

PrivacyLens: A Framework to Collect and Analyze the Landscape of Past, Present, and Future Smart Device Privacy Policies

Aamir Hamid^{1,*}, Hemanth Reddy Samidi¹, Tim Finin¹, Primal Pappachan^{2,†}, Roberto Yus¹
¹University of Maryland, Baltimore County, ²Portland State University
 {ahamid2, finin, hsamidi1, ryus}@umbc.edu, primal@pdx.edu

Abstract

As the adoption of smart devices continues to permeate all aspects of our lives, concerns surrounding user privacy have become more pertinent than ever before. While privacy policies define the data management practices of their manufacturers, previous work has shown that they are rarely read and understood by users. Hence, automatic analysis of privacy policies has been shown to help provide users with appropriate insights. Previous research has extensively analyzed privacy policies of websites, e-commerce, and mobile applications, but privacy policies of smart devices, present some differences and specific challenges such as the difficulty to find and collect them. We present PrivacyLens, a novel framework for discovering and collecting past, present, and *future* smart device privacy policies and harnessing NLP and ML algorithms to analyze them. PrivacyLens is currently deployed, collecting, analyzing, and publishing insights about privacy policies to assist different stakeholders of smart devices, such as users, policy authors, and regulators. We show several examples of analytical tasks enabled by PrivacyLens, including comparisons of devices per type and manufacturing country, categorization of privacy policies, and impact of data regulations on data practices. At the time of submitting this paper, PrivacyLens had collected and analyzed more than 1,200 privacy policies for 7,300 smart devices.

1 Introduction

The Internet of Things (IoT) has gained rapid popularity in recent years. *Smart* IoT devices are now utilized in, among others, transportation, industrial processes, smart homes, and health care. Smart devices have advanced capabilities that include, in general, the collection of information, usage of cutting-edge technologies, such as Artificial Intelligence (AI) to process such data, and automation of tasks to provide personalized user experiences [1]. More than 40 million households in the US have adopted smart home devices, and this number is expected to reach 64.1 million by 2025 [2]. The

growing use of smart technology also poses privacy risks [3,4]. IoT devices collect large amounts of diverse data, and consumers sometimes do not understand what data is being collected in their environment [5]. Data collection from IoT devices could lead to unbounded profiling of customers by businesses or disruption of regular operations by malicious entities [6, 7]. As a consequence, consumers are concerned about the privacy risks of owning and using IoT devices [8].

Privacy policies have traditionally provided information on the data collected/used/shared by services such as e-commerce and mobile applications, and extensive analysis has been performed on them [9–11]. However, to the authors’ knowledge, privacy policies of IoT devices (e.g., smart home devices) have not received as much attention. Arguably the most important work is that of the Mozilla Privacy Not Included project [12] in which human analysts study smart devices w.r.t their privacy, security, and usage of AI. However, manual analysis is not scalable to the increasing number of IoT devices in the market. Mozilla reported that their human analysts spent 68,160 minutes (47 days) reading and analyzing privacy policies in 2022. Hence, automatic collection and analysis of IoT privacy policies is required.

To bridge the gap in the privacy policy landscape of smart devices, we introduce PrivacyLens, a framework for automatic IoT privacy policy collection, analysis, and publication of insights. As there does not exist a centralized repository of IoT devices and their privacy policies, in contrast with mobile applications (e.g., Google Play Store), PrivacyLens searches e-commerce websites (e.g., Amazon, Walmart) continuously for IoT devices from which it extracts metadata, such as the manufacturer. Then, it uses a technique to retrieve the manufacturer’s website and, from it, the device’s privacy policy. Using the Wayback Machine [13] which is a digital archive of the world wide web, PrivacyLens retrieves, when available, historical privacy policies for each IoT device to generate a longitudinal analysis of the changes due to, for instance, privacy regulation, such as General Data Protection Regulation (GDPR). Using natural language processing (NLP) and machine learning (ML) techniques, PrivacyLens determines

different features of each privacy policy, including their overall quality, readability, and ambiguity. Finally, PrivacyLens stores both the raw data and inferences made for each IoT privacy policy and publishes it online. In summary, the main contributions of PrivacyLens are as follows:

- Discovers IoT devices in e-commerce websites and extracts their current and past privacy policies. Our evaluation shows that PrivacyLens achieves high precision and recall in collecting IoT privacy policy text.
- Extracts insights about an IoT privacy policy, including privacy as well as readability features, using NLP and ML-based techniques. Our evaluation shows that PrivacyLens achieves good results in every feature extracted.
- Publishes all the information collected online to enable different stakeholders (e.g., IoT customers, researchers, regulators) to make informed decisions based on insights about IoT device privacy policies.

PrivacyLens has been deployed and has so far generated a comprehensive data set encompassing privacy policies, key features, relevant information, and ambiguity levels of more than 1,200 privacy policies for 7,300 smart devices. Additionally, the paper includes a sample study of a subset of the PrivacyLens data to understand the landscape of IoT device privacy policies.

The rest of the paper is structured as follows. Section 2 reviews the state of the art in privacy policy collection and analysis. Section 3 describes motivating use cases for PrivacyLens and overviews the framework’s architecture. Section 4 and Section 5 explain finding and collecting IoT privacy policies and extracting insights. Section 6 evaluates the framework’s performance. Section 7 analyzes a subset of the data generated by PrivacyLens. Section 8 concludes the paper.

2 Related Work

We review the state of the art in privacy policy collection and analysis. We note that the IoT domain is still fairly unexplored in this context since most prior work focuses on domains such as e-commerce and mobile applications.

Automated frameworks for privacy policy extraction and analysis. Several automated frameworks have been proposed to analyze and evaluate privacy policies from different domains. The approach in [14] builds upon prior work in NLP, privacy preference modelling, crowdsourcing, and privacy interface design. It extends existing research on user preference modeling in privacy policies and incorporates innovative approaches such as semi-automated feature extraction and privacy notice design based on extracted policy features. Polisis [15] uses a neural network hierarchy to extract high-level privacy practices and precise data from

the privacy policies of websites. Then, it offers an interface for both structured and free-form querying of privacy policies. PI-Extract [16] is a fully automated system to extract fine-grained personal data phrases and their corresponding practices from the privacy policies of websites. PI-Extract is based on a neural model which outperforms rule-based baselines in accurately extracting privacy practices. In [17], the authors address the challenge of users being asked to release personal information without full awareness of the data collection practices. They propose an automated solution that utilizes Information Extraction techniques to analyze privacy policy text and highlight the personal information being collected. In [18], the authors deal with the problem of extracting transparency information from website privacy policies by proposing a ‘Human-in-the-Loop’ approach that combines machine learning-generated suggestions with human annotation decisions. Their prototype system streamlines the annotation process by providing meaningful predictions to users, resulting in improved performance compared to other extraction models for legal documents.

Kuznetsov et al. [19] present the only other framework focused on collecting privacy policies from IoT devices in the literature. As PrivacyLens, their approach collects information about IoT devices from e-commerce and finds their privacy policies online. Their system collected 592 distinct privacy policies from various IoT device manufacturers. Their analysis conducts a detailed statistical and semantic analysis to improve policy transparency. In contrast with our approach, we do not limit ourselves to current privacy policies; instead, we also collect and analyze past versions, and since the framework is continuously running, we will also collect future versions. This allows us to track the evolution of these policies over time and assess the impact of changes in regulatory guidelines and data collection practices. Additionally, the analysis performed by PrivacyLens extends also that of [19] by employing NLP and ML techniques to conduct a more in-depth examination (including extracting privacy insights and other textual features such as readability and ambiguity).

Annotated privacy policy datasets. Prior work has published privacy policy datasets pertaining to policies from mobile applications, web services, and IoT devices. The OPP-115 dataset [20] contains 115 website privacy policies collected using Amazon Alexa [21]. The dataset includes annotations and labels which provide examples of how personal data are used and details on the experts who annotated the texts. Another annotated dataset is the APP-350 corpus [22] which includes more than a million privacy policies for Android applications available on the Google Play store. Amos et al. [23] curated a longitudinal data set of more than a million privacy policies, exposing troubling trends in transparency and accessibility. For example, policies have doubled in length and increased in complexity over two decades and often fail to disclose common tracking technologies and third parties.

None of the previous works focuses on IoT/smart devices, hence, the only other dataset in this domain, up to the author’s knowledge, is the one containing 592 policies described before [19]. In contrast, the dataset extracted from PrivacyLens, at the time of writing this paper, contains 1,200 policies from more than 7,300 smart devices.

3 Framework Overview

PrivacyLens automatically and periodically collects and analyzes the privacy policies of IoT devices, including products such as smart home devices and wearable technology. In this section, we first present motivating use case scenarios supported by PrivacyLens. Then, we overview the high-level architecture of the framework and its main components.

Motivating use cases. PrivacyLens can enable different stakeholders to perform further analysis, such as:

- *IoT customers* can understand how their personal data will be collected, used, and protected by the smart device. Using PrivacyLens, they can evaluate whether the device functionalities align with their privacy preferences and determine whether the privacy protection mechanisms implemented are sufficient. The comparison of insights for different devices of the same category (e.g., smart watches), or even of the same manufacturer (e.g., Fitbit) over the years, empowers customers to choose devices from manufacturers that prioritize privacy and security.
- *Product privacy lawyers* can compare their product’s policy and ensure that it meets the industry best practices/standards and expectations of consumers. PrivacyLens allows them to examine the strengths and weaknesses of competitor policies, identify areas for improvement or differentiation, and adopt better privacy protection mechanisms or create a privacy policy that can more easily be understood by the customers.
- *Data protection regulators* can evaluate the influence of data protection regulations on various smart device manufacturers. PrivacyLens offers insights into the changes instigated by these laws, shedding light on whether companies adapted their policies, as well as the timeline involved in such adaptations. By filling in these gaps, data protection regulators can also understand the trends in privacy policy updates, enabling more informed decision-making for future regulations.

