

Volume I: Technical and Management Proposal

Cover Page

BAA Number	DARPA-BAA-15-29
Technical Area	(2) Human Data Interaction (HDI)
Proposal Title	Sunlight is the Best Disinfectant: Increasing Privacy through Awareness with the Hubble Scalable Web Transparency Infrastructure
Lead Organization	The Trustees of Columbia University in the City of New York
Type of Business	Other Educational
Contractor's Reference Number	RASCAL PT-AABL4408
Technical Point of Contact	Prof. Roxana Geambasu Department of Computer Science, Mail Code 0401 Columbia University, 1214 Amsterdam Avenue New York, NY 10027-7003 212-939-7099 (v) roxana@cs.columbia.edu
Administrative Point of Contact	Kammy Lou Cabral Director Sponsored Projects Administration 615 West 131st Street, Room 254, Mail Code 8725 New York, NY 10027-7003 (212) 854-6851 (v) ms-grants-office@columbia.edu
Subcontractor Information Technical Point of Contact	University of Washington Prof. Franziska Roesner Department of Computer Science & Engineering, Box 352350 University of Washington, Paul Allen Center, 185 Stevens Way Seattle, WA 98195-2350 206-221-8248 (v) 206-543-2969 (f) franzi@cs.washington.edu
Administrative Point of Contact	Lynette Arias Director of Sponsored Research University of Washington, Office of Sponsored Programs 4333 Brooklyn Ave. NE Seattle, WA 98195 206-543-4043 (v) 206-685-1732 (f) osp@u.washington.edu
Award Instrument Requested	Grant
Period of Performance	09/01/2015 – 08/31/2019
Places of Performance	New York, NY; Seattle, WA
Proposal Validity Period	120 days
Prime DUNS Number	049179401
Prime TIN	13-5598093
Prime CAGE Code	1B053

Contents

1	Executive Summary	1
2	Goals and Impact	2
3	Collaborative Research Team Concept	5
4	Technical Plan	5
4.1	Thrust 1: The Hubble Transparency Infrastructure and Abstractions	6
4.1.1	The Hubble System and Core Abstractions	6
4.1.2	Statistical Correlation and Causal Inference	7
4.1.3	Privacy-Preserving Transparency Protocols	9
4.2	Thrust 2: Transparency and User Awareness Tools	12
4.2.1	AdObservatory: Revealing Targeting in Online Advertising	12
4.2.2	DiscriminationObservatory: Revealing Online Price Discrimination	12
4.2.3	CollectionObservatory: Revealing Third-Party Content and Tracking . . .	13
4.2.4	LocationObservatory: Revealing Privacy Implications of Location Tracking	14
4.3	Thrust 3: User Studies and Transparency Tool Measurements	15
5	Personnel and Management Plan	16
5.1	Personnel	16
5.2	Integration and Evaluation	18
6	Capabilities	18
7	Statement of Work	20
7.1	Phase 1 (Months 1-18)	20
7.2	Phase 2 (Months 19-36)	21
7.3	Phase 3 (Months 37-54)	23
8	Schedule and Milestones	24
9	Cost Summary	25
10	Appendix A	31
10.1	Team Member Identification	31
10.2	Government or FFRDC Team Member Proof of Eligibility to Propose	31
10.3	Government or FFRDC Team Member Statement of Unique Capability	31
10.4	Organizational Conflict of Interest Affirmations and Disclosure	31
10.5	Intellectual Property (IP)	31
10.6	Human Subjects Research (HSR)	31
10.7	Animal Use	31
10.8	Representations Regarding Unpaid Delinquent Tax Liability or a Felony Conviction under Any Federal Law	31
10.9	Cost Accounting Standards (CAS) Notices and Certification	31

1 Executive Summary

Motivation and Goal: Today’s web services – such as Google, Amazon, and Facebook, as well as third-party advertisers less visible to users – collect and leverage user data for varied purposes, including personalizing recommendations, targeting advertisements, and adjusting prices. At present, users have little insight into how their data is being collected or used and how that affects them. This lack of awareness prevents them from making informed choices about the services they use, what they should be revealing to those services and what not, or what protection tools they should use to prevent misuse. Our goal is to develop *user awareness tools* that will help users gain a better understanding of the implications of their online actions by revealing to them concretely how their data is being collected and used by the services to target them. For example, one user awareness tool could reveal what specific data within a user’s profile – such as emails, prior browsing behaviors, etc. – are being targeted by each advertisement they receive. Another tool could reveal to a user that she is seeing a differentiated price, and specifically which data within her profile triggered that differentiation. In support of such tools, we propose to build *Hubble*, an extensible, generic, and scalable infrastructure that provides the necessary scientific methods and programming abstractions to facilitate the building of many such user awareness tools. Using Hubble, we will develop and evaluate several user awareness tools, and will study how transparency and awareness can help shape user actions and enable them to better manage their online privacy. Our effort targets *Technical Area #2* (Human Data Interaction) and is *fundamental research*.

