

# PrivOnto: a Semantic Framework for the Analysis of Privacy Policies

**Editor(s):** Name Surname, University, Country

**Solicited review(s):** Name Surname, University, Country

**Open review(s):** Name Surname, University, Country

Alessandro Oltramari <sup>a</sup>, Dhivya Piraviperumal <sup>a</sup>, Florian Schaub <sup>a</sup>, Shomir Wilson <sup>a</sup>, Norman Sadeh <sup>a</sup>, Joel Reidenberg <sup>b</sup>

<sup>a</sup> *Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, USA*

*E-mail: aoltrama@andrew.cmu.edu*

<sup>b</sup> *Fordham University School of Law, New York, NY 10023, USA*

**Abstract.** Privacy policies are intended to inform users about the collection and use of their data by websites, mobile apps and other services or appliances they interact with. This also includes informing users about any choices they might have regarding such data practices. However, few users read these often long privacy policies; and those who do have difficulty understanding them, because they are written in convoluted and ambiguous language. A promising approach to help overcome this situation revolves around semi-automatically annotating policies, using combinations of semantic technologies, machine learning and natural language processing to analyze them. In this article, we introduce PrivOnto, a semantic framework to represent annotated privacy policies with an ontology developed in collaboration with privacy experts. PrivOnto has been applied to a corpus of over 23,000 annotated data practices, extracted from a dataset of 115 privacy policies. We designed a collection of 57 SPARQL queries to extract information from the PrivOnto knowledge base, with the dual objective of (1) answering privacy questions users often have and (2) supporting researchers and regulators in the analysis of privacy policies at scale. We present respective findings, after examining the process of developing PrivOnto. Finally, we outline future research and open challenges in using semantic technologies for privacy policy analysis.

**Keywords:** Privacy policies, privacy technologies, ontology-based data access, SPARQL

## 1. Introduction

The quality of the digital world we live in depends on the quality of the information we have access to: this is the leitmotif of *Being Digital*, written by Nicholas Negroponte 20 years ago [25]. Since then, the frontiers of cyberspace have expanded considerably, and the technologies we use to organize and consume information have flourished. But in the present age of Big Data and the Internet of Things, the amounts of information are becoming increasingly difficult for individuals to manage, diminishing their situational awareness of the digital environment.

Given this trend, it is unsurprising that Internet users often skip reading the privacy policies of the websites, apps, and social media services they use [24,29]. The reasons for users' inattention are rooted in the often convoluted language of privacy policies, which are too complex and too voluminous to be read in a timely manner. This situation creates a disconnect between service providers and their consumers: privacy policies are legally binding documents, and their stipulations apply regardless of whether users read them. This disconnect between Internet users and the practices that apply to their data has led to the assessment that the "notice and choice" legal regime of online privacy is ineffective in the status quo [31]. Additionally,

policy regulators—who are tasked with assessing privacy practices and enforcing standards—are unable to assess privacy policies at scale. These shortcomings motivate work to retrieve salient details about privacy policies and reason about the practices they contain.

In this paper, we demonstrate how semantic technologies (STs) can be a viable and scalable approach to help address some of the problems that affect user privacy: STs can help consumers better understand the implications of their online activities, and support policy regulators in facing the intertwined challenges of preventing privacy abuse and reducing the information asymmetry between consumers and companies. Accordingly, we describe a semantic framework to model and analyze salient data practices from a corpus of annotated privacy policies developed in the Usable Privacy Policy (UPP) project.<sup>1</sup> The UPP project integrates machine learning, natural language processing and crowdsourcing to improve the analysis of privacy policies and facilitate the development of more accessible privacy notices by extracting and highlighting those data practices that are most relevant to users.

The rest of this article is structured as follows. First, we provide an overview of related work in Section 2. In Section 3, we describe an ontology of privacy policies populated with about 23,000 annotations of data practices. In Section 4, we illustrate the analysis of the obtained knowledge base with suitable SPARQL queries, designed to pinpoint relevant patterns of privacy practices in the annotated corpus. In Section 5, we conclude the paper with a discussion of open challenges and directions for future research.

## 2. Related Work

Privacy-enhancing technologies (PETs) can be defined as the ensemble of technical solutions that preserve the privacy of individuals in their interactions with technological systems. In a recent overview, Heurix et al. [16] categorize PETs along relevant dimensions of privacy, such as the types of data being processed or communicated, application scenarios, grounding in security models, presence of a trusted third party, etc. What their classification fails to account for, however, is the *knowledge dimension* in PETs: without empowering users with the adequate resources to better understand data collection, use and

sharing practices, their privacy awareness—the first barrier against any kind of violation—is hindered. In this regard, STs can be considered as knowledge-enabling solutions for PETs, and as support tools for developing context-aware applications [14,38,19,37].

According to Grau [13], to be used as effective privacy-preserving systems STs need to embody the following functionalities: (F1) *policy representation*, namely a declarative representation of policies in a system; (F2) *models of interaction*, i.e., a set of queries that can extract relevant information from the system; and (F3) *policy violation*, which formalizes the cases when user preferences and data practices collide, leading to consequences that put users' data at stake. These interconnected functionalities can emerge only when system development follows certain design stages, characterized by Grau as: identification of clear privacy requirements and translation into a suitable formal language; realization of the formalized requirements in a computational system; and analysis and verification of the instantiated requirements [22].

*PrivOnto*, the semantic framework we propose, strives to realize all three functionalities described above, adhering to the related design stages. To the best of our knowledge, most of the existing work on leveraging STs as PETs focuses on defining formal languages for privacy policy representation. For instance, Duma et al. [9] and De Coi & Olmedilla [7] have compared policy languages on the basis of theoretical (e.g., language expressivity) and empirical principles (action execution, extensibility, etc.). More recently, Bartolini et al. [1] created a legal domain ontology for data protection and privacy, and Breux et al. proposed 'Eddy' [3], a description logic designed to model privacy requirements, comparing it with alternative – yet less articulated – proposals like KAoS [39], ExPDT [33] and Rein [20]. Eddy has been used to detect conflicts in the specifications of privacy policies, but not yet at large scale. Formalizing policies in the context of description logics is also a goal of 'PeopleFinder' [34], a semantic web environment in which policies are expressed with a rule extension of OWL, and that enables users to selectively share their locations with others. More proposals for privacy specification languages exist, such as P3P [6], XACML [23], and EPAL [10], but they lack formal semantics.

Policy languages and domain ontologies are necessary to implement (F1) and (F3), but are not sufficient to realize (F2). Enabling (F2), namely identifying suitable queries to extract privacy information, is a data-intensive task. In the UPP project we address this is-

<sup>1</sup>Usable Privacy Policy Project: <https://www.usableprivacy.org/>

sue with an extensive data annotation effort conducted by domain experts. The centrality of (F2) is recognized by Kagal et al. [20] when outlining Rein. Rein is a semantic web framework for representing and reasoning over policies in domains that use different policy languages and knowledge expressed in OWL and RDF-S. Rein realizes a basic version of (F2): a rule-based inference engine checks for relations between a *requester*, a *resource* and some *access properties*. If a relation holds, the output will state whether the *request* is either *valid* or *invalid*. Kagal et al. note that to enhance the privacy and security of web applications more complex, yet user-friendly, query mechanisms need to be implemented. In the next sections, we articulate how this objective is being accomplished in our work by outlining *PrivOnto*'s architecture and core features. We illustrate how this semantic web framework can be used to model relevant data practices described in natural language privacy policies and augment context-awareness accordingly. We further discuss how *PrivOnto* can support privacy engineers and regulators in policy analysis, and become a powerful tool for end-users.