Figure 1 shows an overview of the PrivacyLens framework that can support these use cases, as well as others. The system and the data it generated can be accessed at <https://privacy-lens.web.app/home>¹. The framework

¹The source code of the system is available at *URL removed for anonymization of the manuscript*.

is structured around three stages:

Collection. PrivacyLens finds and collects IoT device privacy policies from prominent e-commerce platforms (such as Walmart and Amazon), processes the documents extracted, and stores the information in a database. The input for this task is a list of 24 categories of smart IoT devices (e.g., smart cameras, lights, TVs, etc.) and the output is a cleaned and structured database of IoT device privacy policies. This output includes both current and past policies (extracted from a web archive) to enable the longitudinal study of IoT device privacy. We describe the methodology used to collect privacy policies in detail in Section 4.

Analysis. PrivacyLens enables a comprehensive analysis of the database of IoT device privacy policies, focusing on two key aspects: textual analysis and privacy analysis. The goal of this stage is to generate metadata that enables different stakeholders to gain deeper insights into the privacy landscape of smart devices. The input is the cleaned and structured privacy policy data from the collection stage, and the output consists of metadata and in-depth analysis of similarity, ambiguity, keywords, and privacy insights. Further details on the analysis stage can be found in Section 5. The results of these analyses are presented using a variety of exploratory tools, allowing stakeholders to synthesize complex data, draw meaningful conclusions, and develop targeted strategies to address privacy challenges more effectively.

Publication. The results from the analysis stage are published on a dedicated website, which is regularly updated every month. This update frequency ensures that users are always presented with the most current and accurate data available. In addition to the metadata, the raw data collected is also published for the purpose of transparency. 1 includes a representative screenshot of our website, showcasing its user-friendly interface and various sections, including but not limited to raw data, analysis results, and weekly updates. By making this wealth of information accessible and actionable, PrivacyLens contributes to the fostering of a more transparent and responsible digital ecosystem.

4 Policy Collection

IoT device information extraction. We designed and implemented a multi-threaded web scraper specializing in extracting IoT device data from e-commerce platforms, i.e., Amazon and Walmart. The web scraper code establishes a link, using *WebDriverManager*², with a Firefox browser to extract the data. Then, it uses *Selenium*³ to control the web browser from the code and search IoT devices using as a

²<https://bonigarcia.dev/webdrivermanager/>

³<https://www.selenium.dev>

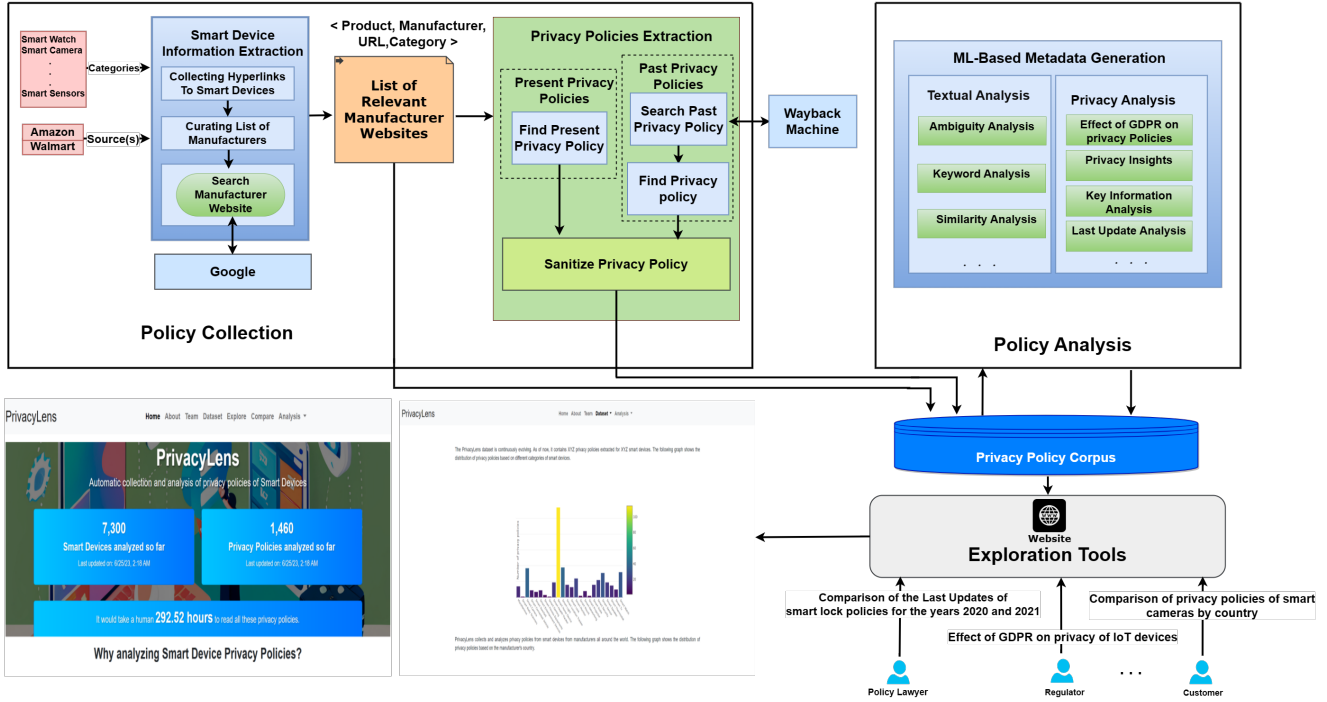


Figure 1: A high-level overview of PrivacyLens.

query, a list of relevant categories of smart devices. The list includes the following keywords with “smart” as a prefix: body scanner, camera, connected vehicle, doorbell, entertainment device, fitness equipment, gaming technology, health tracker, home device, light, location tracker, lock, monitor, mount, networking device, projector, scale, security system, sensor, speaker, thermostat, TV, and watch.

Based on the category, the system constructs a URL specific to the e-commerce site for web parsing. This enables the corresponding web page to appear when that particular category is searched in the e-commerce site’s search bar. Then PrivacyLens parses that web page using *BeautifulSoup*⁴ to obtain the product URLs by finding certain “href” element. This process will be executed for each input category. PrivacyLens then navigates to individual product URLs using the web driver to collect the manufacturer’s name from the HTML markup. This information is gathered by finding a certain class attribute value in the markup (DOM structure of the web page) that is consistently marked for rendering the manufacturer’s name on that web page. At the end of this step, the parser returns the e-commerce URL, name of the product, and manufacturer for each product.

Manufacturer’s website extraction. PrivacyLens executes a search query on the Google search engine using a composite term (manufacturer’s name and device type). Next, the

scraper analyzes the HTML markup of the first page of results obtained to find all URLs. For each result, PrivacyLens calculates a value that represents the likelihood that the URL is the official manufacturer’s website. This is done by comparing the manufacturer’s name to the domain using the Ratcliff/Obershelp algorithm [24]. Then, PrivacyLens selects the URL with the highest value as long as it scores higher than a threshold (0.8) which we determined experimentally.

Privacy policy extraction and cleaning. After manually analyzing dozens of such websites we note that, in general, sites have a hyperlink in their main page that links to the privacy policy and is labelled “privacy”, “privacy policy”, “privacy-statement”, “privacy-notice” or “privacy-policy”. PrivacyLens incorporates another scraper that searches such keywords within the HTML code of the website to identify the URL of the privacy policy. Then, it fetches its HTML, parses it, and cleans it to remove irrelevant HTML tags and website content.

The cleaning task uses regular expressions to remove elements from the HTML document such as navigation bars, headers, and footers. Then, it uses the Bleach library⁵ for further sanitization to remove tags and attributes. Finally, PrivacyLens takes care of formatting to make the result more readable by adding appropriate indentation, eliminating unnecessary spaces including any leading or trailing spaces.

⁴<https://www.crummy.com/software/BeautifulSoup/>

⁵<https://pypi.org/project/bleach>

Past privacy policy extraction. To collect the past privacy policies, PrivacyLens uses the Wayback Machine [13], an internet archive containing older snapshots of websites. For each manufacturer URL obtained in the previous step, PrivacyLens queries the Wayback Machine, using its API which takes a year and a URL as parameters, to obtain snapshots at different points in time⁶. PrivacyLens executes one query per year for the last 10 years as well as additional queries for the months before/after the moment when important data privacy regulation (such as the General Data Protection Regulation (GDPR) [25] and the California Consumer Privacy Act (CCPA) [26]) went into effect. Once the different snapshots are retrieved, PrivacyLens uses the process explained in *Privacy policy extraction and cleaning* to find the privacy policy within the website and clean it.

5 Privacy Policy Analysis

PrivacyLens incorporates a suite of NLP and ML techniques to perform different types of analyses and extract various insights from a privacy policy. The goal of this stage is to determine the readability of the policy document and evaluate its acceptability with respect to privacy.

5.1 Privacy Analysis

PrivacyLens extracts the following four privacy features of an IoT device policy.