Key Technical Challenges: Constructing user awareness tools raises significant and unresolved challenges. First, once data is given out to a service, how can one still track its use? Tracking data in a controlled environment, such as a modified operating system, language, or runtime, is an old problem with a well-known solution: taint tracking systems [26,37,47,93]. However, is it possible to track data in an uncontrolled environment, such as the Web? Can robust, generic mechanisms assist in doing so? What kinds of data uses are trackable and what are not? How would the mechanisms scale with the amount of data being tracked? Second, constructing user awareness tools that do not themselves create new privacy challenges is a difficult challenge. Intuitively, to reveal how data is being used, a user awareness tool needs to monitor that user’s data, and perhaps share it with a third party that aggregates data from multiple users. Why should the users trust those tools and the third-party that run them, and how can we minimize that trust? Third, quantifying the effect of user awareness tools on the end-users is an open question. For user awareness tools to be effective, they must not only help educate users – and watchdog organizations like the Federal Trade Commission (FTC) or the Electronic Frontier Foundation (EFF) – about data collection and use, but they must provide useful and auditable actions that users can take to manage their privacy.

Review of Proposed Technologies: Our project will develop both the tools and the necessary scientific methods and infrastructures necessary to increase users’ awareness over what happens with their data once they share it with web services. The center piece of our proposal is *Hubble*, the first scalable and extensible infrastructure that provides a series of abstractions for accurately detecting data collection and use for targeting and personalization within and across (largely) arbitrary web services. Hubble’s three main research contributions are: (1) leveraging *state-of-the-art statistical methods in unique ways* to accurately detect targeting in black-box services based on experiments with differentiated user profiles, (2) providing an *extensible and dynamic architecture* that enables automatic validation, refinement, and attribution of targeting inferences, and (3) providing primitives for doing targeting inference in a privacy-preserving way.

Using Hubble, we will build a number of useful user awareness tools, which reveal specific aspects about data collection and targeting on the web which may now escape notice by the end users. Those tools include: (1) *CollectionObservatory*, a tool that reveals third-party web content that invisibly collects information about users’ browsing behaviors; (2) *AdObservatory*, a tool that reveals how third-party web trackers leverage the information they collect about the users – such as visited pages, Facebook Likes, or explicitly shared information – to target ads at them; (3) *DiscriminationObservatory*, a tool that reveals to the end users any differential treatment in the prices or offers they get from ecommerce, lending, and mortgage websites; and (4) *LocationObservatory*, a tool that communicates to users how multiple location records associated with them both relate to their sensitive data (e.g., age, ethnicity) and personalization offered by services.

Current Approaches and Limitations: Our project will create *robust, generic user awareness tools to track the use of personal data at fine granularity (e.g., individual emails, photos, or visited websites) within and across arbitrary Web services*. At present, hardly any such tools exist, and the science of tracking the use of personal Web data at scale and at a fine granularity is extremely limited. Our own recent system, XRay [60], includes some preliminary results that transparency at fine granularity is possible, but does not address any of the significant scaling, privacy, and usability challenges defined above. We have also previously developed TrackingObserver [81] to detect third-party trackers on the web, but it remains limited in terms of the types of data collection it can detect (notable, omitting fingerprint-based trackers) and does not provide information directly useful to end users. Other transparency systems, such as Bobble [90], AdFisher [30], and OpenWPM [38], are either not generic (e.g., Bobble reveals personalization of news and search results on based on a few user attributes) or operate at small scale [30, 90].

Expected Impact: The greatest impact of our work will be to increase user awareness about the implications of their online actions and to provide building blocks for tools that empower users to manage the information that they reveal. We believe that a vital part of protecting private data that users knowingly provide to third parties is to enable non-expert users to *know more, take action, and verify the results of their actions*. Moreover, we believe that by empowering users, as well as third-party privacy watchdogs, with transparency tools we will help transition the web services world toward a more privacy-aware future. In Louis Brandeis’ own words – “Sunlight is said to be the best of disinfectants; electric light the most efficient policeman” [?]. Hubble will help bring a new level of oversight and accountability into a very obscure Web world, thereby putting pressure on web services to be more privacy aware. Finally, while this proposal focuses on awareness tools and building blocks for the Web, we believe that our technologies will be applicable more broadly to track how sensitive information – be it users’ personal data, proprietary enterprise information, or classified defense data – is being used (or abused) by the parties that obtain it (such as web services, partner enterprises, or foreign governments). We thus expect that extended versions of Hubble could be applicable to use cases of national importance beyond protecting and increasing end-user privacy on the web.

Cost, Duration, and Team: Our proposed effort will last 4.5 years (starting on 09/01/2015), with a total cost of \$3,960,419. The team are from Columbia University and University of Washington.