### 3. *PrivOnto*: Knowledge Base of Privacy Policies

The *PrivOnto* knowledge base is comprised of 913,544 RDF triples, obtained by populating a suitable domain ontology with 23,000 annotated data practices from a corpus of 115 privacy policies. *PrivOnto* merges a *bottom-up* and a *top-down* approach for ontology creation [26,36]: the former is illustrated in Section 3.1, where we describe the main categories and attributes identified by domain experts to capture data practices expressed in privacy policies; the latter is presented in Section 3.2, where we show how those conceptual structures are formalized as a domain ontology, which has been subsequently populated with a corpus of about 23,000 annotations of data practices. The corpus is described in Section 3.3.

#### 3.1. Domain Expert Frame Analysis of Privacy Policies

In order to extract the data practices described in privacy policies, a small group of domain experts (privacy experts, public policy experts, and legal scholars) first analyzed what practices are expressed in privacy policies and how they are commonly described. This analysis was partially informed by existing privacy frame-

works and prior work, such as the FTC's Fair Information Practices [12], the Platform for Privacy Preferences (P3P) [6], and specific privacy notice requirements prescribed by legislation, such as notice requirements in CalOPPA [27], COPPA [5], and the HIPAA Privacy Rule [17].

The domain experts analyzed multiple privacy policies to develop a collection of frames that codify the different data practice categories, their descriptive attributes, and typical attribute values expressed in privacy policies. Each frame has its own respective structure of frame-roles and values [11]. These frames were refined over multiple iterations involving their application to additional privacy policies and extensive discussions among the domain experts. The resulting collection of frames represents ten categories of data practices, which are defined as follows:

**First Party Collection/Use:** Privacy practice describing data collection or data use by the service provider operating the service, website or mobile app a privacy policy applies to.

**Third Party Sharing/Collection:** Privacy practice describing data sharing with third parties or data collection by third parties. A third party is a company or organization other than the first party service provider operating the service, website or mobile app.

**User Choice/Control:** A practice describing general choices and control options available to users.

**User Access, Edit, & Deletion:** A practice describing if and how users may access, edit or delete the data that the service provider has about them.

**Data Retention:** A practice specifying the period and purposes for which collected user information is retained.

**Data Security:** A practice describing how user data is secured and protected, e.g., from confidentiality, integrity, or availability breaches.

**Policy Change:** A practice on whether and how the service provider informs users about changes to the privacy policy, including any choices offered to users.

**Do Not Track:** A practice specifying if and how Do Not Track signals (DNT)<sup>2</sup> for online tracking and advertising are honored.

**International & Specific Audiences:** A Practice that pertains only to a specific group of users, e.g., children, California residents, or Europeans.

<sup>2</sup>W3C Tracking Protection WG:  
<https://www.w3.org/2011/tracking-protection/>

**Other:** Additional sub-labels for introductory or general text in the privacy policy, contact information, and practices not covered by other categories.

A *data practice* statement belongs to one of these categories, and is characterized by a category-specific set of *attributes*. The frames define a set of potential values for each attribute. Each attribute is supported by a text *fragment* in the privacy policy, which serves as the natural language evidence for the annotated attribute value.

For example, a *First Party Collection/Use* practice is represented by four mandatory and five optional attributes. The mandatory attributes are whether the practice is a positive or negated statement (*Does* or *DoesNot*), how the first party obtained information (*action-first-party*), what kind of information is collected (*personal-information-type*), and for what purpose (*purpose*). In addition, a first party practice statement may indicate whether information is collected implicitly or if the user explicitly provides information (*collection-mode*), whether collected information is linkable to a user’s identity (*identifiability*), whether the practice applies to registered users only (*user-type*), and if a user choice is offered explicitly for this practice (*choice-type* and *choice-scope*). Data practices in other categories are represented with similar sets of attributes.

Mandatory and optional attributes reflect the level of specificity with which a specific data practice is typically described in privacy policies. Optional attributes are less common, while mandatory attributes are essential to a data practice. However, privacy policies are often ambiguous on many of these attributes [30]. Therefore, a valid value for each attribute is *Unspecified*, allowing annotators to express an absence of information. For instance, the fragment “we disclose information to third parties only in aggregate or de-identified form” exemplifies vagueness in data practices.

This collection of data practice frames constitutes the semantic foundation for the *PrivOnto* ontology, described in the next section.

### 3.2. Domain Ontology for Privacy Policies

The *PrivOnto* ontology is a formal model of the data practices identified by domain experts. It represents unstructured policy contents according to frame-based structures specified using OWL-DL. In *PrivOnto*, each

data practice category is modeled as a class characterized by a wide spectrum of Object and Datatype properties (see Figure 1): we used the latter to represent the specific attributes of each category, which essentially correspond to the backbone of the collection of frames presented in the previous section; conversely, the former were used to represent the conceptualization of the domain, and delineate the semantic relations holding between the defined classes.

The Object property *denote* holds between the class *ANNOTATION* and the class *SEGMENT*: the resulting pattern captures the difference between annotations, namely the entities that emerge from tagging discrete parts of privacy policies with suitable frames and roles, and the specific text they refer to. Accordingly, individual annotations *denote* individual segments (policy paragraphs) and their constituent parts or *fragments*. The class *SEGMENT* and the class *FRAGMENT* are linked by the *part\_of* relation, which is axiomatized as asymmetric and ir-reflexive. This semantic structure reflects the compositionality of paragraph-length segments: fragments can span from single words to well-formed sentences, whereas segments correspond to syntactically and semantically coherent sequences of fragments. By means of the *part\_of* relation, the same segment can instantiate multiple data practices via its fragments.