Keyword Usage. We employed a policy annotation approach [27] to capture the data practices described in privacy policies. The annotation scheme, in addition to information from previous works [28], was used to classify the keywords derived from each privacy policy into ten clusters which represent concepts related to the management of data by the device/manufacture such as user choice or third party sharing (see Table 1). PrivacyLens counts the number of keywords related to each category found in each paragraph of the privacy policy.

Last Update. Policy updates are crucial to keeping users informed about changes in the handling of their data. PrivacyLens identifies explicit mentions to the date in which the privacy policy was updated.

Regulations Compliance. Over 150 countries around the world have adopted new data protection regulations [29]. An automated cross-country and manufacturer-wise assessment of regulatory compliance is essential for analyzing the policy landscape. PrivacyLens employs a deep learning approach,

⁶We use the manufacturer URL instead of the privacy policy URL since we observed that in a significant amount of websites, the privacy policy URL changed over the years.

utilizing the Bidirectional Encoder Representations of Transformers (BERT) model [30], fine-tuned on the Semantic Textual Similarity Benchmark dataset (STS-B) [31], to evaluate the changes in IoT device privacy policies before and after specific regulations (e.g., GDPR) came into effect. This approach generates a semantic similarity score between policy pairs (i.e., the policy before and after certain event), providing a numeric gauge of the policy changes.

Overall Assessment. PrivacyLens summarizes the results of its privacy analysis of an IoT policy into an overall assessment (i.e., *acceptable* or *unacceptable*) inspired by the Mozilla Privacy Not Included (PNI) project [12]. To this end, we created a dataset consisting of the 172 IoT devices analyzed by PNI (at the time of writing this paper). We partitioned the dataset into an 80% training set (138 policies) and a 20% testing set (34 policies). Notably, the dataset was imbalanced, with "acceptable" instances numbering 135 and "unacceptable" ones only 17. To address this, we employed random oversampling, effectively balancing the "unacceptable" class. We use the logistic regression classifier which estimates the probability of an outcome based on independent variables. This model is particularly effective for categorical target variables, such as policy categorization [32]. We considered 19 features that include the ones described above along with the features about readability that we will explain next. To optimize the model further, we used GridSearchCV from the *sci-kit-learn*⁷ library for hyperparameter tuning, ultimately enhancing the performance of our logistic regression model.

5.2 Readability Analysis

PrivacyLens extracts the following eight features as a result of an analysis of the readability of the policy.

Entropy. In the realm of privacy policy analysis, entropy serves as a critical metric, offering a quantitative measure of the policy's textual uncertainty and complexity, thereby shedding light on interpretive challenges that users may face. The entropy of a language is a statistical quantity that measures how much information is generated on average for each letter in a text in that language. Hence, it defines the uncertainty or disorder in a text document. Shannon [33] presented a technique to calculate the entropy of English text using the following equation:

$$H(X) = - \sum_{i=1}^n P_i \log_2(P_i) \quad (1)$$

where $H(X)$ represents the entropy of a discrete random variable X and P_i is the probability of the i th outcome of the random variable X . PrivacyLens computes the entropy of a

⁷<https://scikit-learn.org/>

Privacy Attribute	Definition	Keywords
First Party Collection	The methods and purposes used by a service provider to get user data.	collect, gather, use, we collect
Third Party Sharing	The methods used by third parties to share or acquire user information.	third party, third parties, third party sharing, third party collection
Access, Edit, & Deletion	If users may access, edit, or remove their information, and how.	access, edit, delete, modify, revise
Data Security	How user data is safeguarded.	security, secure, safety, protect, infosec
Policy Change	Whether and how users will be informed of privacy policy changes.	change, modify, policy change, policy modification
Do Not Track	Whether and how internet tracking and advertising using Do Not Track signals are handled.	dnt, do not track
Opt-Out	The user's ability to choose not to participate in certain online activities, such as internet tracking and advertising.	opt out, opt-out , optout
Legislation	The legal frameworks that empower individuals to control the collection, usage, and distribution of their personal information by businesses and organizations.	GDPR, CCPA, General Data Protection Regulation
User Choice	User's right to make decisions about how their personal data is collected, used, and shared by a service or platform.	correct, review, change, update
Data	Information collected about users by a company.	identifier, name, email, address, IP address, number, biometric, activity, sleep, geolocation, location, GPS, photo, friends, voice, video

Table 1: Privacy attributes extracted from a privacy policy based on keywords.

document based on Equation 1 by iterating over each word and noting its frequency of occurrence. Then, it divides the frequency by the total number of words to get an estimate of the probability of each word. Next, it calculates the average length of each word in bits by multiplying its probability by the negative logarithm (base 2) of that same probability. The entropy of the document is then the sum of these calculated values for all words.

Reading Time. Privacy policies are hard to read and therefore do not help customers make informed decisions due to the fact that they are very lengthy and time-consuming [34]. Reading time is calculated using an individual's typical reading pace (roughly 238 WPM) [35]. PrivacyLens counts all the words in the document and divides the total by 238 to compute the estimated number of minutes that it would take a person to read the full privacy policy.

Unique Words. Privacy policies tend to use technical jargon to convey data usage and control to consumers. Unique, low-frequency words are crucial to understanding these documents, as they provide key context and learning aspects [36]. Despite their scarcity, they make comprehension challenging, as understanding the content fully requires a strong grasp of this specialized vocabulary.

The process begins by converting the text to lowercase and eliminating stop words, punctuation, and numbers followed by standardized through tokenization [37] using the Spacy library [38]. Following standardization, the PrivacyLens proceeds to identify unique words, which are words that have distinct character sequences. These unique words are counted, providing a distinct vocabulary size. Furthermore, the system calculates the ratio of the unique word count to the total word count in the document, offering a quantitative measure of the text's lexical diversity.

Coherence Score. Topic modeling employs a coherence score to measure how well a topic is understood by people, as established by Syed et al. [39]. It evaluates word similarity within a topic based on their frequency in a document. The Latent Dirichlet Allocation (LDA) technique, a type of unsupervised machine learning, aids in text analysis by identifying the best topics to represent the data (Yu et al., 2001) [40]. In this approach, a Dirichlet distribution is created first for documents in the subject space, and topics and words are selected from multinomial distributions. The coherence score, calculated as the sum of scores between every pair of words, is used to gauge the quality of the topics learned. The measure used in this case is CV, which computes scores via cosine similarity and normalized pointwise mutual information (NPMI) based on word co-occurrences. Then, the overall coherence of the privacy policy is computed as $\sum_{i < j} \text{score}(w_i, w_j)$, Where w_i and w_j represent words at positions i and j respectively within a given text. In the context of readability, a high coherence score suggests a well-structured and clear flow of ideas, thereby increasing readability, while a low coherence score often indicates a disjointed or unclear progression of thoughts, potentially making the text more difficult to understand.

Frequency of Imprecise Words. Imprecise words, such as "commonly" or "normally" can create ambiguity, making it hard to understand a service provider's operations. PrivacyLens employs NLTK to tokenize the text and regex to count the frequency of imprecise words to measure their prevalence in the privacy policy. Table 10 shows the list of imprecise words considered.

Connective Words Frequency. While connective words (such as "and" or "then") are useful to create coherent sentences, their overuse can make the text complex. PrivacyLens,

in a similar way than for the previous feature, counts the frequency of connective words using the list of words in Table 11.

Grammatical Errors. The integrity of a work depends on proper grammar, much as it does on word spelling [41]. PrivacyLens takes privacy document as input and then uses the NLTK library for tokenization (i.e., breaking the text into sentences) and the `language_tool_python` library to check for grammatical errors. It counts the total number of sentences and the number of sentences that contain at least one mistake providing a measure of the grammatical correctness of the input text.

Readability. Readability signifies how easily a text, like a policy, can be understood, based on its vocabulary, syntax, and sentence structure. Various readability tests exist [42], devised by linguists, each considering different text aspects. In our study, we employed the Flesch-Kincaid Grade Level [43], which presents the score as a U.S. grade level. This metric represents the educational level required to understand a text and is computed using the following formula:

$$FKGL = 0.39 \left(\frac{\text{totalwords}}{\text{totalsentences}} \right) + 11.8 \left(\frac{\text{totalsyllables}}{\text{totalwords}} \right) - 15.59 \quad (2)$$

In the Flesch-Kincaid Grade Level (FKGL) formula, 0.39 and 11.8 are weights for the average sentence length and syllables per word respectively. The constant -15.59 calibrates the score to U.S. grade levels.

Ambiguity. While privacy policies should clearly state the data handling practices, privacy policies often contain ambiguous language [44]. PrivacyLens incorporates the supervised learning approach to classify a privacy policy based on a scale with three ambiguity levels (*not ambiguous*, *somewhat ambiguous*, and *very ambiguous*) presented in [45]. We annotated 100 IoT device policies extracted from PrivacyLens and trained a random forest classifier [46] and a logistic regression [32] classifier. PrivacyLens applies both classifiers to each privacy policy and stores their output labels.

6 Framework Evaluation

In this section, we assess the effectiveness and performance of PrivacyLens by conducting an evaluation of both its policy collection and analysis capabilities.

6.1 Evaluating Policy Collection

Current Privacy Policy Extraction. We assessed the quality of our web-scraped data from Amazon against manually validated “truth” values for each of the top 30 records in ten categories of smart devices (i.e., 300 total devices). Our

initial evaluation is focused on two key aspects: extraction of manufacturers and websites of manufacturers. We used three metrics in our evaluation: accuracy, recall, and precision. Table 2 shows the results of the experiment. For the manufacturer data, we observed a high overall F1 score (0.98), recall (0.96), and precision (0.99). For the collection of the website URL, PrivacyLens achieves also a high overall F1 score (0.95), recall (0.91), and precision (0.99). These results highlight that PrivacyLens can accurately obtain information about IoT device manufacturers and their websites, which is required to find their privacy policies.