2 Goals and Impact

Many of today’s pervasive practices that collect and leverage user data are invisible, or at best unexpected, to users. For example, web and mobile applications commonly collect and aggregate

information about users (including browsing behaviors, location, and unique identifiers) for the purposes of targeted advertising or other types of personalization [60, 78]. Many of today’s users have some notion that this data collection is happening (e.g., through extensive media reporting on the topic [85]) and that they are exchanging some amount of private information for the use of free services (email, search, social media). Indeed, these practices are typically disclosed in terms of service agreements, to which users must technically agree to use an application or service. However, users’ understanding of the extent of this data collection, as well as its use and implications, remains limited and abstract [87]. **Thus, a necessary goal on the path to protecting private data that users knowingly provide to third parties is to help non-expert users *know more, take mitigating actions, and verify the results of their actions.***

To this end, we propose the design, development, and evaluation of a new generation of **user awareness tools** that help non-expert users better understand and monitor the data collected about them and how it is (or might be) used. We identify a set of goals for effective user awareness tools:

1. *Actionability*: Beyond just displaying information about private data collection and use to users, an effective user awareness tool must be actionable — that is, users must be able to do something with the information they learn. Though it can be useful to simply inform users about the amount of data invisibly collected about them to build support for broader efforts to manage such collection, such solutions have limited effect on individual end users at present.
2. *Auditable results*: Once a user takes an action to mitigate data collection or use based on increased awareness, it is important that the user be able to audit the results of his or her action. In other words, users should be able to answer the question: “Are my tools, actions, and mitigation strategies actually doing what I expect?”
3. *Attribution*: An effective user awareness tool should allow users to attribute data collection and use to the specific entities responsible. For example, when multiple third-party trackers are loaded on a web page, an effective tool would allow users not just to identify their presence but to trace back particular page content (e.g., ads) to the responsible third party. This attribution helps with both actionability and auditable results, as it helps users understand who is (or is not) doing what.
4. *Awareness about use, not just collection*: We must help users understand not just what data is collected about them but also the potential uses of that data. We cannot expect that non-expert users will be able to extrapolate all possible implications of revealing or allowing certain data to be collected, particularly when multiple third parties collecting data interact in unexpected ways. Thus, our user awareness tools must help users understand and anticipate these implications in order to help them make informed decisions about which data they are willing to share with whom.

There are many aspects of personal data on the Web that are interesting to reveal. For example: can we build tools to reveal to users how their data is leveraged to target ads or recommendations, whether shopping or mortgage sites are using their browsing histories or Facebook profiles to adjust their prices, whether their purportedly encrypted email service is actually decrypting their emails and using the data for its marketing purposes, or whether a service shares their data with third parties – and then how those third parties use the data? For each case, can we reveal exactly which specific data item (or items) that they share with their services – such as emails, documents, locations, or previously visited websites – trigger the specific ads, recommendations or prices? Such

visibility, we believe, would be beneficial to the end users to better understand the implications of their online actions, as well as to assess the effectiveness of any defenses they apply.

Unfortunately, constructing user awareness tools that reveal these and many other potentially interesting aspects about the data’s journey on the web is extremely difficult today, due to a lack of scientific methods to both *detect* data collection and use and *surface* it to end users in effective and actionable ways. For example, a number of tools exist that visualize third-party web trackers (e.g., Ghostery [46], Lightbeam [67]). While these tools can help users understand how many trackers they encounter in their browsing experience, and allow users to block individual trackers, they lack desirable properties including attribution — that is, users may know that a tracker is present on a webpage, but not which parts of the page were affected, e.g., which ads were placed by that tracker. The lack of attribution also limits the auditability of effectiveness, as it can be hard for non-expert users to verify that anything is different when a tracker is reported blocked. Finally, hardly any tools exist today, which can expose to the users how their data is being used by the services that collect it. A few efforts have recently been made (e.g., AdReveal [62], Bobble [89], AdFisher [30], and our own XRay system [60]), but they are all primitive in both detecting data use by Web services and effectively surfacing that information to the end users.

Thus, **our specific goal is to develop not only the first *effective and actionable user awareness tools* that reveal specific aspects of personal data collection and use on the web, but also the science and infrastructural support for building many such tools in the future.** More specifically, as part of this program, we will develop *Hubble*, an extensible, generic, and scalable infrastructure that will provide the necessary scientific methods and programming abstractions to facilitate the building of a new generation of user awareness tools for the web. Hubble’s two main scientific contributions are: (1) providing an *extensible, scalable, and dynamic architecture* leverages statistical methods in unique ways to accurately detect tracking, targeting, personalization, and discrimination in black-box services based on observations of differentiated user profiles, and (2) providing primitives for *effectively surfacing to end users* information about detected data collection and use.