Fragments are labeled with a unique identifier (UID), consisting of the policy number, the segment number, and the start and end indexes of the selected text. In the same way, we assigned UIDs to instances of practice categories. Thanks to this modeling strategy, we can refer to different annotations of the same fragment, so that the “raw” policy content is kept distinct from all the annotations that refer to it. For example, a fragment stating that “by use of our websites and games that have advertising, you signify your assent to SCEA’s privacy policy” is annotated as an instance of **First Party Collection** and as an instance of **User Choice**, reflecting different aspects of the policy text. This situation can be represented in *PrivOnto* by two instances of *ANNOTATION*, each exemplifying different data practice categories, and referring to the same individual of *FRAGMENT*. The actual content of a fragment is expressed in the form of ‘string’ values in the range of the *annotated\_text* datatype property, whose domain is the *FRAGMENT* class. For example fragment 3819-3-95-203 is associated with the following statement “The information we learn from customers helps us personalize and continually

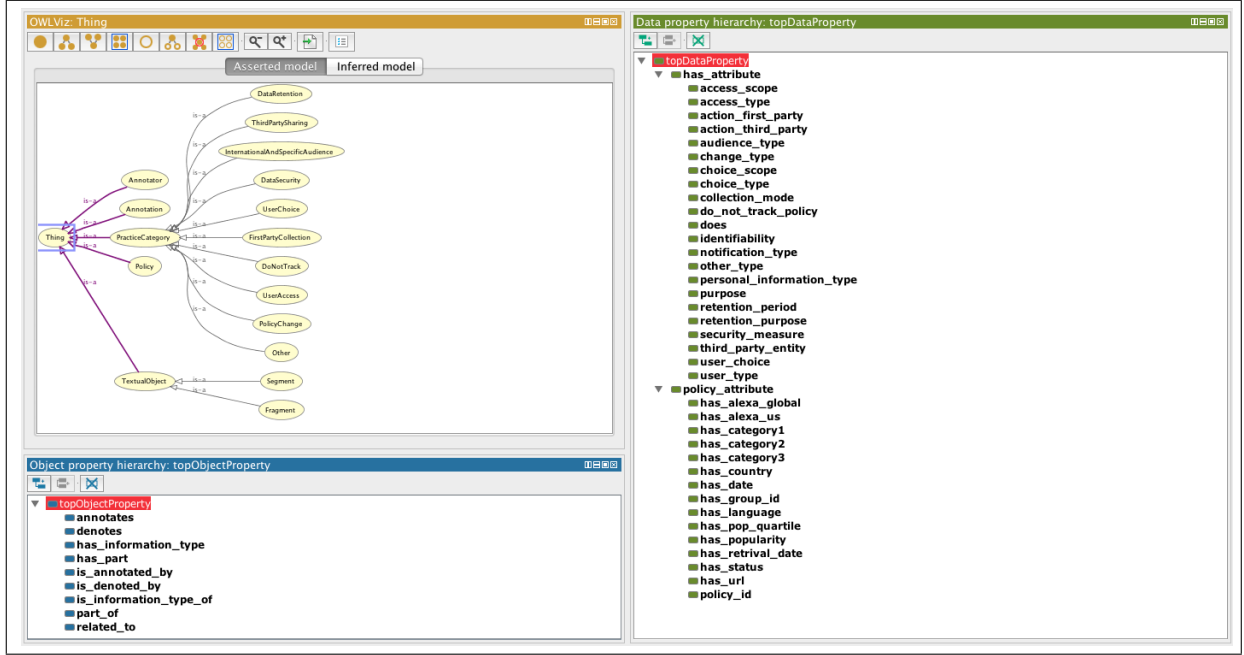


Fig. 1. Protégé visualization of *PrivOnto* hierarchies of Classes, Object properties and Datatype Properties.

improve your Amazon experience." This fragment is used in Figure 2, which shows how annotations, data practice categories and fragments are connected in the ontology. The *PrivOnto* framework does not directly address the linguistic structures of a given policy, but it pinpoints them only insofar as they instantiate a data practice category: we demonstrate in Section 4 how this is actually a key strength of our approach.

The ontology also includes *ANNOTATOR*, a class whose instances denote the individuals involved in the annotation task: the relation *executed\_by* between *ANNOTATION* and *ANNOTATOR* preserves the traceability of the identified data practices.

*PrivOnto* also includes general information about the website where the privacy policy can be found: the date when it was crawled, contact information of the company to which the policy belongs, the company's website, the associated Alexa's traffic ranking information,<sup>3</sup> etc. Note that some of this 'meta-information' is subject to change, and thus needs to be regularly monitored and documented: to this end, *PrivOnto* supports *xsd:dateTime* values, which serve as temporal indexes for policies' meta-information. Privacy policies may vary over time as well: in this case it is not only important to record changes, but also to in-

vestigate their implications: policies are systematically updated by companies for a variety of reasons, and analyzing the consequences of these modifications to enforced data practices is of key importance to regulators and users. The privacy policies obtained for annotation were collected at the same time, thus policy changes do not occur in our dataset. Nevertheless, future expansion of our corpus will include the addition of new privacy policies along with updates to already represented policies. We therefore plan to extend *PrivOnto* with *OWL-Time*<sup>4</sup> to enable qualitative and quantitative temporal reasoning [18].

### 3.3. Corpus of Annotated Privacy Policies

Privacy policies vary along many dimensions of analysis, including length, legal sophistication, readability, coverage of services, and update frequency. Large companies' policies may cover multiple apps, services, websites, and retail outlets, while privacy policies of smaller companies may have narrower scope. Accordingly, privacy policies were chosen for inclusion in the UPP corpus using a procedure that encouraged diversity.

<sup>3</sup><http://www.alexacom/topsites/countries/US>

<sup>4</sup><https://www.w3.org/2001/sw/BestPractices/OEP/Time-Ontology-20060518>

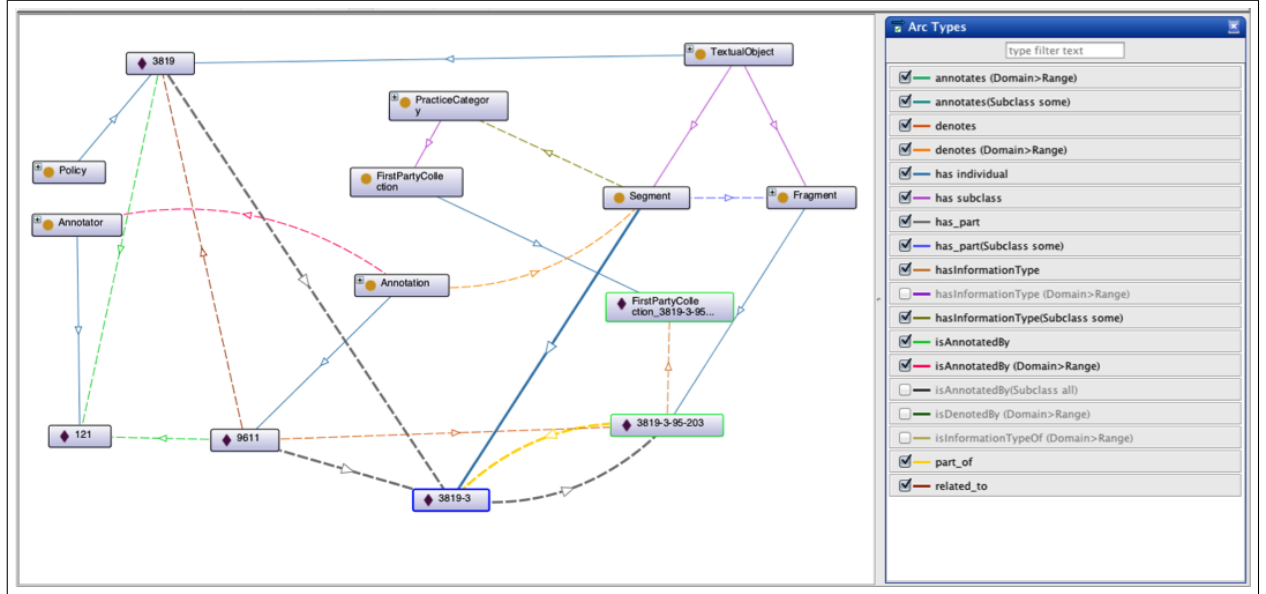


Fig. 2. LEFT: an example that shows how *PrivOnto* structures are used to model the semantic relations between data practices, fragments and segments of policies. RIGHT: legenda of semantic relations (redundant arcs are grayed-out to simplify the figure).