Category	Manufacturer Collection			Website Collection		
	Recall	Prec.	F1	Recall	Prec.	F1
Sensor	0.90	1.00	0.95	0.95	1.00	0.98
Projector	1.00	1.00	1.00	0.94	0.94	0.94
Bulb	1.00	1.00	1.00	1.00	1.00	1.00
Speaker	0.97	1.00	0.98	0.97	1.00	0.98
Alarm	0.97	1.00	0.98	0.89	1.00	0.94
Camera	0.97	1.00	0.98	0.85	1.00	0.92
Scale	0.90	1.00	0.95	0.92	1.00	0.96
Watch	1.00	1.00	1.00	1.00	1.00	1.00
Lock	0.97	1.00	0.98	0.88	1.00	0.93
Tracking	0.97	0.97	0.97	0.74	1.00	0.85
Overall	0.96	0.99	0.98	0.91	0.99	0.95

Table 2: Evaluation of the extraction of manufacturers (Manufacturer Collection) and their websites (Website Collection) for IoT devices.

To evaluate the collection of privacy policies by PrivacyLens we randomly selected 100 extracted manufacturer websites from the previous step. Then, we manually identify the policy URL on each website and compare it to the policy URL automatically extracted by PrivacyLens. The results show that PrivacyLens exhibited strong performance, achieving an F1 score of 0.83 for locating URLs linked to privacy policies, and a 0.76 F1 score for extracting the policy text from the web pages associated with the extracted URLs. This indicates that the system is highly precise and accurate in its extraction capabilities. PrivacyLens encountered difficulties when the privacy policy was either present in another child component that renders on performing an event, or the website had a special download option for the privacy policy, which would require specialized parsers for certain websites and would not scale. These results show that the web parsing technique in PrivacyLens is a feasible method for collecting privacy policies.

Past Privacy Policy Extraction. In this experiment, we selected randomly another 100 manufacturers (and their associated websites) from the previous step making sure that the companies existed on or before 2020. Then, for each of them, we randomly select a date between 2020 and 2022 and use PrivacyLens to automatically retrieve the archived version of the manufacturer’s website on that date (if available). Our results indicate that PrivacyLens was able to retrieve archived

websites for 64% of them. After manually analyzing these retrieved websites, we observed that 8% of them pointed to empty homepages. This can be attributed to the snapshots taken by the Wayback Machine, which sometimes do not capture any information at all. Delving further into the results, PrivacyLens was successful in finding the archived privacy policy on the archived website of 30 out of 64 (47%) manufacturers. For the remaining manufacturers (44%), the failures were due to the absence of a privacy link in the snapshot. The average distance between the queried snapshot date and the retrieved snapshot date was 87 days.

Privacy Policy Cleaning. We evaluate PrivacyLens’ cleaning performance by manually cleaning 10 policies and automatically comparing the number of tokens that can be extracted from both automatic and manually cleaned policies. We use the *punkt* tokenizer model from NLTK [47] to break the policy text into individual tokens for comparison. The comparison results show that PrivacyLens’s automatic cleaning achieves high precision (0.94), recall (0.97), and F1-score (0.95), demonstrating its effectiveness in retaining essential information while removing unnecessary content from the privacy policy.

6.2 Evaluating Policy Analysis

To evaluate the insights extracted from each privacy policy, we leveraged existing ground truth datasets from the literature. As there does not exist a human-annotated dataset of privacy policies for IoT devices, we choose datasets of privacy policies of websites that include annotations. The dataset in [45], which is based on the OPP-115 dataset [20], contains annotations on grammatical errors, frequency of imprecise words, and connectivity words. Then, for the remaining features of PrivacyLens’s readability analysis, we annotate 10 policies of the OPP-115 using popular and free web-based resources: *Readable*⁸ for readability, *Planetcalc*⁹ for entropy, and *The Read Time*¹⁰ for reading time. We show the results next for each of the features.

W.r.t. **grammatical errors**, We followed the benchmark study’s [45] methodology for our comparative analysis. Despite the benchmark’s range of results, we were able to compare our specific scores effectively. Our results for the policy with the least grammatical errors nearly aligned with their lowest range, with a negligible difference of -0.06. However, for the policy with the most errors, we observed a higher difference, exceeding their range by about +0.7. This indicates our grammatical error detection tool, *language_tool_python*, has a stringent approach toward grammatical correctness. For **connective words** We compared our results with the general ranges provided in study [45]. For the policy with the least

connective words, the deviation is minuscule at 0.005, aligning our results closely with their lower range. With the policy utilizing the most connective words, our figures slightly exceed theirs by 0.016, still affirming the consistency between our method and the reported study. For **imprecise words** we also compared the use of imprecise words in policies with the ranges from study [45]. In the policy with the most imprecise words, our results were slightly higher, with a small difference of 0.05. For the policy with the fewest imprecise words, our results were a bit lower, with a difference of 0.09, still affirming the consistency between our method and the reported study.

Our Flesch-Kincaid **readability** scores closely matched those from Readable, indicating agreement between the methods. For instance, RedOrbit’s privacy policy, our score was 10.26 compared to Readable’s 10.4, and for sci-news policy, our score was 9.4 while Readable’s was 8.9. Scores for Uptodate and Earthkam were also similar, with our approach producing scores of 12.9 and 14.6, compared to Readable’s 12.9 and 15.5, respectively. Overall, our results were consistent with those from readable.com, affirming the reliability of our approach. PrivacyLens’s measure **entropy** between 4.1-4.2, shows a high consistency with Planetcalc’s range of 4.1-4.3. Specifically, Redorbit’s privacy policy scored 4.1 with our method, versus 4.2 with Planetcalc. For Earthkam, Uptodate, and Amazon, both our method and Planetcalc agreed on scores of 4.2, with the exception of Amazon, where Planetcalc reported a slightly higher score of 4.3. This close alignment affirms the precision and reliability of our method. In evaluating ten policies, our method calculated a total **reading time** of 120 minutes and 5 seconds, closely matching The Read Time’s estimate of 120 minutes and 40 seconds. The small variance of 120 seconds between both methods underscores the precision of our evaluation technique. The **unique words**, **keyword usage**, and **last update** features are based on searching specific keywords in the document and the search function is a well-tested Python library.

Next, we evaluate the performance of the **ambiguity** feature. The results in Table 3 show that the classification of policy texts is significantly more challenging for more ambiguous policies. We observed a decrease in the F1 score from 86% to 67% as the ambiguity of the policy increased. Additionally, in the case of logistic regression, we observed that the F1 score dropped from 77% to 69% when moving from non-ambiguous to ambiguous policies. These findings indicate that classification algorithms have less accurate results for highly ambiguous policies. This can be attributed to the fact that highly ambiguous policies have less well-defined requirements and guidelines, making it harder for classification algorithms to accurately assess them.

Finally, we evaluate the performance of the classifier trained to predict a human analyst **overall assessment** using the subset of 172 privacy policies corresponding to devices analyzed by the Mozilla PNI initiative (with their corresponding

⁸<https://readable.com>

⁹<https://planetcalc.com>

¹⁰<https://thereadtime.com/>

Ambiguity Class	Number of Policies	Random Forest Classifier	Logistic Regression
Not Ambiguous	283	0.86	0.77
Somewhat Ambiguous	90	0.67	0.72
Very Ambiguous	89	0.67	0.69

Table 3: F1-score of ambiguity determination models.

human analyst assessment). We applied cross-validation for a reliable estimate of the model’s effectiveness. The results (see Table 4) show that PrivacyLens successfully predicts the human analyst assessment with a high F1-Score (0.91). PrivacyLens performs better for the “acceptable” label than for the “unacceptable” label. This is due to the scarcity of unacceptable instances in the dataset which leads to low prediction accuracy, due to model bias and difficulty distinguishing this minority class despite using random oversampling. Figure 2 shows the importance of analysis of the different features used to train the model. Positive feature importance coefficient means the model’s prediction accuracy increases with the feature value, while a negative one means the prediction accuracy decreases, assuming other features are constant. The results show that, in general, the features related to the privacy analysis tend to have a higher impact on the prediction than those related to readability. However, features like reading level, unique words, and reading time, score high importance which might indicate that the analyst’s assessment is influenced by how much effort it requires to understand the policy. We would also like to highlight that PrivacyLens, matching the manual analysis of Mozilla PNI that took thousands of hours, showcases its proficiency in analyzing IoT device privacy policies.

Class	Precision	Recall	F1-score
acceptable	0.96	0.96	0.96
unacceptable	0.67	0.67	0.67
Weighted Avg	0.92	0.90	0.91

Table 4: Performance of the "overall assessment" model.

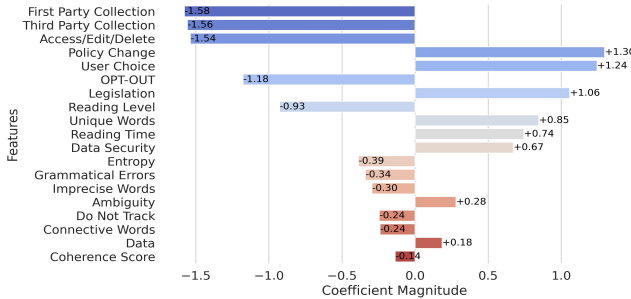


Figure 2: Feature Importance of the "overall assessment" model.