To drive Hubble’s design, we will develop and evaluate at least four user awareness tools, which leverage and inform Hubble’s programming abstractions to detect and visualize various aspects about online data collection and use for targeting, personalization, and discrimination. The specific tools are: (1) *CollectionObservatory*, a tool that detects and visualizes third-party web content that invisibly collects information about users’ browsing behaviors; (2) *AdObservatory*, a tool that detects and visualizes how third-party web trackers leverage the information they collect about the users — such as visited pages, Facebook Likes, or explicitly shared information — to target ads at them; (3) *DiscriminationObservatory*, a tool that detects and visualizes personalized content present on arbitrary websites, with a particular focus on personalized prices or offers on e-commerce, lending, and mortgage websites; and (4) *LocationObservatory*, a tool that detects geo-based targeting and interpret it in terms of demographic attributes.

If successful, our work will lay the first scientific foundations and technology for comprehensive tracking of data collection and use within and across the Web. We foresee multiple domains of impact for our technology. First, by increasing user awareness of how their data is being used on the Web, we hope to make users more mindful of service selection and usage. Second, by enabling robust and scalable transparency tools, we can empower privacy watchdogs — such as journalists, Federal Trade Commission (FTC) investigators, consumer protection agencies, or internet freedom groups (e.g., EFF) — to keep this giant, complex, and ever-changing Web in constant

check to discover any abuses. Third, by enabling transparency at scale and from the exterior, we hope to usher in a new era of voluntary transparency and responsible data behaviors by the web services themselves. Fourth, we believe that our work can integrate well with personal data protection technologies that will be developed as part of the Brandeis program, including TA1 and TA2 technologies. We discuss our vision of such integration in Section 3. Finally, we believe that our technologies will be applicable more broadly to track how sensitive information – be it users’ personal data, proprietary enterprise information, or classified defense data – is being used (or abused) by the parties that obtain it (such as web services, partner enterprises, or foreign governments). We thus expect that extended versions of Hubble could be applicable to use cases of national importance beyond protecting and increasing end-user privacy on the Web.

3 Collaborative Research Team Concept

The transparency and awareness tools we propose to build and evaluate will integrate very closely with other TAs. In particular, our tools can help incentivize adoption of TA3 systems that are built using protection mechanisms developed by other TA1 and TA2 performers; this is achieved by effectively informing users of privacy implications/threats of systems driven by user data. Through extensive user studies, we will evaluate how users’ mental models of digital threats and their behaviors are impacted by increased transparency and awareness of risks, and by the presentation of viable alternatives to existing privacy-risky systems. To support this effort, we will require specifications of how TA1 and TA2 protection tools can prevent specific privacy risks that we may expose via our transparency and awareness tools, and how they are deployed in TA3 systems. We will also require specifications of how TA3 systems make use of sensitive user data. For instance, does a privacy flag in a web browser mask sensitive user attributes such as race, religion, and sexual orientation from being exposed to advertising networks via browsing behavior? Can the protection tools prevent recommender systems from leaking user-supplied preferences to other users? If we can expose such specific privacy risks in existing systems using our tools, and at the same time guarantee that alternative privacy-hardened systems mitigate those risks, then we will be in a better position to inform end-users of concrete risks and viable alternatives. Our findings will also help assess the effectiveness (or ineffectiveness) of TA1 and TA2 protection tools, and will ultimately help iterate on and improve those designs.

4 Technical Plan

Our project will develop both the first *tools* and the first *extensible, scalable, and robust infrastructure* needed to track data use in the uncontrolled Web. While others have previously studied various specific data uses (e.g., price discrimination in Orbitz [34, 50], coarse-grained ad targeting studies [13, 63], or advertising discrimination studies [83]), to the best of our knowledge we are the first to actively seek generic, robust (a.k.a., accurate), and scalable systems to track personal data use on the Web.

Our specific plan involves efforts in three thrusts, which we will execute in parallel. First, we will develop the *Hubble infrastructure and programming abstractions*; it will provide a set of highly reusable and scalable components that will facilitate the building of transparency and user awareness tools to lift the curtain on how personal data is being used. Second, we will build a set of robust, scalable, and usable *transparency and user awareness tools* that leverage those abstractions and enable users, journalists, and investigators to obtain visibility into Web services’ data uses. Third, we will leverage these tools to run *measurement studies* of various data-driven platforms,

such as targeted advertising ecosystems, online trackers, and online price discrimination. These studies will increase awareness, and perhaps help uncover examples of data mistreatments, which will provide the grounds for an informed societal argument on the need for increased voluntary transparency by the services. Moreover, our studies will evaluate the impact of transparency tools on end-user mental models of the privacy implications of online data sharing.

4.1 Thrust 1: The Hubble Transparency Infrastructure and Abstractions

Hubble and its abstractions will support the development of transparency and user awareness tools that reveal aspects about data use on the web. More specifically, Hubble will support the development of any tool that aims to reveal which specific data *inputs* – such as emails, documents, Facebook Like’s, or previously visited websites – are being used to target which specific service *outputs*, such as advertisements, recommendations, or prices. Hubble offers an infrastructure and programming abstractions that can identify targeting at great scale (in the number of inputs, outputs, and services that are being audited), with solid statistical guarantees, and with privacy-preserving properties. Examples of tools that could benefit from Hubble support include AdObservatory, DiscriminationObservatory, and a number of web transparency tools that others have recently built (e.g., [29, 49, 50, 65, 90]).

