b>: (Integer) The internal id that refers to the privacy concern defined in the platform</li> <li>• <b>tagIds</b>: (List of Integers) A list of integer ids that refer to the tag values under the same privacy concern defined by the platform</li> <li>• <b>comment</b>: (String, optional) A text string with an optional comment to be included in the annotation</li> </ul>

Table C.22: View - Create Annotations

<b>Title</b>	Get User Details
<b>Description</b>	Renders and returns the HTML part that fills the user details modal
<b>URL Path</b>	/api/view/user/{id}
<b>Method</b>	GET
<b>Path params</b>	{id}: the user id that identifies a specific user on the system
<b>URL params</b>	-
<b>Data params</b>	

Table C.23: View - Get User Details

<b>Title</b>	User Login
<b>Description</b>	Logins the user in platform and returns the HTML that replaces the login part in the header bar.
<b>URL Path</b>	/api/view/user/login
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.24: View - User Login

<b>Title</b>	User Logout
<b>Description</b>	Logout the current logged user from the platform and invalidates the corresponding cookie value. Returns the HTML that replaces the login part in the header bar. The cookie should be included to the request
<b>URL Path</b>	/api/view/user/logout
<b>Method</b>	POST
<b>Path params</b>	-
<b>URL params</b>	-
<b>Data params</b>	-

Table C.25: View - User Logout



# Bibliography

- [1] JOEL R REIDENBERG. Privacy Harms and the Effectiveness of the Notice and Choice Framework. 11:40, 2014.
- [2] Pedro Giovanni Leon, Blase Ur, Yang Wang, Manya Sleeper, Rebecca Balebako, Richard Shay, Lujo Bauer, Mihai Christodorescu, and Lorrie Faith Cranor. What matters to users?: factors that affect users' willingness to share information with online advertisers. page 1. ACM Press, 2013.
- [3] Cory Hallam and Gianluca Zanella. Online self-disclosure: The privacy paradox explained as a temporally discounted balance between concerns and rewards. *Computers in Human Behavior*, 68:217–227, 2017.
- [4] Spyros Kokolakis. Privacy attitudes and privacy behaviour: A review of current research on the privacy paradox phenomenon. *Computers & Security*, 64(Supplement C):122–134, January 2017.
- [5] France Bélanger and Robert E. Crossler. Privacy in the Digital Age: A Review of Information Privacy Research in Information Systems. *MIS Q.*, 35(4):1017–1042, December 2011.
- [6] Protecting Consumer Privacy in an Era of Rapid Change: Recommendations For Businesses and Policymakers, March 2012.
- [7] Official California Legislative Information. The Online Privacy Protection Act of 2003, 2003.
- [8] PricewaterhouseCoopers. PDPA 2010. Laws of Malaysia, Act 709. personal Data Protection Act (PDPA), 2010.
- [9] Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation - GDPR) (Text with EEA relevance), May 2016.
- [10] Yuanxiang Li, Walter Stewart, Jake Zhu, and Anna Ni. Online Privacy Policy of the Thirty Dow Jones Corporations: Compliance with FTC Fair Information Practice Principles and Readability Assessment. 12(1):27, 2012.

- [11] Lorrie Faith Cranor, Candice Hoke, Pedro Giovanni Leon, and Alyssa Au. Are They Worth Reading? An In-Depth Analysis of Online Advertising Companies' Privacy Policies. page 23, 2014.
- [12] Hui Na Chua, Anthony Herbland, Siew Fan Wong, and Younghoon Chang. Compliance to personal data protection principles: A study of how organizations frame privacy policy notices. *Telematics and Informatics*, 34(4):157–170, July 2017.
- [13] The Platform for Privacy Preferences 1.0 (P3p1.0) Specification. page 104.
- [14] Paola Benassi. TRUSTe: An Online Privacy Seal Program. *Commun. ACM*, 42(2):56–59, February 1999.
- [15] Aleecia M McDonald, Robert W Reeder, Patrick Gage Kelley, and Lorrie Faith Cranor. A Comparative Study of Online Privacy Policies and Formats. page 19, 2009.
- [16] Aleecia M. McDonald and Lorrie Faith Cranor. The Cost of Reading Privacy Policies. *I/S: A Journal of Law and Policy for the Information Society*, 4:543, 2008.
- [17] Florian Schaub, Rebecca Balebako, Adam L Durity, and Lorrie Faith Cranor. A Design Space for Effective Privacy Notices. page 17.
- [18] Joel R. Reidenberg, Travis Breaux, Lorrie Faith Cranor, Brian French, Amanda Grannis, James T. Graves, Fei Liu, Aleecia McDonald, Thomas B. Norton, and Rohan Ramanath. Disagreeable Privacy Policies: Mismatches between Meaning and Users' Understanding. *Berkeley Technology Law Journal*, 30:39, 2015.
- [19] Carlos Jensen and Colin Potts. Privacy Policies As Decision-making Tools: An Evaluation of Online Privacy Notices. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '04, pages 471–478, New York, NY, USA, 2004. ACM.
- [20] Benjamin Fabian, Tatiana Ermakova, and Tino Lentz. Large-scale readability analysis of privacy policies. pages 18–25. ACM Press, 2017.
- [21] Ewa Luger, Stuart Moran, and Tom Rodden. Consent for All: Revealing the Hidden Complexity of Terms and Conditions. In *Conference on Human Factors in Computing Systems - Proceedings*, page [In Preparation/Press], April 2013.
- [22] Mark A. Graber, Donna M. D&apos, apos, Alessandro, and Jill Johnson-West. Reading level of privacy policies on Internet health Web sites. (Brief Report), July 2002.