Websites were selected using a two-stage process: (1) relevance-based website pre-selection and (2) sector-based subsampling. This first stage consisted of monitoring Google Trends [15] for one month (May 2015) to collect the top five search queries for each trend; then, for each query, the first five websites were retrieved on each of the first ten pages of search results. This produced a selection of 1,799 unique websites. For the second stage, websites were chosen from each of DMOZ.org’s top-level website sectors.<sup>5</sup> However, for uniformity, selection from geographic sectors was restricted to ensure that all privacy policies in the corpus were subject to US legal and regulatory requirements.

For each sector, eight websites were selected based on occurrence frequency in Google search results. More specifically, the eight websites were randomly selected two-apiece from each rank quartile. Each selected website was manually verified to have an English-language privacy policy and to belong to a US company (according to contact information and the website’s WHOIS entry). Websites that did not meet these requirements were replaced with random redraws from the same sector and rank quartile. Notably, some privacy policies covered more than one selected website (e.g., the Disney privacy policy covered

disney.go.com and espn.go.com). The consolidation of the corpus resulted in a final dataset of 115 privacy policies across 15 sectors.

We developed a web-based annotation tool, shown in Figure 3, to facilitate annotation of the UPP corpus’ privacy policies by expert annotators according to our frame-based annotation scheme. Privacy policies were divided into *segments* and shown to annotators sequentially in the tool. Each segment may be annotated with zero or more data practices from each category. To annotate a segment with a data practice, an annotator assigns a practice category and specifies values and respective text spans (*fragments*) as appropriate for each of its attributes.

Each privacy policy was independently annotated by three expert annotators. In total, we hired 10 law students as experts on an hourly basis to annotate the complete set of 115 privacy policies. Note that the average annotation time per policy was 72 minutes. The annotation of the corpus resulted in about 23,000 annotations of data practices, which were used to populate the *PrivOnto* ontology and create the corresponding knowledge base.

#### 4. Query-based Semantic Analysis of Privacy Policies

*PrivOnto* facilitates the elicitation of prominent information from privacy policies in order to gain in-

<sup>5</sup>The DMOZ.org website sectors are notable for their use by Alexa.com.

Current Policy: a\_98\_neworleansonline.com

First Party Collection/Use Third Party Sharing/Collection

User Choice/Control User Access, Edit and Deletion

Data Retention Data Security Policy Change Do Not Track

International and Specific Audiences Other

7/41 Annotated Practices: 1

Previous Next

**Information We Collect**

Whether you access our Online Services from **your computer**, smart phone, tablet or other mobile device, NOTMC and its agents **may** collect some information that **identifies you or relates to you as an individual** ("Personal Information"), such as your **name, mailing address, telephone number, e-mail address, user name and password** (for account administration), device ID, including IP address, geolocation (if using a mobile application and you consent to providing it), and additional personal information necessary for the administration of certain promotional events.

**First Party Collection/Use**

- Does/Does Not: Does
- Collection Mode: Unspecified
- Action First-Party: Collect on website
- Identifiability: Identifiable
- Personal Information Type: Contact
- Purpose: Unspecified
- User Type: Unspecified
- Choice Type: Unspecified
- Choice Scope: Unspecified
- ☐ References another place in the policy

Save

Fig. 3. Web-based tool for expert privacy policy annotation.

sights on the nature of data practices. This knowledge elicitation process leverages a library of 57 SPARQL queries<sup>6</sup> we engineered to retrieve data practice categories, attributes, and values from the annotated corpus.<sup>7</sup> Our work required only marginal effort for translating unstructured natural language questions into formal queries, as our frame-based annotation process embedded ‘saliency’ in the corpus of annotations in the form of ontology categories and attributes. For this reason, the ontology-based analysis of privacy policies proposed in this article did not require dealing with the diversity and ambiguity of natural language text [21]. The queries we present in Section 4.2 match *by design* the privacy questions that domain experts deemed as relevant for policy analysis, and that originated the *PrivOnto* framework in the first place.

#### 4.1. Architecture

Our architecture for mapping the structured annotation corpus to the *PrivOnto* ontology is shown in Figure 4. The mapping process resulted in a .owl file that

captures the corpus (913,544 RDF triples). The obtained knowledge base was then loaded in an Apache Jena Fuseki server<sup>8</sup> for dynamic processing: the server provides a web service framework for different applications to access data through SPARQL queries. Figure 5 shows the *PrivOnto* semantic web environment.

#### 4.2. Library of Queries

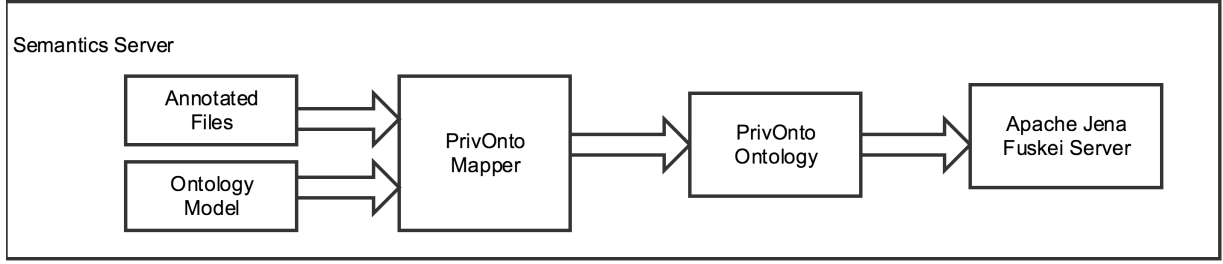
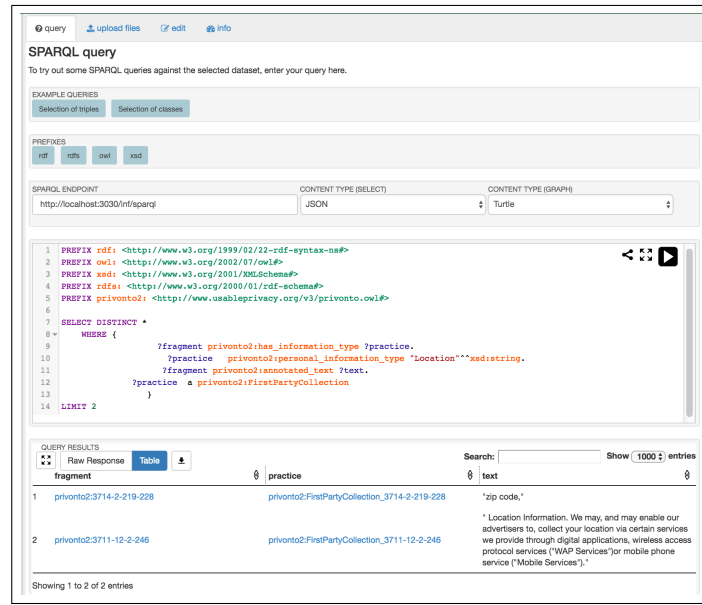
We created 57 SPARQL queries to analyze different aspects of the 115 privacy policies represented in the *PrivOnto* ontology: this method enabled us to build a scalable semantic retrieval system for gaining insights on privacy practices related to the collection, use, and sharing of personal data. The queries in the library can be categorized by two orthogonal dimensions, based on: (1) the type of targeted information (quantitative, qualitative, truth-values) and (2) the selected practice category.