4.1.1 The Hubble System and Core Abstractions

At the highest level, Hubble will operate as follows. To reveal which specific inputs (e.g., emails or visited websites) in a user’s profile were used to target a particular output O that the user is shown (e.g., an ad), Hubble relies upon *observations* of that same output O in the context of other user profiles with different sets of inputs. For example, if an ad O is often shown to user profiles with a particular website W in their histories, but is never shown to profiles that lack W , then it is plausible that the ad targets that website. Hubble applies statistical methods to infer such *targeting hypotheses* of outputs on specific inputs or combinations of inputs.

The decision of what inputs and outputs are interesting to relate, as well as how the observations required to make the targeting inferences are obtained, is entirely dependent upon the developer of the transparency tool. The observations, for example, could be obtained directly from a population of end-users (e.g., through a browser plugin that monitors users’ browsing and the ads they receive and reports them to Hubble) or through a controlled experiment where user profiles are managed automatically and assigned inputs in a controlled way so as to maximize inference power. Hubble supports both use cases, and transparency tool developers decide what use case they wish to support in their own tools. Regardless, Hubble abstracts out these tool-specific aspects into a core abstraction: *a targeting experiment* (or simply *an experiment*). An experiment specifies the inputs, outputs, and differentiated profiles, as well as an observation collection procedure, along with the type of inference to apply (e.g., correlation, causal, and others).

A transparency tool is constructed as *a workflow of experiments* that build upon each other’s findings to reveal increasingly complex aspects about online targeting. For example, AdObservatory reveals not only reveals which websites in a user’s browsing history trigger specific targeted ads but also attributes each targeting to a specific tracker that most likely was responsible for it. While Hubble will not impose a particular workflow structure, we find one design pattern particularly useful in practice and propose to develop programming abstractions to support it.

Fig.1 shows this pattern and two abstractions designed to support it. The developer structures her tool as a set of large-scale, survey experiments each followed by several finer-grained, reactive

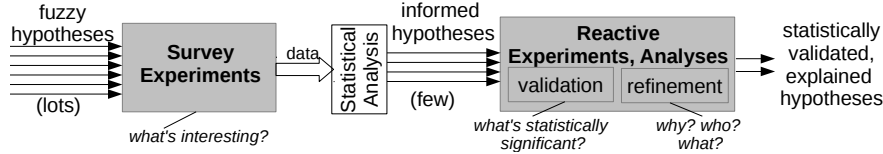


Fig. 1: Example Experiment Workflow in Hubble-based Transparency Tool.

experiments. Hubble’s survey experiment abstraction lets the developer simultaneously evaluate many possible targeting hypotheses, using powerful ideas from compressed sensing [16, 32] to maximize statistical efficiency with limited and sparse observations. Data from the experiment (reports of output observations in particular user profiles) feeds into the statistical analysis engine, which yields a set of *plausible, informed hypotheses* (specific inputs or sets of inputs that appear targeted by specific outputs).

Several of these hypotheses may have confidence scores that are high enough to suggest some effect but perhaps not high enough to reveal to an end-user (erroneous targeting may affect tool credibility). These plausible hypotheses are hence used to trigger a set of follow-up experiments, called *reactive experiments*, that focus on specific hypotheses and attempt to either boost their confidence (*validation experiments*) or provide a more detailed investigation (*attribution experiments*).

Validation experiments can typically be less statistically complex than survey experiments and thus afford more statistical power. For example, a validation experiment may focus on just the specific inputs that have are believed to have triggered the output, and this number may be far fewer than the original number of inputs from the survey experiment.

Attribution experiments are also informed by survey results, and typically pipelined after the validations, and they attempt to pose more in-depth questions about the plausible hypothesis. For example, in AdObservatory, we may follow an initial survey experiment with *attribution experiments* that considers different trackers that may be responsible for ad targeting.

Reactive experiments are defined by the developer by implementing Hubble’s API, which lets programmers define under what conditions certain experiments should be run. Hubble executes the workflow in real-time according to the developer’s specification and returns a set of statistically validated and explained targeting hypotheses. Section 4.2.1 shows a concrete example of the experiment workflow that we plan to leverage in AdObservatory.

4.1.2 Statistical Correlation and Causal Inference

The proposed Hubble infrastructure requires mechanisms for both generating and validating plausible targeting hypotheses. The possible causes for ad targeting and tracking in a given system are myriad, and it is intractable—for both human users and computational methods—to exhaustively consider all of the possibilities. Therefore, it is critical to identify and develop methods that efficiently search for the most likely causes, properly evaluate these potential causes, and then succinctly report reliable results in an interpretable fashion. Unfortunately, many existing techniques designed specifically for finding causes of ad targeting in various settings (e.g., [30, 60, 90]) are generally fragile, inflexible, and do not scale with large numbers of potential targeting hypotheses.