7 Analysis Supported by PrivacyLens

In this section, we present a study of a subset of privacy policies of IoT devices collected and analyzed by PrivacyLens and based on the motivating use cases described in Section 3.

7.1 Overview of the Dataset

The dataset considered in this study contains 462 IoT devices from 462 different IoT device manufacturers (see Table 5). We manually identified the country of origin of each manufacturer and added that information to the dataset. The USA was the country with the highest number of manufacturers in the dataset with a total of 251 devices (ranging from 0.4% of the devices being smartwatches and 31% being smart home devices). China was second with 46 devices (ranging from 2% of the devices being smart scales and 23% being smart home devices). Germany, Japan, France, and the United Kingdom also contributed significantly to the number of devices in the dataset.

Device Type	#	Device Type	#
Miscellaneous	15	Smart Camera	37
Smart Body Scanners	1	Smart Connected Vehicle	10
Smart Doorbell	8	Smart Entertainment Devices	10
Smart Fitness Equip.	3	Smart Gaming	1
Smart Health Tracker	18	Smart Home Device	117
Smart Light	37	Smart Location Tracker	15
Smart Lock	13	Smart Monitor	23
Smart Mount	1	Smart Networking	8
Smart Projector	2	Smart Scale	16
Smart Security	22	Smart Sensor	31
Smart Speaker	19	Smart Thermostat	15
Smart TV	10	Smart Watch	30

Table 5: Distribution of IoT device policies collected and analyzed in our study.

Additionally, we manually analyzed each of the privacy policies to check whether the privacy policy makes an explicit mention to the IoT device. 254 policies in the dataset did not make any explicit mentions of the IoT device while 208 policies explicitly mentioned that the privacy policy applied to the data collected by the device. This finding highlights that the state of IoT device privacy policies today lacks transparency. We further analyze this aspect by considering both the type of device as well as the country of origin of the manufacturer.

Figure 3 shows the distribution of policies with/without explicit mention to the smart device by country. We observe that for a significant amount countries, the number of policies without explicit mention is higher while only 12 countries had a higher number of policies with explicit mention. Note also that the majority of those countries contained only one device/manufacturer in our dataset (e.g., Finland, India, Jordan, Poland, Scotland, Singapore, Ireland, Vietnam) hence, this might not apply to a larger dataset. Several European countries, where the GDPR was enacted in 2018, have at least 10 devices in our dataset and also a higher number of policies with explicit mention of smart devices (e.g., France, Germany). Figure 4 shows the distribution of policies with/without explicit mention to the smart device by type of device. The large number of policies that do not explicitly mention the smart device is concerning due to the sensitive nature of the data collected by some of these devices such as smart cameras

and speakers. We observe that for two sensitive devices (i.e., smart health trackers and location trackers) the number of smart location trackers with explicit mention to the device vs. no explicit mention is much higher than for other categories.

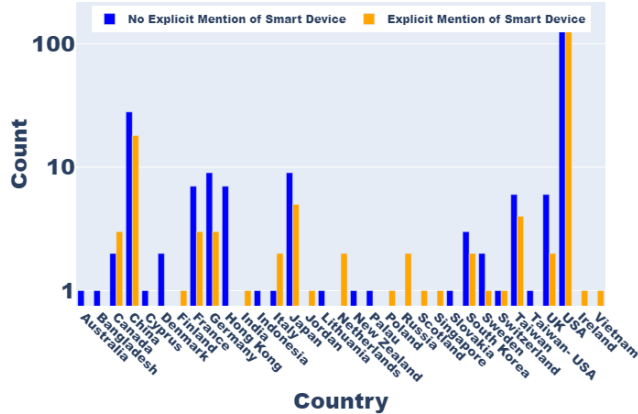


Figure 3: Distribution of policies based on manufacturer's country.

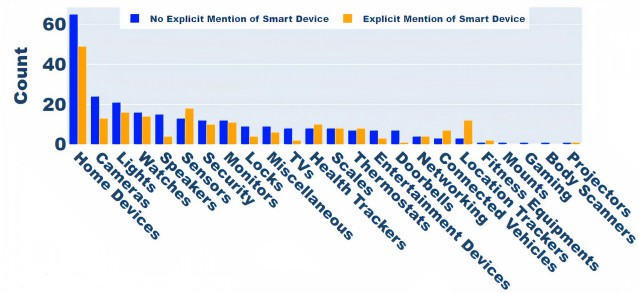


Figure 4: Distribution of policies based on the type of smart device.

Next, we analyzed how updated IoT device privacy policies are based on our dataset (see Figure 5). We observe that 62% of the policies (288 policies) disclosed their last update date, a concern considering how frequently IoT devices and their firm-wares are updated. When we separate policies based on whether there is an explicit mention of the IoT device or not, we observe that policies with explicit mention tend to be more updated for the same category. For instance, for the smart home device category, we note that six policies without an explicit mention of the device had a last update before 2013. For nearly 88% of the policies with an explicit update, the update occurred after the GDPR [48] became effective in 2018 (61% were updated after CCPA [26] became effective in 2020). This suggests that manufacturers are acknowledging and adopting the new requirements introduced by these regulations. Overall, this emphasizes the need for regular policy updates that comply with the latest regulations to ensure proper handling and protection of user data.

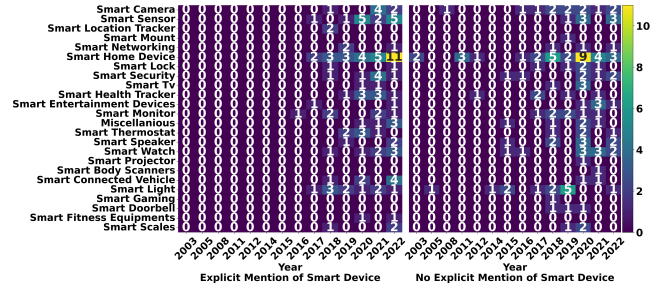


Figure 5: Distribution of privacy policies based on their last update.

7.2 PrivacyLens for IoT Customers

In this section, we will describe a potential analysis enabled by PrivacyLens for a user on the market for an IoT device. In particular, let us consider a customer who wants to make a decision on what device of a specific category to purchase based on privacy information extracted from their privacy policies. To this end, the user will focus on the keyword analysis feature (see Section 5) extracted by PrivacyLens for the devices in the study dataset. First, the user might want to get an overall picture of the mentions to certain privacy attributes (e.g., first/third party collection, do not track) for all the policies. Figure 6 shows a graph with the Interquartile Range that indicates the extent of variability in the values for each group across all devices examined. We summarize some of the results in the following. First, IoT device policies mention "Third Party Collection" about five times, in general. However, some policies refer to it significantly more, with up to 74 mentions, far exceeding the upper quartile of 11. The "First Party Collection" category has a significant number of policies mentioning it very frequently, as indicated by a high third quartile (103.5) compared to the median (53). "Do Not Track" has a median and Q1 of 0, with only some policies mentioning it, as reflected by the maximum value of 4. This suggests that it is not a widely addressed topic.

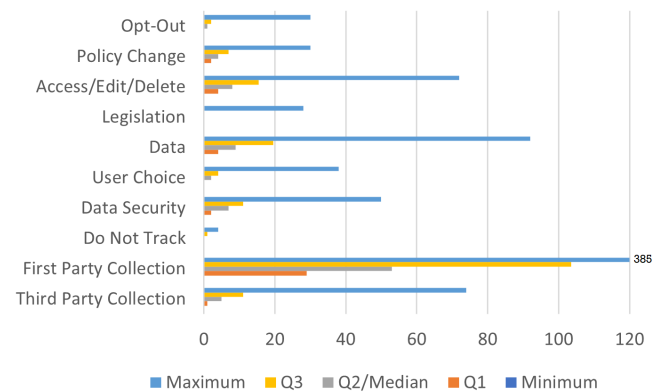


Figure 6: Distribution of privacy attributes in the study dataset.

Next, imagine that the user is especially interested in the "First Party Collection" attribute. Figure 7 shows the number of mentions to such privacy attribute IoT device privacy

policies make (for the subset of privacy policies that make an explicit mention to the device). If the user was interested in buying a smart watch, scale, tracker, or home device, the data shows that some policies did not make any explicit mention to their first party collection process at all while others made a large number of mentions. In contrast, if the user was interested in purchasing a smart doorbell, the results show that all manufacturers mention their first party collection process roughly the same number of times. The starkest difference between manufacturers is for smart lights where (at least) one of the policies had less than 5 mentions of first party collection processes while (at least) one policy had hundreds of mentions. The user would be able to continue to explore the results offered by PrivacyLens further to find a manufacturer for their desired product (e.g., smart light) with a significantly larger mention of their privacy attributes of interest. This way, PrivacyLens would help the user understand the differences between IoT devices w.r.t. their handling of user data and make an informed decision when purchasing one.

