# PrivaSeer

# A Large-Scale, Longitudinal Resource to Advance Technical and Legal Understanding of Textual Privacy Information

**Shomir Wilson & Lee Giles**, Penn State

**Florian Schaub**, University of Michigan

**Gabriela Zanfir-Fortuna**, Future of Privacy Forum

PennState
College of Information
Sciences and Technology

UNIVERSITY OF MICHIGAN

FUTURE OF PRIVACY FORUM

Try our search engine for over 1.4M privacy policies at privaseer.ist.psu.edu

We're building a **large-scale**, **annotated**, and **searchable resource** of **privacy-related texts**
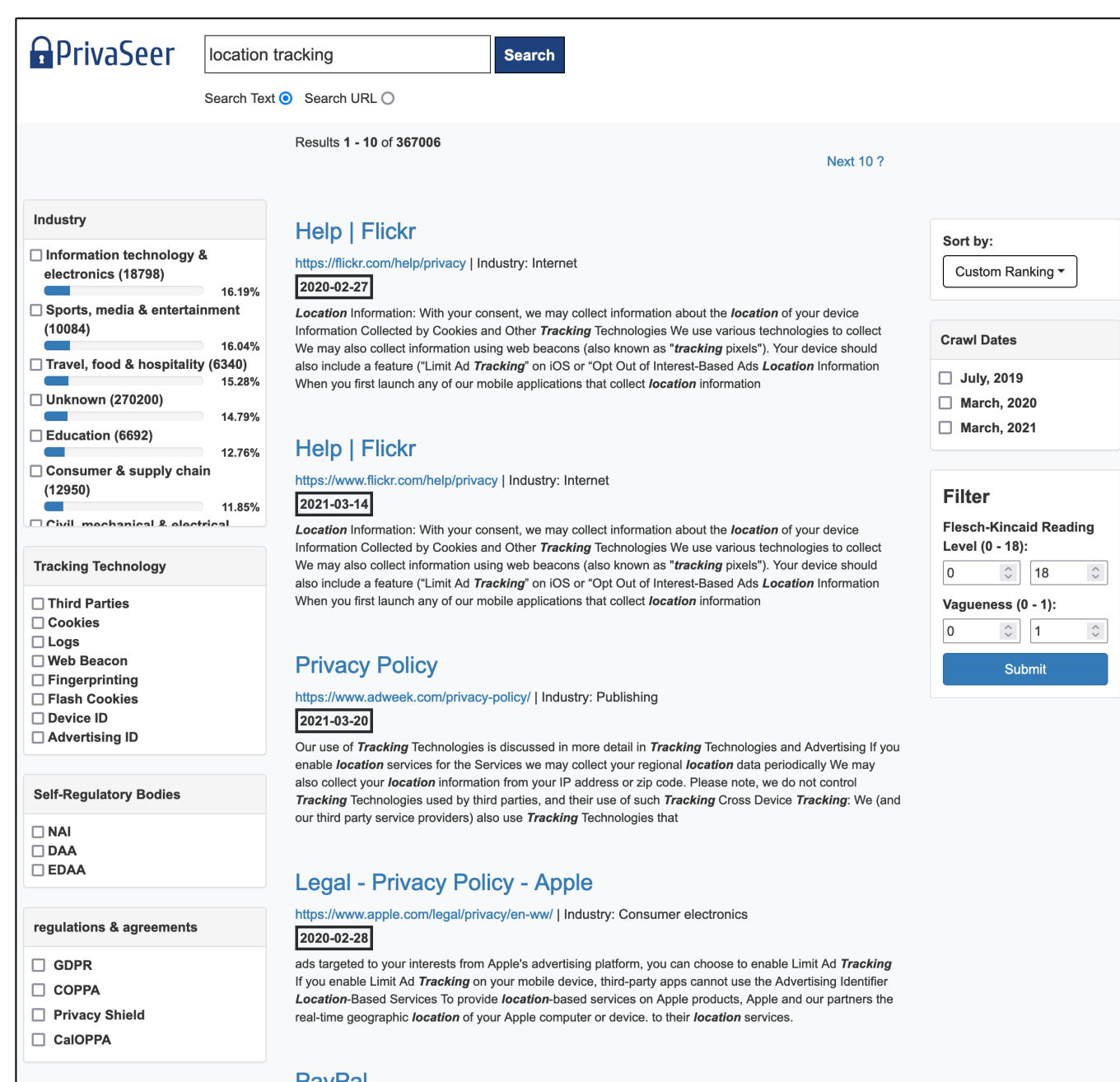
- Advancing natural language processing techniques of privacy documents
- Facilitating research on privacy documents with infrastructure and tools
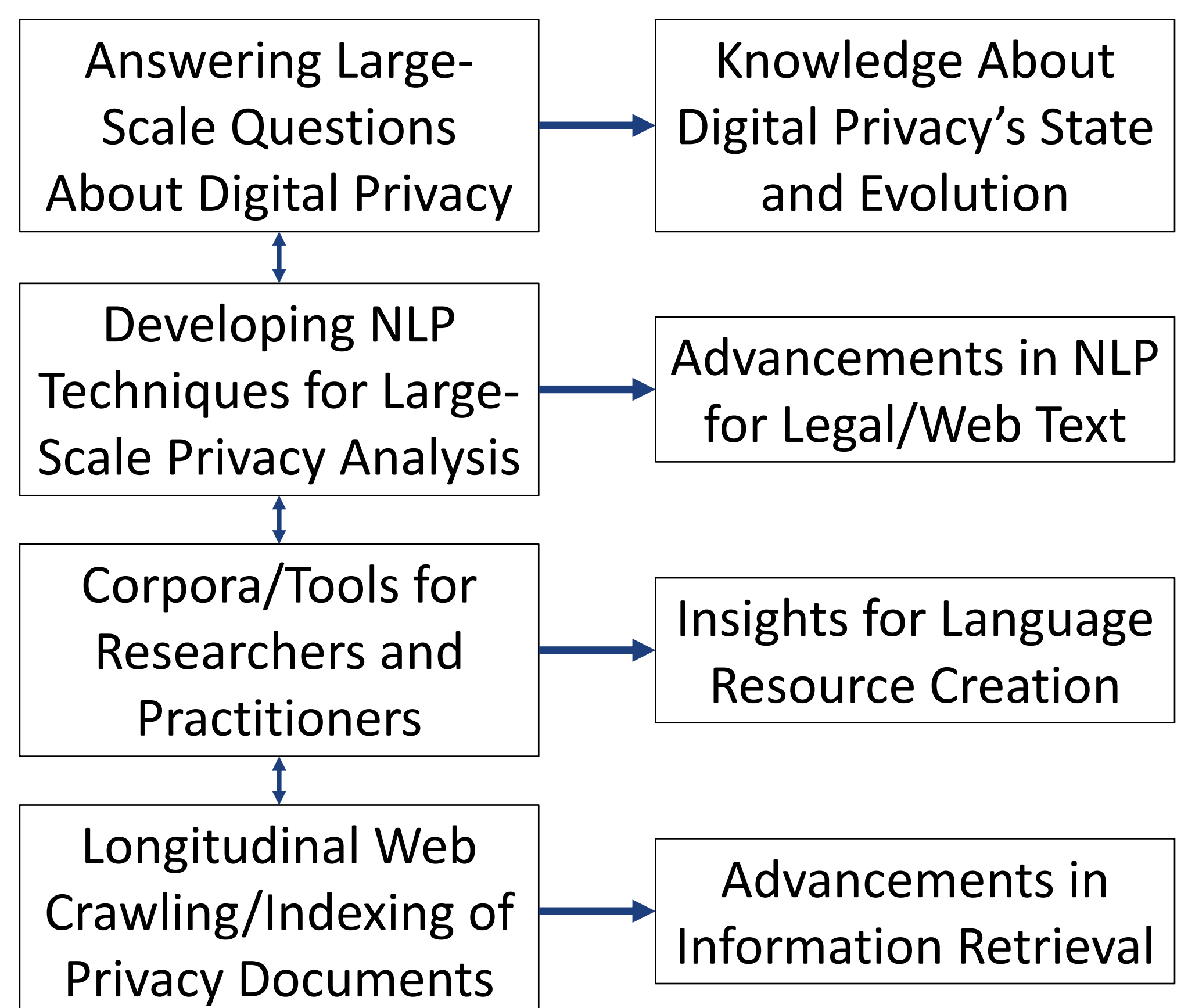- Analyzing the state of privacy at scale