As part of Hubble, we will develop a rigorous and scalable statistical methodology for generating and testing targeting hypotheses based on Hubble’s primitives for conducting *randomized experiments* based on synthetic user profiles—which permit strong causal findings of targeting—as well as using *observational data* of real user profiles. These findings will help inform users of the privacy implications of exposing sensitive information to online systems/trackers.

Basic approach to generating targeting hypotheses. We will first develop a method based on linear regression to discover putative targeting hypotheses from experimental data collected by Hubble. A linear regression model posits that a real-valued *output variable* y is determined by a linear combination of p *input variables* $\mathbf{x} := (x_1, x_2, \dots, x_p)$, plus a random mean-zero noise ε . (Categorical variables are expanded using “dummy variables”.) The linear model is written as $y = \sum_{i=1}^p w_i x_i + \varepsilon$, where $\mathbf{w} := (w_1, w_2, \dots, w_p)$ is the *regression coefficient* vector. Given n vectors of input values $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ together with corresponding output values $y^{(1)}, \dots, y^{(n)}$, the goal of linear regression is to estimate the regression coefficients \mathbf{w} .

In a basic Hubble application for a particular service (e.g., ads targeting), each vector of input values corresponds to a user profile. The variables in \mathbf{x} correspond to either $\{0, 1\}$ -valued indicators for the possible targeting inputs—such as the websites that the user has visited that might be used to target ads—or uncontrolled variables associated with a user profile (e.g., time-of-day of experiment, IP address of client used). For each user profile created, Hubble randomly and independently assigns a value to each targeting input, and also records the value of the uncontrolled variables. The output variable y encodes a particular measured output of an online service or system: for instance, it may represent the number of times a particular ad was displayed to the user, or it may be some aggregate function of all ads displayed to the user. The regression coefficients \mathbf{w} are used to screen the inputs and generating plausible targeting hypotheses based on the coefficients’ magnitudes; the uncontrolled variables are accounted for in the regression and hence may help suppress correlations between the output and irrelevant targeting inputs that would otherwise arise. The targeting inputs associated with large regression coefficients are then, in some combination, hypothesized to be responsible for differences in the observed output. Such a hypothesis will then be vetted using a valid statistical test in a subsequent stage (again, discussed later).

Using sparse linear regression. The basic regression approach for generating targeting hypotheses is not scalable because there are likely many possible targeting inputs—and combinations thereof—that may *a priori* have causal effect on the output. Standard linear regression approaches will require at least as many user profiles as there are possible inputs (i.e., $n \geq p$), regardless of how many of these inputs are actually responsible for the targeting output, and thus result in costly and time-consuming experiments.

We propose to use *sparse linear regression* methods, which are effective at estimating \mathbf{w} even when $p \gg n$, as long as \mathbf{w} is sparse—i.e., has only a few non-zero entries. This sparsity assumption entails that only a few input values are, in combination, correlated with the output. A well-established method for sparse linear regression is Lasso [86]. Under certain conditions on the n input vectors, which we ensure are likely to be satisfied *by construction* of our profiles, Lasso accurately estimates \mathbf{w} as long as $n \geq O(k \log p)$, where k is the number of non-zero entries in \mathbf{w} [15]—i.e., the number of input variables potentially correlated with the output. In fact, this collection of $O(k \log p)$ input vectors supports the *simultaneous* estimation of multiple coefficient vectors for different outputs (e.g., different ads); this remarkable phenomenon (related to compressed sensing [16, 32]) enables highly scalable experiments for generating targeting hypotheses.

Validating targeting hypotheses. To verify whether a targeting hypothesis is valid, we propose to use a two-stage protocol commonplace in statistics and machine learning. We employ two groups of user profiles: the first group (“training set”) is used for generating plausible targeting hypotheses, and the second group (“test set”) is used solely for testing the hypotheses. We will use tests that provide measures of statistical significance in the form of *p-values*. We note that when the

input values involved in a targeting hypothesis are randomly assigned to the user profiles in the test set, then we are able to assess the *causal effects* of these inputs on the output.

Complex targeting hypotheses. In many cases, some additional exploratory experiments can be beneficial for discovering certain complex targeting behaviors. First, targeting inputs may naturally partition into semantically meaningful groups (e.g., health websites, travel websites) that are targeted as a group rather than as individual inputs. Using correlation metrics, we can discover these input groups [9] and then exploit group-level sparsity in our approach for discovering targeting hypotheses; this ultimately may reduce the number of user profiles needed to accurately estimate regression coefficients [55]. Secondly, we may seek out higher-order combinations of inputs that are potentially relevant, and include these combinations in the regression [7]. This would will enable discovery of more complex targeting hypotheses. To support these exploratory experiments, we propose a multi-stage approach whereby groups or higher-order inputs are constructed in a first stage, targeting hypotheses are generated in a second stage, and finally hypothesis testing is conducted in the final stage.