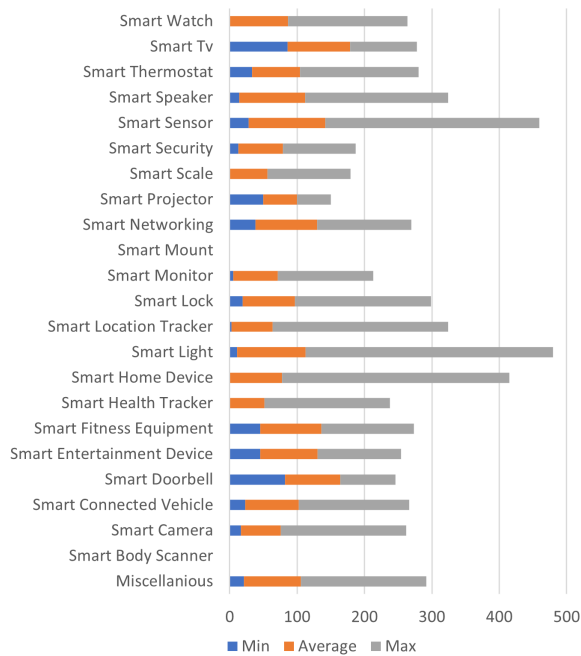


Figure 7: Number of mentions to the first party collection in IoT device privacy policies by type of device.

7.3 PrivacyLens for Privacy Lawyers

Consider a privacy lawyer for a manufacturer who wants to understand the landscape of IoT privacy policies of their competitors. First, the policy lawyer can use PrivacyLens to get insights into the degrees of similarity between different privacy policies which can help identify standard policy formulations and unique clauses, providing them with a foundation for policy writing and revision [3, 6, 49]. For this study, we use

the Word2Vec algorithm [50], which utilizes neural networks to derive meaning from a corpus and create embeddings via the Continuous Bag of Words (CBOW) method [51]. The word embeddings in Word2Vec, derived from hidden layer weights, capture semantic and syntactic similarities. We apply the cosine similarity metric based on the vectors extracted to compute policy similarity.

The results reveal a striking similarity (97%) amongst most policies (see Figure 8). This suggests the existence of a standard template or common elements that many IoT device manufacturers follow while drafting their policies. The remaining 3% of policies showed significant differences, indicating potential areas of innovation or divergence from common policy templates. This is where the privacy lawyer could focus their attention to better understand what sets these policies apart and if their uniqueness offers any strategic advantage, be it in terms of legal compliance, user comprehension, or business objectives.

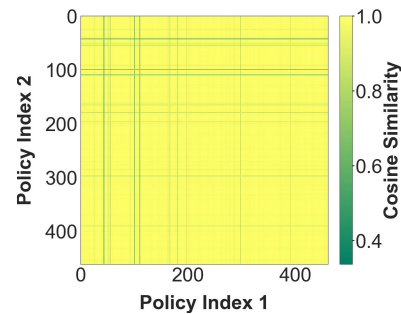


Figure 8: Heatmap of cosine similarities of policy embeddings.

According to Table 6, these features—Coherence Score, Entropy, Frequency of Unique Words, Reading Complexity, Reading Time, Frequency of Imprecise and Connective Words, and Grammar Correctness—highlight unique aspects of company policies. For instance, Fdt’s policy excels in logical structure, while Evapolar’s lacks coherence. Eco4lifhome’s policy shows linguistic diversity, whereas Axis’s is less varied. Nooie’s policy employs a broad range of unique words, indicating conceptual diversity, whereas Eco4lifhome’s policy is conceptually narrower. Eco4lifhome’s and Nooie’s policies require more reading time, while Alarmlock’s, Adero’s, and Umidigi’s policies are more concise. Eco4lifhome’s policy requires higher reading proficiency due to its complexity, while Luxproducts’ policy is simpler. Bulbrite policy contains more imprecise words, risking clarity, while others opt for clarity with fewer imprecise words. Mobvoi and Cablematters use more connectives for better flow, while Axis and Alarmlock use fewer. Fdt and Umidigi maintain better grammar, while Bulbrite’s policy has more errors. However, it is important for privacy lawyers to understand that features that enhance readability, such as simpler language or fewer imprecise words, should not compromise the completeness of the policy. Additionally, privacy policies in our analysis demonstrate a spectrum of distinctive-

ness and readability. Outliers, such as Axis’s policy, are more readable, attributed to their use of simpler language and superior logical flow. However, these user-friendly policies often miss complete details on data usage. This has a potential impact on the text’s flow. However, eco4homelife’s policy, while less readable due to complex language and longer reading time, but it maintains superior flow, making it more detail-oriented. Through a privacy lawyer’s lens, these observations emphasize the necessity of formulating policies that are not only clear and readable but also detailed and legally robust, without sacrificing any aspect. Table 7 presents a summary of the features associated with the analyzed privacy policies. This data can serve as a valuable guide for privacy lawyers seeking to craft or revise privacy policies.

Company Name	Coherence Score	Entropy	Unique Words	Reading Time(min)	Reading Level	Imprecise Words	Connective Words	Grammatical Errors
Eco4lifehome	0.72	10.29	0.27	58	15.10	0	0.01	0.52
Noonie	0.70	5.49	0.67	9	4.39	0	0.03	0.31
Fair	0.92	6.86	0.34	3	5.89	0	0.02	0.15
Axis	0.42	5.41	0.54	2	5.34	0	0.01	0.14
Mobvoi	0.25	7.26	0.45	6	4.24	0.01	0.04	0.18
Alarmlock	0.81	7.14	0.52	2	5.76	0	0.01	0.27
Umidigi	0.59	6.96	0.33	2	4.69	0.01	0.02	0.1
Evapolar	0.22	8.55	0.56	6	13.39	0.01	0.05	0.32
Cablematters	0.74	7.01	0.37	3	6.43	0.01	0.05	0.31
Bulbrite	0.80	7.94	0.40	5	5.96	0.02	0.04	0.52
Airivo	0.61	7.88	0.50	3	6.16	0.01	0.04	0.26
Luxproducts	0.55	7.41	0.44	2	4.38	0.01	0.04	0.19
Adero	0.58	7.56	0.45	2	4.77	0.01	0.05	0.21

Table 6: Readability features extracted for a subset of the policies.

Policy Features	Min Value	Average Value	Max Value
Coherence Score	0.13	0.35	0.92
Freq. of Imprecise Words	0	0.02	0.2
Freq. of Connective Words	0.01	0.04	0.08
Reading Complexity	4.24	11.40	21.73
Reading Time (Min)	2	12.67	107
Entropy	5.41	7.97	10.29
Freq. of Unique Words	0.10	0.30	0.67
Grammatical Errors	0	0.25	1.06

Table 7: Statistics for the readability analysis of the policy corpus.

Next, consider that the privacy lawyer wants to ensure responsible, lawful personal data handling by creating clear privacy policies. The lawyer can use PrivacyLens’s ambiguity analysis results to quickly assess the ambiguity level of the policies (e.g., their own and their competitors). Table 8 shows the distribution of ambiguity levels based on the manufacturer’s country of origin. The European Union (EU), perhaps as a result of the stringent GDPR, had 65% "not ambiguous" policies, the highest proportion. In contrast, China recorded the lowest with 54.3%. In the "somewhat ambiguous" category, the United States led with 27.5%, while the EU had the lowest, 12.5%, reinforcing the influence of strong data privacy laws in reducing policy ambiguity. "Very ambiguous" policies were most prevalent in China, at 23.9%, indicating the need for improved policy clarity, while the US had the least at 18.7%. A unique trend was seen in the smart-connected cars industry, with nearly all their policies being "very ambiguous", suggesting an urgent need for clearer policies. We also analyzed the outlier group with different feature values. While most were classified as "not-ambiguous", a unique exception was

a "Bulbrite" policy, categorized as 'very ambiguous'. These findings can help privacy lawyers push for stricter legislation like the EU’s GDPR in regions with prevalent ambiguous policies, and promote clearer, more specific privacy policies in these sectors.

Manufacturer Country	Not Ambiguous	Somewhat Ambiguous	Very Ambiguous
USA	60.8%	27.5%	18.7%
China	54.3%	21.7%	23.9%
European Union	65%	12.5%	22.5%

Table 8: Percentage of Ambiguity Level across Manufacturers

Finally, analysis of the data in Table 9 revealed that policies not mentioning devices had a higher proportion of 'not-ambiguous' and a lower 'very-ambiguous' category compared to those mentioning devices. This suggests less ambiguity in no device mention policies due to a focus on general data handling and user control, rather than technical specifics. This aligns with research [52] suggesting policies emphasizing high-level privacy principles are more user-friendly.

In conclusion, ambiguity analysis provides valuable insights that privacy lawyers can use to enhance their advocacy efforts, improve privacy policy drafting, and promote better data protection practices in different regions and industries.

Category	Count
Explicit mention of the device - Not Ambiguous	113
Explicit mention of the device - Somewhat Ambiguous	37
Explicit mention of the device - Very Ambiguous	54
No explicit mention of the device - Not Ambiguous	157
No explicit mention of the device - Somewhat Ambiguous	54
No explicit mention of the device - Very Ambiguous	46

Table 9: Device Reference vs Ambiguity

7.4 PrivacyLens for Data Privacy Regulators

In this section, we will describe a potential analysis enabled by PrivacyLens for a Privacy Regulator who wants to monitor the evolution of privacy policies in response to regulatory change. In particular, let us consider a Privacy Regulator who wants to assess policy change pre and post-GDPR. To this end, the regulator will focus on the regulations compliance analysis feature (see Section 5) extracted by PrivacyLens for the devices in the study dataset. First, the regulator might want to get an overall picture of the distribution of similarity scores for policy pairs across all the devices. Figure 9 shows a graph with the distribution of similarity scores with frequency/density across all policy pairs showing variability in scores. We summarize some of the results in the following. For instance, 40% of privacy policies remained unchanged, 14% privacy policies had a similarity score of less than 60% showing significant changes in the privacy policies, while 84% privacy policies showed a score greater than 60% implying prior compliance or minimal changes.