It is important to point out that all 57 queries return the annotated text associated with a policy fragment: this feature realizes a crucial aspect of *model of interaction* (see functionality F2 in Section 2), i.e., the possibility for legal experts and users to understand and

<sup>6</sup>Version 1.1: <https://www.w3.org/TR/2013/REC-sparql11-query-20130321/>

<sup>7</sup>Despite being extensive and detailed, this library is not meant to be exhaustive, and can be further expanded.

<sup>8</sup><https://jena.apache.org/download/index.cgi>

Fig. 4. Semantic server architecture for querying *PrivOnto*.Fig. 5. Screenshot of the Apache Jena Fuseki server used for querying *PrivOnto*: the query in the example returns two policy fragments about collection of location information. Note that the `LIMIT 2` clause was used to fit the results to the window's size.

evaluate the machine-readable semantic models and queries in relation to a privacy policy's original text.

Table 1 shows the different kinds of information that can be extracted from the knowledge base, along with sample queries. Percentage and count type questions help gain an overall understanding of the privacy policy data. For example the query below, which calculates the 'number of policies that allow users to export their data,' returns 1 as the answer. Thus, only one out of 115 policies provides for the export of collected data, which shows the exceptionality of this data practice in the considered dataset.

```

SELECT (COUNT(*) AS ?count) {SELECT DISTINCT ?policy
WHERE {?p a privonto2:UserAccess.
      privonto2:access_type "Export"^^xsd:string.
      privonto2:related_to ?policy.}
  
```

In order to verify facts in the ontology, we can use ASK queries. For instance, the query below, which matches the question 'Does any policy state that personal information is shared or collected as part of a merger?', returns True as output. By replacing the ASK clause with a SELECT clause, we can easily assess that nine policies include that data practice.

```

ASK
WHERE
{?frag privonto2:part_of ?segment.
 ?frag privonto2:has_information_type ?practice.
 ?prc privonto2:purpose "Merger/Acq"^^xsd:string.
 ?prc privonto2:related_to ?policy.
 ?prc a privonto2:FirstPartyCollection.}
  
```

Our SPARQL queries also help gain specific information about different practice categories. For instance, the query exemplified by the question 'How



many websites mention each audience type?’ lead us to discover that clauses are generally added for children (86 out of 115 privacy policies), which suggests that a large number of privacy policies aim to be compliant with the Children Online Privacy Protection Act (COPPA) [5], but also shows that 25% of the privacy policies in our corpus have no provisions specific to children.

The second dimension through which our SPARQL queries can be classified is based on different practice categories. Each practice category provides very specific information about privacy policies. By organizing the queries in this way, we can concentrate on specific characteristics of a policy, and draw parallel conclusions from different categories. Table 2 shows example queries from each category.

While running experiments in the Jena Fuseki environment, we observed that the queries’ processing time depends on the complexity of the SPARQL expression, while being only partially correlated with the number of matches. In particular, Figure 6 represents the proportion between number of matches and retrieval times for a subset of 20 SPARQL queries chosen across all data practice categories to highlight relevant types of information in a policy. For instance, the figure shows that only four queries had processing time higher than 1500 ms: these queries included SPARQL constraints like `OPTIONAL` and `MINUS`. The queries labeled as ‘Financial Information and Purpose’, ‘General Information and Purpose’, ‘Unspecified Information and Purpose’ refer to user’s collected information at different levels of granularity, and specify the purpose of collection only when found in a policy: this condition was expressed in the SPARQL request by an `OPTIONAL` clause on the ‘Purpose’ attribute of the ‘First Party Collection/Use’ category. In the case of the query labeled as ‘Policies with User Choice,’ the high processing time was brought about by the `MINUS` clause, introduced to discard from the results all the policies with no real user choice, but only with take-it-or-leave-it option (this aspect is further analyzed in section 4.3.3).

### 4.3. Results

In this section we overview the quantitative and qualitative results of our query-based semantic analysis of about 23,000 data practices instantiated in the *PrivOnto* knowledge base.

#### 4.3.1. Personal information collection/sharing

For the practice categories *User Choice*, *First Party Collection/Use*, and *Third Party Sharing/Collection*, we observed that privacy policies specify the information collected or shared, though the purpose of data collection is rarely mentioned in the same fragment. Therefore, we collected the purpose information from the other fragments present in the parent segment. We observed that, apart from ‘unspecified,’ ‘basic service’ and ‘additional service’ were the most mentioned purposes. ‘Device information’ and user’s ‘online activity’ are collected from users’ for ‘analytics/research’ purposes, whereas ‘finance’ and ‘contact information’ were collected for ‘marketing’ and ‘advertising purposes.’ Purpose for which information is highly shared is ‘Advertising’ (14.6%), and the purpose for which information is highly collected is for ‘basic service/feature’ (16%).

Table 3 presents the comparison of different personal data types which are collected and shared. We observed that most of the data types collected and shared are unspecified (last row). This result can be explained by the fact that the word “information” is often used with no further description or specification in the policies. As a result, the privacy policies make it difficult for consumers and regulators to determine which information is actually collected or shared by a company. The following text fragments exemplify this vagueness: “the information we learn from customers helps us personalize and continually improve your Amazon experience” and “any information that we collect from or about you.”

Table 3 also shows that ‘device,’ ‘location identifiers,’ and ‘contact information’ are often collected by the websites, but are not explicitly mentioned in statements with respect to third party sharing. Because of the extensive use of generic descriptions for information types, the privacy policies do not indicate whether these data items are actually shared with third parties.

‘Contact information,’ ‘user online activities,’ and ‘general personal information’ are the top referenced types of information. ‘Contact information’ appears frequently as collected information, while ‘general personal information’ is highly shared. ‘General personal information’ is also often ambiguous. The corresponding policy fragments describe this information as “personally identifiable information” or “personal information.” For example, one policy in the corpus shares “any and all personal identifiable information collected from our customers” with third parties.

Table 1  
Targeted information and related query types.

Targeted Information	Query example
Percentage	What percentage of policies apply to websites and mobile apps?
Count on Practices	How many practice statements per policy are unclear about where information are collected from users?
True or False	Is information shared or collected as part of a merger or acquisition?
Count on Policy	How many policies have statements on user choice?
Count on distribution of policies across values in practice category	For each of the security-measure values, how many websites mention them?