## Challenges

- **A huge volume of privacy-related text exists online**: privacy policies, terms of use, cookie policies, do-not-sell notices, privacy laws, bills, regulatory guidance, etc.
- **Collecting and analyzing privacy text is necessary** for privacy research, practice, and effective regulation and enforcement.
- **Lack of infrastructure and tools** leads to repeated and duplicative effort across research projects.



*PrivaSeer search engine prototype (April 2022)*

## Project Objectives:

- **Enable analyses** of the state of privacy at an unprecedented scale
- **Remove barriers** for privacy analysis
- **Provide resources and tools** for researchers, practitioners, and policymakers

| | |
|---|---|
| Answering Large-Scale Questions About Digital Privacy | Knowledge About Digital Privacy's State and Evolution |
| Developing NLP Techniques for Large-Scale Privacy Analysis | Advancements in NLP for Legal/Web Text |
| Corpora/Tools for Researchers and Practitioners | Insights for Language Resource Creation |
| Longitudinal Web Crawling/Indexing of Privacy Documents | Advancements in Information Retrieval |

## Broader Impacts

*Research:* Radically reduce effort needed to analyze privacy documents at scale.

Provide public tools and release corpora for research.

*Industry & Public Policy:* Enable data-driven and evidence-based decision and policy making.

Outreach to policymakers and practitioners.

*Education & Broadening Participation:* Train diverse team of students and fellows.

Foster diverse and interdisciplinary privacy research community.