

Privacy in a Data-Driven World

Our personal privacy – the precious and critical component of our individual identities and behavior – has come under siege in today’s data-driven world. Web services, mobile applications, and third parties collect and use our personal data for their own purposes – some in line with our interests, others potentially not. And they do so with almost no accountability, raising the risk for deceptive and unfair practices. My vision points toward a new model for how we address these personal privacy issues. I see a digital world where users are more aware of the privacy consequences of their online actions and make more informed decisions about the services they use. In my model, services and applications are held accountable for their actions and are explicitly constructed to protect user privacy and best interests.

To forge this new world, I design, build, and evaluate *new transparency tools* that increase users’ awareness of and society’s oversight over how applications use personal data, and *new development abstractions and tools* that facilitate the construction of privacy-mindful applications. Recognition and impact of my work include multiple best paper awards, technology transfers to two companies, invitations to address the Federal Trade Commission (FTC) as well as members of the Cybersecurity Caucus of the U.S. House of Representatives, media coverage such as New York Times and The Economist, and prestigious awards including NSF CAREER, Popular Science Brilliant 10, Microsoft Faculty Fellowship, and Google Research Award.

Transparency Tools for the Data-Driven Web: Are third-party web trackers watching our children’s online activities and targeting them? Do shopping and loan sites tailor their prices based on what they know about us? Today, we have no solid answers to such questions. As end-users, we have little insight into how our data is being used by web services and third-party groups that collect it. And privacy watchdogs – such as FTC investigators and journalists – lack tools to track online personal data flow to discover deceptive or unfair practices.

Under my leadership, my group is building the first scalable, generic, and reliable tools to detect data flows within and across web services. Our initial system, XRay (USENIX Security’14), offers premier system design along with the theoretical building blocks to detect the use of digital personal data for targeting and personalization. The key insight in XRay is to infer targeting by *correlating* user inputs (such as searches, emails, or locations) to service outputs (such as ads, recommendations, or prices) based on observations obtained from user profiles populated with different subsets of the inputs. Our latest tool, Sunlight (CCS’15), leverages rigorous statistical methods to determine not only the correlations, but more significantly the causes of online targeting at great scale and based on solid statistical justification.

Last year, we used our tools to run the largest-scale studies of online ad targeting. One thing we found is evidence that contradicts two of Google’s own privacy statements regarding the lack of targeting on sensitive topics in its networks. Our work has attracted attention from the FTC, with whom we are exploring a collaboration to adapt our tools to advance their investigations into unfair and deceptive web practices. DARPA’s Information Science and Technology (ISAT) advisory group has commissioned me to design a web transparency workshop to galvanize a community around this critical topic.

Development Tools for Privacy-Mindful Applications: At present, programmers lack the development tools to build privacy-mindful applications. In direct response to this critical need, I am building a *privacy testing toolkit* that will support the detection and debugging of many kinds of privacy “bugs.” The first tool in my toolkit is FairTest, which detects unwarranted association bugs, a subtle privacy bug specific to data-driven applications. Unwarranted association bugs were responsible for the racist labels recently found in Google’s image tagging system as well as the unintended, discriminatory effects found in Staples’s online pricing algorithm. FairTest catches such bugs and assists programmers in remedying them.

We applied FairTest to three real-world data-driven applications: a predictive healthcare app, an image labeling system, and a movie recommender. FairTest exposed previously unknown association bugs in all three applications. For example, it revealed that the healthcare application – the winner of a Heritage Health Prize Competition – offers good overall accuracy (85%) but concentrates its error disproportionately on elderly patients with certain conditions, where the error can reach 45%.

In another project, I am defining a new approach to privacy protection in cloud data ecosystems, challenging the conventional approach of collecting enormous amounts of personal data and archiving them forever. Not all data collected and archived by services is needed for those services to function well. Distinguishing “nice to have” data from “need to have” data could enable developers to limit data exposure to cloud hackers and nosy employees. My work develops mechanisms to differentiate data’s utility over time and integrates these mechanisms into popular distributed data processing infrastructures, such as Spark. My approach adds rigor and privacy protection to their data management. Previously, I applied similar concepts to minimize data exposure in theft-prone, mobile devices, where I showed that they can reduce sensitive-data exposure time by 95-97% (OSDI’14, OSDI’12, EuroSys’11 best paper).

In all my work, I bring a passionate vision of how computing can enrich our lives without imposing undue or unknowable personal costs.