Targeting hypotheses from observational data. Thus far, we have discussed approaches to generating and validating causal targeting hypotheses. However, these methods are based on synthesizing user profiles that may be far removed from any given real user’s profile. Therefore, we believe it will be beneficial to also consider targeting hypotheses based solely on actual users’ profiles. We will explore and evaluate techniques for estimating causal effects from this *observational data* based on assumed casual models [70], and also apply techniques for non-causal targeting hypotheses (e.g., hypotheses of correlations or other measures of associations) that are simply annotated with a familiar correlation-vs-causation disclaimer. For these correlation hypotheses, we will also aim to discover latent factors like population stratification structure that may induce or mask correlations between inputs and outputs. Accounting for these latent factors may produce more reliable findings and increase power to discover targeting behavior that only manifest in subpopulations.

4.1.3 Privacy-Preserving Transparency Protocols

Transparency tools so far - including the most generic ones (XRay [60] and AdFisher [30]) - are designed and have been used *in vitro*. Within the confinement of some experimental settings, investigation was made possible using synthetic data that are inspired from real cases but also deliberately distinct. Hubble crosses this gap as it permits web transparency deployed *in vivo* on actual end-users information. One immediate pressing need is to carefully design access to personal data at any stage of the process to prevent from introducing new threats to users privacy. For instance, a user wishes to track the use of the content of her personal email that is hosted within a large online provider. She is in fact well aware that this content is known entirely by this provider, and used internally for data mining. On the other hand, she also expects this provider to have multiple defense mechanisms in place to protect her email from data thieves. To track the use of her information she is willing to use a tool built upon Hubble but *only if* she can has a guarantee that doing so does not increase that risk. In another scenario, multiple users wish to participate within a collaborative tools to answer a similar data usage issue, but they do not wish to reveal to each other their personal information. Addressing those concerns require us to build transparency mechanisms that can extract how various data items present in a users population are used, without revealing any of those users individual and specific record.

The Hubble architecture includes a series of tools and steps that we conduct to address that need. First, we note an important design choice: Hubble is *not* meant to keep data private *from*

the service that is to track. Hence we do not allow a user to preemptively conduct an experiment to test how her data *will* be used would she choose to disclose it to the system. Note however, that general information that Hubble make available on the use of multiple sensitive information like race, sexual orientation, marital status within services can inform a user’s choice before. Our design choice stems from multiple reasons. First, in all the motivating scenarios we consider, the data is already shared with this service. Second, Hubble necessarily interacts with the service to perform experiment and this poses extraordinary challenges to do so without revealing the inputs that are tested. Third, we argue that if the user considers that merely revealing her information to the provider is already too much of a risk, other privacy techniques which focus on any generic data use would be more appropriate. As a consequence of our choice, we now primarily focus on the two other factors of risk: Hubble itself through its building blocks and APIs, and the tools built on top of Hubble.

Three needs and models of a privacy preserving transparency architecture. Our design choice (discussed immediately below) leverage our first experience with privacy preserving transparency to address three multiple dimensions:

- **strength of the privacy protection:** Ideally all personal data remains undecipherable from the way they are accessed. A less stringent requirement is that, although individual data item can be read in raw form (*i.e.*, attributes like gender, age, webpages visited), those cannot be linked to the same user.
- **client-side and infrastructure-side performance cost:** in order to maintain an acceptable accuracy level, experiments may run on larger set of examples, adding to the cost of Hubble and client may have to perform costly data transformation.
- **flexibility of the design:** Due to the stringent constraint mentioned above, design choice may restrict which data and targeting type can be tracked. Ideally we would like to compare the cost of custom-made solutions to increasingly general abstractions that provide similar protection and accuracy guarantees.

For example, in a recent work, developed in parallel to Hubble outside of the scope of this project, we explored how to track keyword based targeting inside a radical deployment model that assumes users trust no code except the one they execute. Asking each participant to run experiments under their control is prohibitively expensive. Hence we turned to a custom made design that leverages efficient zero knowledge protocols for vector addition and property testing [33], and privacy preserving statistics [39]. That architecture offers some formal guarantees, at significant performance cost for both users and our infrastructure, and it offers no flexibility.

While we hope to leverage lessons learned from this experiment, our preliminary experience motivates us to follow a different design choice for Hubble. First, even in this restricted trust model, users have to assume that the code behave correctly, which is hard to guarantee in general. Ideally, the user should not trust any of our custom-made code at all. Second, we wish to reduce the cost on the user side even further. Finally, our design ought to embrace naturally multiple data formats and uses. Hubble handle those requirement by a new departure that redefines the trust model and reformulates our objective.