Table 2

Queries are sent to the Apache Jena Fuseki server that runs the *PrivOnto* framework: quantitative results shown in the table indicate the number of fragments, number of policies, and percentages related to specific data practices.

Category	Type of Queries	Result
First Party Collection	Fragments that collect finance information and for what purpose?	231
Third Party Sharing	Fragments that denote user information is shared with external third parties	2,220
User Choice	How many policies have statements on user choice?	106
User Access	Percentage of policies that allow users to delete their account	0.18
Data Retention	Percentage of statements where a period is stated for data retention	0.09
Data Security	For each of the security-measure values, how many websites mention them?	10
Policy Change	How many websites specify a user choice on policy change?	91

Table 3

Queries on information collected from users or shared about users. Number of fragments are visualized, as well as coverage across policies,

Question	First Party Collection	% Policies	Third Party Collection	% Policies
Fragments that collect/share location information and for what purpose?	265	59.13	61	26.09
Fragments that collect/share contact information and for what purpose?	736	90.43	246	57.39
Fragments that collect/share device identifier and for what purpose?	319	76.52	75	25.22
What kind of Fragments are especially negated	199	67.83	313	78.26
Fragments that collect/share finance info and for what purpose?	231	63.48	102	35.65
Fragments that collect/share user's online activities info and for what purpose?	559	87.83	294	66.96
Fragments that collect/share user's general personal information info and for what purpose?	587	88.70	730	91.30
Fragments that collect/share user's unspecified info and for what purpose?	936	85.22	820	88.70

Out of 115 policies, 90 privacy policies state that the service providers do not share some information with

third parties, and 78 policies explicitly state what in-

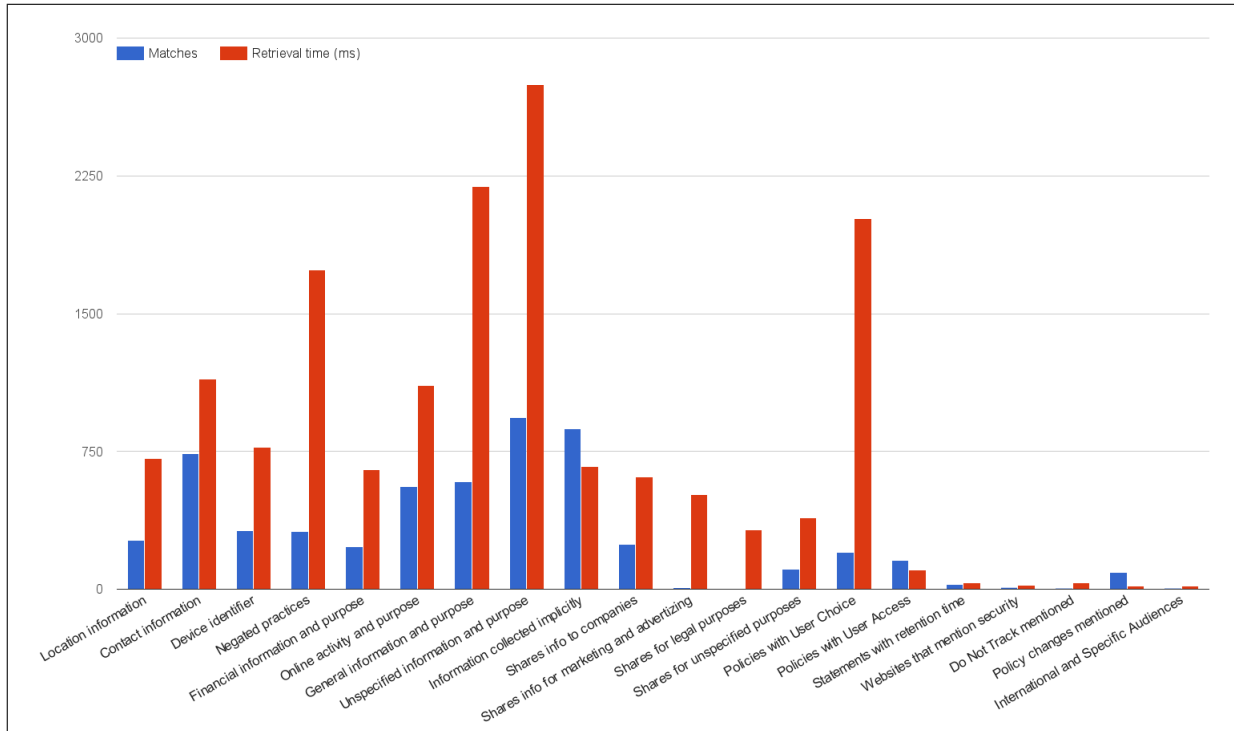


Fig. 6. Proportion between number of matches and processing times for a subset of 20 queries. The labels in the *x-axis* represent types of information collected, shared, or mentioned in a policy and returned by suitable SPARQL queries. The *y-axis* represents the corresponding number of matches (blue histograms) and the retrieval time in milliseconds (red histograms).

formation they do not collect from users. The top categories of information type reportedly not collected or not shared are ‘generic personal information,’ ‘cookies and tracking elements,’ and ‘contact’ information. While this appears to contradict the previous finding that contact information is frequently collected and general personal information is widely shared, the contradiction reflects that privacy policies are explicit when they do not share data.

#### 4.3.2. Marketing and Advertising

There were 886 fragments which described the collection of information for ‘Marketing’ and ‘Advertising’ purposes. Information collected for advertising purposes is typically identified as the user’s ‘online activities’ or ‘cookies and tracking elements’. Users’ ‘contact’ information is typically used for ‘marketing’ purposes.’ By contrast, ‘financial’ information is often identified for sharing with third parties when these are partners or affiliates.

#### 4.3.3. User’s choice on enabling service

Almost all privacy policies (92%) have statements describing User Choices. But, of these privacy policies, 48% have statements that merely describe a take-

it-or-leave-it choice. Instead of a real choice, users are told not to use the service or feature if they disagree with the privacy policy or with certain data practices. Examples are: “if you choose to decline cookies, you may not be able to fully experience the interactive features of this or other Web sites you visit” or “if you do not agree to this privacy policy, you should not use or access any of our sites.”

#### 4.3.4. User Data Retention

About half of the privacy policies (56%) specify for how long they store user data. In 40% of these policies a retention period is explicitly ‘stated’ (e.g., 30 days) or the retention period is at least ‘limited’ (e.g., stored as long as needed to perform a requested service); while 7% express that the data will be stored indefinitely. The distinction between ‘Limited’ and ‘Stated’ retention periods is sometimes blurred due to drafting vagueness and annotator interpretation. For instance, the fragment “we will retain your data for as long as you use the online services and for a reasonable time thereafter” has been annotated both as “limited period” or as “stated period.” This creates ambiguity with re-

spect to the duration that user data will remain in a service's database.

#### 4.3.5. Data Export

As mentioned in the previous section, only one policy in our knowledge base describes how users can export data. The respective annotated fragment states: "California Civil Code Section 1798.83, also known as the Shine The Light law, permits our users who are California residents to request and obtain from us once a year, free of charge, information about the personal information (if any) we disclosed to third parties for direct marketing purposes in the preceding calendar year."