Second, the regulator might want to see GDPR impacts across different countries. Figure 10 shows a graph with pri-

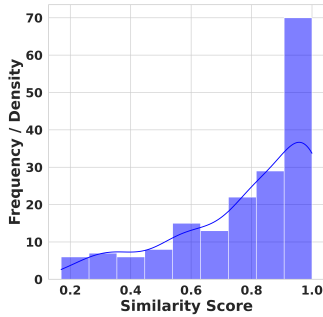


Figure 9: Distribution of similarity scores.

vacancy policy similarity scores across countries, highlighting patterns and trends in the data. For instance, the United States, Japan, and China made substantial adjustments to their privacy policies following the GDPR enactment, a reflection of their enhanced compliance with the new privacy regulations. In contrast, Poland and Hong Kong showed high similarity scores, suggesting a uniform approach to GDPR compliance. The United Kingdom and Germany, too, showed a high similarity score, indicating a similar trend. Surprisingly, Canadian manufacturers exhibited the most changes post-GDPR, suggesting a significant policy shift. With a global average similarity score of 0.77 between pre-and post-GDPR policies, it's evident that while organizations have made GDPR-related adjustments, there's room for further revisions.

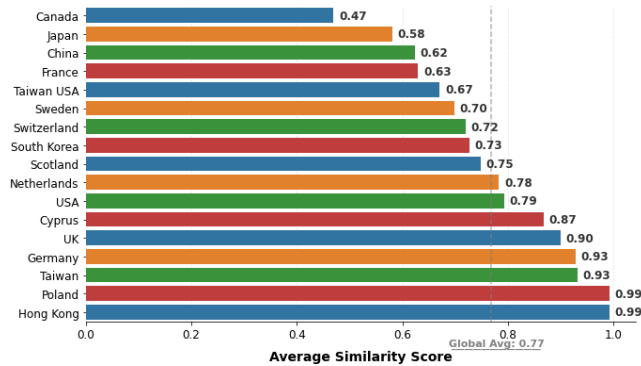
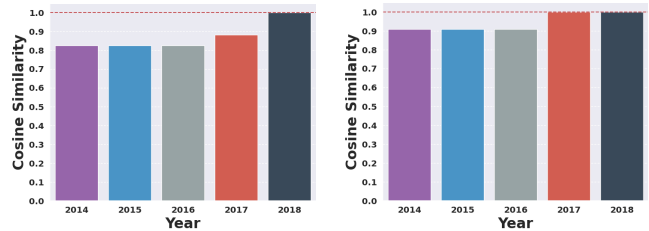


Figure 10: Bar Graph: Visualizes privacy policy similarity scores across countries, highlighting patterns and trends in the data.

Next, imagine that the regulator is especially interested in checking for the timeline where changes towards compliance with GDPR took effect. Figure 11 shows the timeline of changes in privacy policies. If the regulator was interested in assessing this timeline for European countries like Germany (Schluter from 2014-2018), the data shows that a significant shift was observed in 2017, implying that the firm adjusted its policy for GDPR compliance during this time. In contrast, if the regulator was interested in the US (Mielusa from 2014-2018), the data shows that a significant policy shift was observed in 2017, while making a more prominent shift in 2018, implying the firm adjusted its policy for GDPR com-

pliance during this timeline. The regulator would be able to continue to explore the results offered by PrivacyLens. This exploration could help identify geographic areas facing compliance challenges. Consequently, this insight could direct enforcement focus toward addressing prevalent or substantial issues.



(a) Cosine similarities across for policies of Mieleusa (b) Cosine similarities across for policies of Schluter

Figure 11: Visualizations of similarity across policies

Overall, a systematic comparison of pre and post-GDPR privacy policies can provide regulators with a data-driven approach to ensuring and enhancing compliance with the regulation. Also, they can pinpoint common trends, and potentially challenging compliance areas, and guide future regulatory efforts.

8 Conclusions

In this paper, we presented PrivacyLens, a framework that collects and analyzes privacy policies of IoT devices. It incorporates a module to automatically find IoT devices in e-commerce sources, search for their current privacy policies, and even extract their privacy policies from recent years, if available. It also incorporates machine learning and natural language processing techniques to extract insights about privacy policies. This encompasses identifying and interpreting keyword usage, deciphering readability aspects like reading time, reading level, and entropy, and examining ambiguity. The system also notes the last update of the policy, privacy-related insights, the country of the manufacturer and any device mentions within the policy. PrivacyLens is currently deployed and continuously collecting and analyzing privacy policies (1,200 policies from more than 7,300 IoT devices have been collected at the time of submitting this paper). All the information collected and produced by PrivacyLens is publicly available to offer insights to customers, researchers, policymakers, and regulators, enabling them to make well-informed decisions when dealing with IoT privacy practices. In future work, we plan to extend the input sources that the system collects to include information extracted from user reviews and reported data breaches for IoT devices.

References

- [1] Anil Alter, Michele M Tugade, and Barbara L Fredrickson. Smartness as a continuous variable: identifying dimensions of intelligent environments. *Frontiers in psychology*, 7, 2016.
- [2] Oberlo. Smart home statistics. <https://www.oberlo.com/statistics/smart-home-statistics#:~:text=Smart%20home%20statistics%20show%20that%20in%202018%2C%2029.5,the%20coming%20years%2C%20reaching%2064.1%20million%20by%202025.>, 2018. Accessed: May 24, 2023.
- [3] Noah Apthorpe, Dillon Reisman, and Nick Feamster. Always on (even when we’re off the grid): Privacy risks and conservation benefits associated with the internet of things. In *IEEE Symposium on Security and Privacy (SP)*, 2017.
- [4] Rolf H Weber. *Internet of Things*. Springer, 2010.
- [5] Jayavardhana Gubbi, Rajkumar Buyya, Slaven Marusic, and Marimuthu Palaniswami. Internet of things (iot): A vision, architectural elements, and future directions. *Future Generation Computer Systems*, 29(7), 2013.
- [6] Paul Biocco, Mahsa Keshavarz, Patrick Hines, and Mohd Anwar. A study of privacy policies across smart home companies. In *An Interactive Workshop on the Human aspects of Smarthome Security and Privacy (WSSP 2018), Symposium on Usable Privacy and Security (SOUPS)*, 2018.
- [7] Luigi Atzori, Antonio Iera, and Giacomo Morabito. The internet of things: A survey. *Computer Networks*, 54(15), 2010.
- [8] Alessandro Acquisti. Privacy and data protection in the age of big data: A time for big decisions. *Computer Law & Security Review*, 31(6), 2015.
- [9] Benjamin Fabian, Tatiana Ermakova, and Tino Lentz. Large-scale readability analysis of privacy policies. In *International Conference on Web Intelligence*, 2017.
- [10] Christina L Madden and Douglas B Mc Donald. The scoring of digital rights: a preliminary analysis. *International Journal of Law and Information Technology*, 17(2), 2009.
- [11] Barbara Krumay and Jennifer Klar. Readability of privacy policies. In *IFIP Annual Conference on Data and Applications Security and Privacy*, 2020.
- [12] Mozilla Foundation. Privacy not included. <https://foundation.mozilla.org/en/privacynotincluded/>, 2021. Accessed: 29th April 2023.
- [13] Internet Archive. Wayback machine, 1996.
- [14] Norman M. Sadeh, Alessandro Acquisti, Travis D. Breaux, Lorrie Faith Cranor, Aleecia M. McDonald, Joel R. Reidenberg, Noah A. Smith, Fei Liu, N. Cameron Russell, Florian Schaub, Shomir Wilson, Jim Graves, Pedro Giovanni Leon, Rohan Ramanath, and Ashwini Rao. Towards usable privacy policies: Semi-automatically extracting data practices from websites’ privacy policies. 2014.
- [15] Shomir Liu, Yang Liu, Yuan Li, Shuqin Li, and Fei Niu. Polisis: Automated analysis and presentation of privacy policies using deep learning. *USENIX Security Symposium*, 27(3), 2018.
- [16] Duc Viet Bui, Kang G. Shin, Jong-Min Choi, and Junbum Shin. Automated extraction and presentation of data practices in privacy policies. *Proceedings on Privacy Enhancing Technologies*, 2021, 2021.
- [17] Elisa Costante, Jerry den Hartog, and Milan Petković. What websites know about you. In Roberto Di Pietro, Javier Herranz, Ernesto Damiani, and Radu State, editors, *Data Privacy Management and Autonomous Spontaneous Security*, 2013.
- [18] Michael Gebauer, Faraz Mashhur, Nicola Leschke, Elias Grünwald, and Frank Pallas. A human-in-the-loop approach for information extraction from privacy policies under data scarcity. *arXiv preprint arXiv:2305.15006*, 2023.
- [19] Mikhail Kuznetsov, Evgenia Novikova, Igor Kotenko, and Elena Doynikova. Privacy policies of iot devices: Collection and analysis. *Sensors*, 22(5), 2022.
- [20] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2018.
- [21] Joshua Upp. Understanding amazon alexa’s technology. *Journal of Technology and Science Education*, 8(2), 2018.
- [22] Rachel Story, Michael L Paul, and Dennis S Dye. Natural language processing and machine learning for identifying smartphone applications that can collect personal information from children on the google play store. In *8th International Conference on Educational and Information Technology (ICEIT)*, 2019.
- [23] Ryan Amos, Gunes Acar, Elena Lucherini, Mihir Kshirsagar, Arvind Narayanan, and Jonathan Mayer. Privacy policies over time: Curation and analysis of a million-document dataset. pages 2165–2176, 2021.