- [23] Tatiana Ermakova, Benjamin Fabian, and Eleonora Babina. Readability of Privacy Policies of Healthcare Websites. page 16.
- [24] Gabriele Meiselwitz. Readability Assessment of Policies and Procedures of Social Networking Sites. In *Online Communities and Social Computing*, Lecture Notes in Computer Science, pages 67–75. Springer, Berlin, Heidelberg, July 2013.
- [25] Gitanjali Das, Cynthia Cheung, Camille Nebeker, Matthew Bietz, and Cinnamon Bloss. Privacy Policies for Apps Targeted Toward Youth: Descriptive Analysis of Readability. *JMIR mHealth and uHealth*, 6(1):e3, January 2018.
- [26] Janet Prichard, Bryant University, Kevin Mentzer, and Bryant University. AN ANALYSIS OF APP PRIVACY STATEMENTS. 18(4):10, 2017.
- [27] Norman Sadeh, Alessandro Acquisti, Travis D Breaux, Lorrie Faith Cranor, Noah A Smith, Fei Liu, Florian Schaub, et al. The Usable Privacy project, 2013.
- [28] A. Dara S.K. Cherivirala S. Zimmeck M.S. Andersen P.G. Leon E. Hovy N. Sadeh S. Wilson, F. Schaub. Demystifying Privacy Policies with Language Technologies: Progress and Challenges, 2016.
- [29] Shomir Wilson, Florian Schaub, Aswarth Abhilash Dara, Frederick Liu, Sushain Cherivirala, Pedro Giovanni Leon, Mads Schaarup Andersen, Sebastian Zimmeck, Kanthashree Mysore Sathyendra, N. Cameron Russell, and others. The Creation and Analysis of a Website Privacy Policy Corpus. In *ACL (1)*, 2016.
- [30] M.S. Andersen S. Wilson N. Sadeh J.R. Reidenberg S.K. Cherivirala, F. Schaub. Visualization and Interactive Exploration of Data Practices in Privacy Policies, 2016. SOUPS '16 Poster Session.
- [31] Alessandro Oltramari, Dhivya Piraviperumal, Florian Schaub, Shomir Wilson, Sushain Cherivirala, ThomasB. Norton, N.Cameron Russell, Peter Story, Joel Reidenberg, and Norman Sadeh. PrivOnto: A semantic framework for the analysis of privacy policies. *Semantic Web*, 9(2):185–203, January 2018.
- [32] Dhiren Audich. Enhancing Readability of Privacy Policies Through Ontologies. page 109, 2018.
- [33] Jaspreet Bhatia and Travis D. Breaux. Towards an information type lexicon for privacy policies. In *Requirements Engineering and Law (RELAW), 2015 IEEE Eighth International Workshop on*, pages 19–24. IEEE, 2015.
- [34] John W. Stamey and Ryan A. Rossi. Automatically identifying relations in privacy policies. ACM Press, 2009.