#### 4.3.6. Policy Change

Privacy policies typically provide that users are notified about changes to the privacy policy through some form of general notice or through a website. Only 30% of the privacy policies containing descriptions of change in notification practices mention a notification of individual users (e.g., via email). The lack of personal notice for policy changes means that users are unlikely to be aware of changes to the privacy policy, although such changes may alter how information about them is collected, used, or shared by a service.

#### 4.3.7. Data Security

The major security measures which most websites describe are the use of 'secure user authentication,' the existence of a 'privacy/security program,' and the communication of data with 'secure data transfer.'

The analysis above shows that query-based analysis of the *PrivOnto* knowledge base can provide insights on privacy policy data both on a semantic and textual level. We can both verify information and collect statistics on privacy policies by means of the *PrivOnto* semantic framework. Ontology-driven analysis can help distill the content of a privacy policy, as well as help compare the target policy with similar policies. In this respect, *PrivOnto* can help users gain insights on the stated practices of services they use and help them make more informed privacy choices.

## 5. Discussion and Future work

In this paper we described *PrivOnto*, a semantic web framework used to represent data practices in privacy policies and support knowledge elicitation. *PrivOnto* is an essential tool for regulators and can

also enable more usable privacy notices by exposing semantic reasoning results to users. We show the utility of *PrivOnto* by instantiating it with a corpus of 115 privacy policies which have been annotated by domain experts as part of the Usable Privacy Policy project.

The *PrivOnto* ontology model formalizes a frame-based annotation scheme that helps experts identify data practices in policy text. As a result, each relevant fragment of a policy has been mapped to suitable ontology categories and attributes, generating a knowledge base of about 23,000 annotated data practices. Each fragment may be associated with different categories and attributes, on the basis of interpretations by multiple annotators. In this regard, consolidating alternative and potentially conflicting interpretations is a relevant challenge for our work, which we are currently addressing using natural language processing and machine learning techniques.

To the extent that contradictions have a logical nature, state-of-the-art inference engines like Pellet [28] would be sufficient to flag them. Preliminary results show that there's complete agreement between annotators on whether the *Do Not Track* data practice is 'honored' or 'not honored' by a given policy: but in cases when those two mutually exclusive values are selected for the same fragment, automatic reasoning with *PrivOnto* should detect the inconsistency, letting regulators and users make a final decision by reviewing the considered fragment. We plan to conduct comprehensive experiments to properly assess the impact of logical inconsistencies in the annotated corpus.

We also plan to evaluate the feasibility of using the Semantic Web Rule Language (SWRL<sup>9</sup>) to apply DL-safe consolidation rules between alternative annotations. For instance, in the eventuality that an annotator selects 'cookies and tracking elements' as value of 'personal information type' attribute associated to a fragment, and a different annotator selects 'location' for the same fragment and data practice, one way to consolidate these two annotations would be by applying a conservative rule, merging the two values into 'location AND cookies and tracking elements'. This rule-based consolidation process requires domain expert support on identifying salient rules, similarly to the process of building the library of 57 SPARQL queries.

Semantically-labeled privacy policies constitute an important resource for privacy analysts and regulators,

<sup>9</sup><https://www.w3.org/Submission/SWRL/>

but scaling the process of annotating natural language privacy policies accordingly can be challenging. As part of the efforts in the UPP project, we investigate the potential of crowdsourcing privacy policy analysis to non-experts, in combination with machine learning, in order to enable semi- or fully automated extraction of data practices and their attributes from privacy policy documents [2,4,40]. These efforts show promise for scaling up our analysis, which would enable further expansion of *PrivOnto*'s knowledge base.

*PrivOnto* can supply privacy researchers and regulators with tangible evidence of the effectiveness of semantic technologies, helping reduce the complexity of policies, and bypass their convoluted language, while retaining the ability to analyze semantic reasoning results in conjunction with textual evidence from the policy. More compact and transparent high-level representations of policy contents also benefit end-users: in this regard, we plan to integrate query-driven search functionality into the UPP project's data exploration portal.<sup>10</sup> This portal already visually integrates the data practice annotations with a privacy policy's original text in an easy-to-use web interface (see Figure 7), and enables users to filter for attributes and values of specific frame categories, although currently in a limited manner without the support of semantic technologies. The notion of "sociotechnical system" [8], according to which the interaction between people and technology is a central aspect of our society, can be a useful paradigm to understand privacy in the Digital Era. In this respect, semantic technologies (STs) are powerful tools for investigating how the human and machine factors of privacy interconnect: as shown in this article, STs provide computational semantic models of privacy policies and requirements, and enable transparent information-access, contributing to users and regulators' contextual awareness. But the complexity of the privacy domain, along with the technological and legal implications that underlie privacy-enabling solutions, are far from being solved issues: on the contrary, they will likely constitute some of the most interesting problems that the semantic web community will have to face in the years to come.

## 6. Acknowledgments

This research has been partially funded by the National Science Foundation under grant agreements

CNS-1330596 and CNS-1330214. The authors would like to acknowledge the entire Usable Privacy Policy Project team for its dedicated work; and especially thank Pedro Giovanni Leon, Mads Schaarup Andersen, and Aswarth Dara for their contributions to the design and validation of the annotation scheme, as well as the corpus creation.

## References

- [1] Bartolini, C., Muthuri, R., and Santos C.: Using Ontologies to Model Data Protection Requirements in Workflows. In Intl. Workshop Juris-informatics (2015)
- [2] Bhatia, J., Breaux, T.D., and Schaub, F.: Mining Privacy Goals from Privacy Policies using Hybridized Task Re-composition. ACM TOSEM (forthcoming)
- [3] Breaux, T.D., Hibshi, H., and Rao, A.: Eddy, a formal language for specifying and analyzing data flow specifications for conflicting privacy requirements. Requirements Engineering, 19.3, 281-307 (2014)
- [4] Breaux, T.D., and Schaub, F.: Scaling requirements extraction to the crowd: Experiments with privacy policies. In Intl. Req. Eng. Conf., IEEE (2014)
- [5] Children's Online Privacy Protection Rule (COPPA). 16 CFR Part 312 (1998)
- [6] Cranor, L.: Web privacy with P3P. O'Reilly Media, Inc. (2002)
- [7] De Coi, J.L., and Olmedilla, D.: A review of trust management, security and privacy policy languages. In Int. Conf. Security and Cryptography (2008)
- [8] De Greene, K. B.: Sociotechnical systems: factors in analysis, design, and management. Prentice-Hall, (1973).
- [9] Duma, C., Herzog, A., and Shahmehri, N.: Privacy in the semantic web: What policy languages have to offer. POLICY '07: IEEE Int. Workshop Policies for Dist. Sys. Netw., 109-118 (2007)
- [10] Ashley, P., Hada, S., Karjoth, G., Powers, C., and Schunter, M.: Enterprise privacy authorization language (EPAL 1.2). Submission to W3C, 156. (2003)
- [11] Fillmore, C. J.: Frame semantics and the nature of language. Annals of the New York Academy of Sciences 280, no. 1, 20-32. (1976)
- [12] Federal Trade Commission: Privacy Online: Fair Information Practices in the Electronic Marketplace. Report to Congress (2000)
- [13] Cuenca Grau, B.: Privacy in ontology-based information systems: A pending matter. Semantic Web 1, 2, 137-141 (2010)
- [14] Gandon, F.L., and Sadeh, N.M.: Semantic web technologies to reconcile privacy and context awareness. Web Semantics: Science, Services and Agents on the World Wide Web. 241-260 (2004)
- [15] Google: Google Trends. Accessed: March 15, 2016.
- [16] Heurix, J., Zimmermann, P., Neubauer, T.: A taxonomy for privacy enhancing technologies. Computers & Security 53, 1-17 (2015)
- [17] HIPAA Privacy Rule. 45 CFR Part 160 (2002)
- [18] Hobbs, J. R. and Pan, F.: An Ontology of Time for the Semantic Web. ACM Transactions on Asian Language Processing (TALIP) Vol. 3, No. 1, pp. 66-85. (2004)