- [24] John W Ratcliff and David E Metzener. Pattern-matching-the gestalt approach. *Dr Dobbs Journal*, 13(7):46, 1988.
- [25] Peter Voigt and Arndt von dem Bussche. The eu general data protection regulation (gdpr): A practical guide. *Springer*, 2017.
- [26] California consumer privacy act (ccpa) of 2018, 2018.
- [27] 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, et al. The creation and analysis of a website privacy policy corpus. In *54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016.
- [28] Alessandro Oltramari, Dhivya Piraviperumal, Florian Schaub, Shomir Wilson, Sushain Cherivirala, Thomas B. Norton, N. Cameron Russell, Peter Story, Joel R. Reidenberg, and Norman M. Sadeh. Privonto: A semantic framework for the analysis of privacy policies. *Semantic Web*, 9(2), 2018.
- [29] David Banisar. National comprehensive data protection/privacy laws and bills 2023. *Privacy Laws and Bills*, 2023.
- [30] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. *Conference of the North*, 2019.
- [31] Daniel Cer, Mona Diab, Eneko Agirre, Iñaki Lopez-Gazpio, and Lucia Specia. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. In *11th International Workshop on Semantic Evaluation (SemEval-2017)*, 2017.
- [32] David W Hosmer Jr, Stanley Lemeshow, and Rodney X Sturdivant. *Applied logistic regression*. John Wiley & Sons, 2013.
- [33] Claude E Shannon. Prediction and entropy of printed english. *Bell system technical journal*, 30(1), 1951.
- [34] Aleecia M McDonald and Lorrie Faith Cranor. The cost of reading privacy policies. *Isjlp*, 4:543, 2008.
- [35] Marc Brysbaert. How many words do we read per minute? a review and meta-analysis of reading rate. *PsyArXiv Preprints*, 2019.
- [36] Elfrieda H Hiebert. Unique words require unique instruction.
- [37] Adam Gettgey. Natural language processing is fun! *Medium*, July, 18, 2018.
- [38] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. 2017.
- [39] Shaheen Syed and Marco Spruit. Full-text or abstract? examining topic coherence scores using latent dirichlet allocation. In *IEEE International conference on data science and advanced analytics (DSAA)*, 2017.
- [40] Hua Yu and Jie Yang. A direct lda algorithm for high-dimensional data—with application to face recognition. *Pattern recognition*, 34(10), 2001.
- [41] Daniel Naber et al. A rule-based style and grammar checker. 2003.
- [42] Julien B Kouame. Using readability tests to improve the accuracy of evaluation documents intended for low-literate participants. *Journal of MultiDisciplinary Evaluation*, 6(14), 2010.
- [43] Rudolf Flesch. Flesch-kincaid readability test. *Retrieved October*, 26(3), 2007.
- [44] Joel R Reidenberg. Privacy policies as decision-making tools: An evaluation of ayres & braithwaite, irobot, and norton. *Journal of Empirical Legal Studies*, 13(4), 2016.
- [45] Anantaa Kotal, Anupam Joshi, and Karuna Pande Joshi. The effect of text ambiguity on creating policy knowledge graphs. In *2021 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCLOUD/SocialCom/SustainCom)*, 2021.
- [46] Leo Breiman. Random forests. *Machine learning*, 45(1), 2001.
- [47] Steven Bird, Ewan Klein, and Edward Loper. Natural language processing with python. In *O'Reilly Media, Inc.*, 2009.
- [48] 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), 2016.
- [49] Joel R. Reidenberg, Jaspreet Bhatia, Travis D. Breaux, and Thomas B. Norton. Ambiguity in privacy policies and the impact of regulation. *The Journal of Legal Studies*, 45(S2), 2016.

- [50] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [51] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, 2013.
- [52] Didier Masha, Lek Chaisorn, and Santi Phithakkitnukoon. Smart homes privacy policies analysis: A critical information analysis. *Journal of Ambient Intelligence and Humanized Computing*, 11(10), 2020.

Appendix A Privacy Policy Analysis Details

Table 10 contains the taxonomy of imprecise words and Table 11 contains the taxonomy of connective words (both extracted from [45]) that PrivacyLens uses in its analysis of an IoT device privacy policy.

Imprecise Words	
Modal Words	may,might,likely,can could,would
Usable Words	easy,adaptable familiar,extensible
Probable Words	probably,possibly,optionally
Numeric Words	anyone, certain everyone, numerous some,most,few much,many,various including but not limited to
Condition Words	depending,necessary inappropriate,appropriate as needed,as applicable otherwise reasonably from time to time
Generalization Words	generally,mostly,widely commonly,usually,general Normally,typically,largely often,primarily among other things

Table 10: Taxonomy of imprecise words extracted from [45].

Connective Words	
Copulative Words	and, both, as well as, not only, but also
Control Flow Words	if, then, while
Anaphorical Words	it, this, those

Table 11: Taxonomy of connective words extracted from [45].

Appendix B Feature Evaluation

Privacy Document	Value (Our Approach)
Minimum Correct Grammar	0.00
Minimum Imprecise Words	0.00
Minimum Connective Words	0.02
Maximum Correct Grammar	0.95
Maximum Imprecise Words	0.04
Maximum Connective Words	0.06

Table 12: Approach

Privacy Document	Value (Ground Truth)
Minimum Correct Grammar	0.06
Minimum Imprecise Words	0.09
Minimum Connective Words	0.025
Maximum Correct Grammar	0.23
Maximum Imprecise Words	0.09
Maximum Connective Words	0.076

Table 13: Ground Truth

Appendix C Country distribution

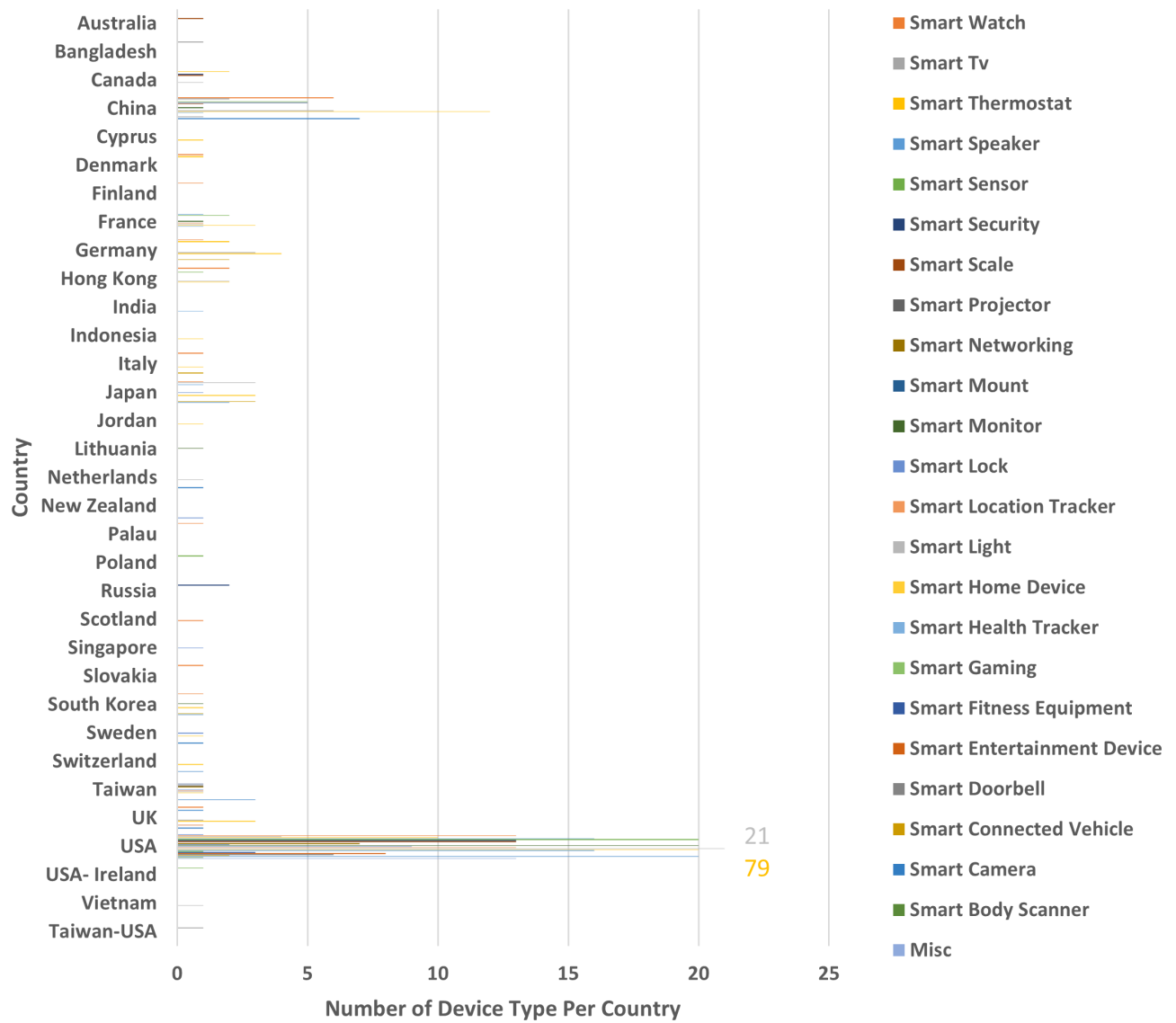


Figure 12: Distribution of Country for each Device Type