- [35] Rohan Ramanath, Fei Liu, Norman Sadeh, and Noah A. Smith. Unsupervised Alignment of Privacy Policies using Hidden Markov Models. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 605–610, Baltimore, Maryland, June 2014. Association for Computational Linguistics.
- [36] Fei Liu, Rohan Ramanath, Norman Sadeh, and Noah A. Smith. A Step Towards Usable Privacy Policy: Automatic Alignment of Privacy Statements. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 884–894, Dublin, Ireland, August 2014. Dublin City University and Association for Computational Linguistics.
- [37] Patrick Gage Kelley, Lucian Cesca, Joanna Bresee, and Lorrie Faith Cranor. Standardizing privacy notices: an online study of the nutrition label approach. page 1573. ACM Press, 2010.
- [38] Terms of Service; Didn’t Read, 2012.
- [39] Sebastian Zimmeck and Steven M Bellovin. Privee: An Architecture for Automatically Analyzing Web Privacy Policies. page 24, 2014.
- [40] Welderufael B. Tesfay, Peter Hofmann, Toru Nakamura, Shinsaku Kiyomoto, and Jetzabel Serna. PrivacyGuide: Towards an Implementation of the EU GDPR on Internet Privacy Policy Evaluation. pages 15–21. ACM Press, 2018.
- [41] Elisa Costante, Yuanhao Sun, Milan Petković, and Jerry den Hartog. A machine learning solution to assess privacy policy completeness: (short paper). page 91. ACM Press, 2012.
- [42] Paul André, Aniket Kittur, and Steven P. Dow. Crowd synthesis: extracting categories and clusters from complex data. pages 989–998. ACM Press, 2014.
- [43] Lydia B. Chilton, Greg Little, Darren Edge, Daniel S. Weld, and James A. Landay. Cascade: crowdsourcing taxonomy creation. page 1999. ACM Press, 2013.
- [44] Forbes Ave. CrowdForge: Crowdsourcing Complex Work. page 10.
- [45] Matteo Negri, Luisa Bentivogli, Yashar Mehdad, Danilo Giampiccolo, and Alessandro Marchetti. Divide and Conquer: Crowdsourcing the Creation of Cross-Lingual Textual Entailment Corpora. page 10.
- [46] T. D. Breaux and F. Schaub. Scaling requirements extraction to the crowd: Experiments with privacy policies. In *2014 IEEE 22nd International Requirements Engineering Conference (RE)*, pages 163–172, August 2014.
- [47] Shomir Wilson, Florian Schaub, Rohan Ramanath, Norman Sadeh, Fei Liu, Noah A. Smith, and Frederick Liu. Crowdsourcing Annotations for Websites’ Privacy Policies: Can It Really Work? pages 133–143. ACM Press, 2016.

- [48] Sylvie Szulman, François Lévy, and Eve Paul. OMTAT annotation tool: semantical enrichment for legal document search \*. page 8.
- [49] Kalina Bontcheva, Hamish Cunningham, Ian Roberts, Angus Roberts, Valentin Tablan, Niraj Aswani, and Genevieve Gorrell. GATE Teamware: a web-based, collaborative text annotation framework. *Language Resources and Evaluation*, 47(4):1007–1029, December 2013.
- [50] Brett Drury, Paula C. F. Cardoso, Jorge Valverde-Rebaza, Alan Valejo, Fabio Pereira, and Alneu de Andrade Lopes. An Open Source Tool for Crowdsourcing the Manual Annotation of Texts. In *Computational Processing of the Portuguese Language*, Lecture Notes in Computer Science, pages 268–273. Springer, Cham, October 2014.
- [51] Alexander J. Quinn and Benjamin B. Bederson. Human Computation: A Survey and Taxonomy of a Growing Field. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '11, pages 1403–1412, New York, NY, USA, 2011. ACM.
- [52] Ivo Blohm, Jan Marco Leimeister, and Helmut Krcmar. Crowdsourcing: How to Benefit from (Too) Many Great Ideas. page 14, 2013.
- [53] Lena Mamykina, Bella Manoim, Manas Mittal, George Hripcsak, and Björn Hartmann. Design lessons from the fastest q&a site in the west. page 2857. ACM Press, 2011.
- [54] Nguyen Quoc Viet Hung, Nguyen Thanh Tam, Lam Ngoc Tran, and Karl Aberer. An Evaluation of Aggregation Techniques in Crowdsourcing. In David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Terzopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, Xuemin Lin, Yannis Manolopoulos, Divesh Srivastava, and Guangyan Huang, editors, *Web Information Systems Engineering – WISE 2013*, volume 8181, pages 1–15. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.
- [55] Amiangshu Bosu, Christopher S. Corley, Dustin Heaton, Debarshi Chatterji, Jeffrey C. Carver, and Nicholas A. Kraft. Building reputation in StackOverflow: An empirical investigation. pages 89–92. IEEE, May 2013.
- [56] Dave Lee. Facebook scandal 'hit 87 million users'. *BBC News*, April 2018.