<sup>10</sup><https://explore.usableprivacy.org/>

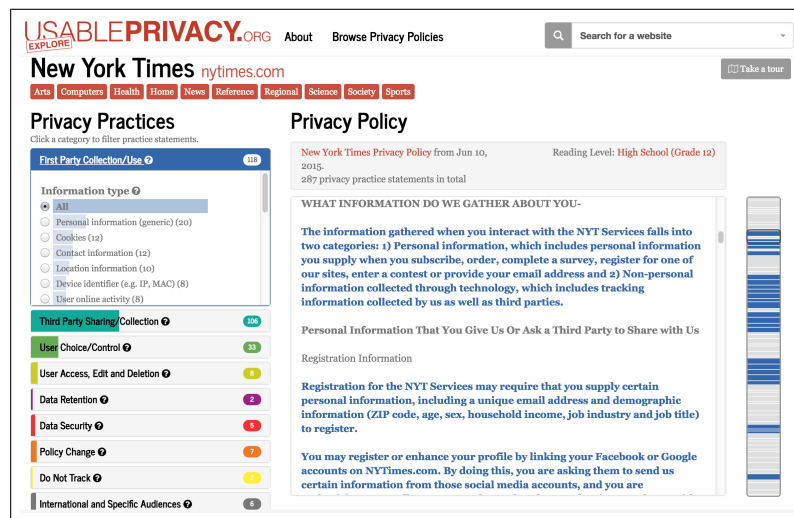


Fig. 7. A screenshot of the UPP Explore website that visualizes the First Party collection data practice of the New York Times' privacy policy.

- [19] Jutla, D.N., Bodorik, P., and Zhang, Y.: PeCAN: An architecture for users' privacy-aware electronic commerce contexts on the semantic web. *Information Systems* 31, 4-5 (2006)
- [20] Kagal, L., Berners-Lee, T., Connolly, D., and Weitzner, D.: Using semantic web technologies for policy management on the web. *Proceeding of the National Conference on Artificial Intelligence*, Vol. 21. No. 2. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press (1999)
- [21] Kaufmann, E., and Abraham B.: Evaluating the usability of natural language query languages and interfaces to Semantic Web knowledge bases. *Web Semantics, Science, Services and Agents on the World Wide Web* 8.4, 377-393 (2010)
- [22] Kost, M., Freytag, Christoph, Kargl, F., Kung, A.: Privacy verification using ontologies. *ARES '11: Int. Conf. Availability, Reliability and Security*, IEEE (2011)
- [23] Lorch, M., Proctor, S., Lepro, R., Kafura, D., and Shah, S.: First experiences using XACML for access control in distributed systems. In *ACM workshop on XML security* (2003)
- [24] McDonald, A. M., and Cranor, L.F.: The Cost of reading privacy policies. *I/S J Law & Policy Info. Soc.*, 4(3) (2008)
- [25] Negroponte, N.: *Being Digital* Random House Inc., New York, NY, USA (1995)
- [26] Niles, I. and Pease, A.: Origins of the IEEE standard upper ontology. In *Working notes of the IJCAI-2001 workshop on the IEEE standard upper ontology*, pp. 37-42. (2001).
- [27] Online Privacy Protection Act of 2003. *California Business and Professional Code*, 22575-22579 (2004)
- [28] Parsia, B., and Evren, S.: Pellet: An OWL DL reasoner. *Int. Semantic Web Conf.*, Poster (2004)
- [29] President's Council of Advisors on Science and Technology: *Big Data and Privacy: a Technological Perspective*. Executive Office of the President, USA (2014).
- [30] Reidenberg, J.R., Bhatia, J., Breaux, T.D., and Norton, T.B.: Automated Comparisons of Ambiguity in Privacy Policies and the Impact of Regulation. *J Legal Studies* 47 (forthcoming).
- [31] Reidenberg, J.R., Russell, N.C., Callen, A.J., Qasir, S., and Norton, T.B.: Privacy harms and the effectiveness of the notice and choice framework. *I/S: J. Law & Policy Info. Soc.* 11, 2 (2015)
- [32] Renars, L., Cerans, K., and Sprogis, A.: Visualizing and Editing Ontology Fragments with OWLGrEd. *I-SEMANTICS (Posters & Demos)* (2012)
- [33] Sackmann, S., and Kahmer, M.: ExPDT: A policy-based approach for automating compliance. *Wirtschaftsinformatik* 50.5, 366 (2008)
- [34] Sadeh, N., Gandon, F., and Oh Buyng, K.: Ambient Intelligence: The MyCampus Experience. In *Ambient Intelligence and Pervasive Computing*. Eds. T. Vasilakos and W. Pedrycz, ArTech House (2006)
- [35] Schaub, F., Balebako, R., Durity, A.L., and Cranor, L.F.: A Design Space for Effective Privacy Notices. In *Symposium on Usable Privacy and Security* (2015)
- [36] Shvaiko, P., Oltamari, A., Cuel, R., and Pozza, D.: Generating innovation with semantically enabled TasLab portal. *The Semantic Web: Research and Applications*. Springer Berlin Heidelberg, 348-363 (2010)
- [37] Toninelli, A., Montanari, R., Kagal, L., and Lassila, O.: Proteus: A Semantic Context-Aware Adaptive Policy Model. In *POLICY '07: Int. Workshop Policies Distr. Sys. Netw.*, IEEE (2007)
- [38] Tonti, G., Bradshaw, J.M., Jeffers, R., Montanari, R., Suri, N. and Uszok, A.: Semantic Web Languages for Policy Representation and Reasoning: A Comparison of KAoS, Rei, and Ponder. *Int. Semantic Web Conf.* (2003)
- [39] Uszok, A., et al.: KAoS policy and domain services: Toward a description-logic approach to policy representation, deconfliction, and enforcement. In *POLICY '03: Int. Workshop Policies Distr. Sys. Netw.*, IEEE (2003)
- [40] Wilson, S., Schaub, F., Ramanath, R., Sadeh, N., Liu, F., Smith, N.A., and Liu, F.: Crowdsourcing Annotations for Websites' Privacy Policies: Can It Really Work? In *WWW '16: Int. World Wide Web Conf.* (2016)