# The PrivaSeer Project: Large-Scale Resources for Analysis of Privacy Policy Text

**Shomir Wilson[1], Florian Schaub[2], Lee Matheson[3], Shahriar Shayesteh[1], and Lu Xian[2]**

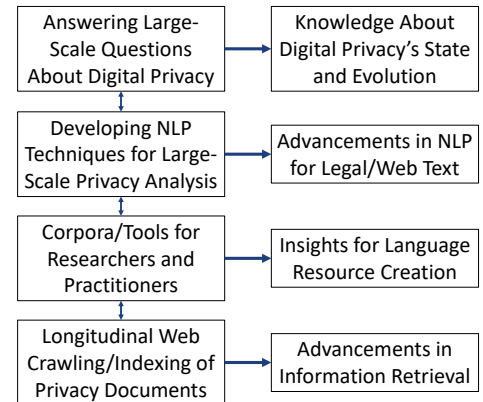[1]College of Information Sciences and Technology, Pennsylvania State University
[2]School of Information, University of Michigan, [3]Future of Privacy Forum

**You can download the PrivaSeer Corpus of 3,967,487 privacy policies and use the PrivaSeer Search Engine at https://privaseer.ist.psu.edu/**



*PrivaSeer Search interface*



*Overall project structure*

## Motivation and Goal

The wealth of privacy policies available on the web contrasts with the challenges of understanding the state of digital privacy at scale.

The PrivaSeer Project builds large-scale, longitudinal, annotated, and usable resources for the study of website privacy policies.
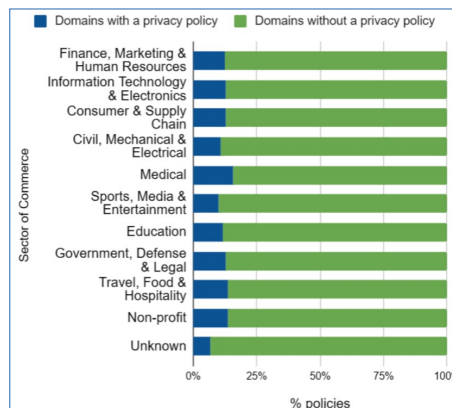
## Resources

The **PrivaSeer Corpus** consists of 3,967,487 website privacy policies, complete with metadata such as sector of activity, readability metrics, and mentions to tracking technologies and regulations.

**PrivaSeer Search** is a search engine for exploring the corpus. This allows potential corpus users to test-drive the corpus prior to downloading it. It also makes the corpus accessible to privacy researchers with limited data analysis or programming skills.

The corpus and search engine are available to the public on the project website. The corpus is subject to a Creative Commons BY-NC-SA license.

## Analysis

See the project website for our publications in ACL, DocEng, ICWE, and PETS.



Percentages of domains with privacy policies by sector of commerce



Distribution of dates that appear in privacy policies (up to September 2021, as limited by data collection)