Transparency with trusted statistical database. First we precise our new trust model: we assume that, in addition to the provider, and Hubble, a third entity exists to which users can donate data in full confidence. However, this entity is *not* a transparency agency and it may not even be aware

Workflow: E1 → E2 → E3

E1: Website targeting experiment: <ul style="list-style-type: none"> – inputs: websites with ads and/or sensitive material – outputs: ads observed – uncontrolled_vars: time, ip. – const: pages visited on websites. 	E2: Targeting validation: <ul style="list-style-type: none"> – registered after all E1 hypotheses – let: w_in = website ad targets; w_out = website where ad appears. – inputs: w_in, w_out. – outputs: ads observed. – uncontrolled_vars: time, ip. – const: always visit all pages w_out. 	E3: Tracker attribution: <ul style="list-style-type: none"> – registered OnNewHypothesis(ad, w_in, confidence) if confidence ≥ 99%. – inputs: trackers on w_in, w_out. – outputs: ads on w_out. – uncontrolled_vars: time, ip. – const: always visit w_in, w_out.
--	--	---

Fig. 2: AdObservatory Experiment Design. First conducts a large scale survey for cross website ads, then runs validation experiments to improve confidence on the targeted ads, and a tracker stage to attribute ads to specific trackers.

of the existence of Hubble. This “data bank only deals with the maintenance of a database and it provides answer to queries whose disclosure guarantees at anytime users privacy according to a predefined level. One first immediate advantage is that each user has minimum performance cost: it only need to upload its data once and later access Hubble to retrieve the results¹. Above all, the key advantage of this architecture is a unique *flexible modularity*. No matter how respectively Hubble and the data bank decides to operate, as long as Hubble can prove the accuracy of its result and the bank operates correctly, transparency and privacy are achieved. The key challenge, however, is to prove that their correct operation is compatible.

From a users standpoint, the combined system offered an ideal case: The accuracy of a transparency engine can be proved and externally tested. The soundness of a privacy guarantee can even be certified by brands. In fact, following the recent trend to consider personal data as a new class of digital asset [3], multiple solutions exist as research prototype [31, 68], grassroot initiatives `lockerproject.org` and even commercial services `aircloak.com`. Whether or not such solutions become successfully deployed and widely used by online services is not our primary concern. It only suffices that one such service is available to complement Hubble. Moreover, transparency is one of the leading argument encouraging users to rely on those. As part of this project, we will evaluate how Hubble can be made fully compatible with the service offered by `aircloak.com`. We obtained, for the scope of this research, a free access to this service for testing purpose, and we obtained agreement to offer free accounts for volunteers in our first transparency trials. Note that our research will consider the merit of this architecture more generally as described below.

Research hypotheses and methods. We will work under the generalized definition of privacy guarantee [?]. Assuming all inputs databases D_1, D_2, \dots are indexed by j , a databank is private if its return mechanism $\mathcal{M}(\cdot)$ satisfies

$$\forall i = 1, \dots, m, \sum_j c_i(j) \mathbb{P}(\mathcal{M}(D_j) = \omega) \leq 0.$$

where m and $c_i(j)$ can be chosen to fit a privacy definition preliminary accepted. (e.g., ϵ -differential privacy requires $\mathbb{P}(\mathcal{M}(D_j) = \omega) - e^\epsilon \mathbb{P}(\mathcal{M}(D_j) = \omega) \leq 0$ whenever j and j' are identical except for one record).

4.2 Thrust 2: Transparency and User Awareness Tools

4.2.1 AdObservatory: Revealing Targeting in Online Advertising

The first tool we propose to build atop Hubble is *AdObservatory*, which leverages Hubble’s abstractions to reveal to the users how they are being targeted by online advertisers. For each ad that a user encounters while surfing the Web, AdObservatory, a browser plugin, will tag the ad with two pieces of information: (1) which specific website(s) in the user’s browsing history that caused that ad to be shown and (2) which specific tracker(s) that witnessed the users’ visits caused the ad to be shown. These pieces of information correspond directly to the goals we established for user awareness tools (see Section 2): (1) enables *targeting awareness* and (2) provides *attribution*.

Fig.2 shows an experiment workflow that we might use for AdObservatory. It consists of three Hubble experiments: E_1 is a broad survey experiment to formulate rough targeting hypotheses for each ad; E_2 is a smaller experiment to validate the targeting hypotheses generated by the survey and prune out ad erroneously labeled as targeted; E_3 is an attribution experiment to determine which specific trackers contributed to the targeting of cross-domain ads discovered by E_1 and confirmed by E_2 . In more detail, E_1 aims to identify what ads are targeting from a huge range of possibilities. In E_1 each input is one of hundreds or thousands of websites in a user’s web history. The data collection is either an automated browser (e.g., using Selenium) in the case of controlled experiments using AdObservatory, or a plugin running in users’ browser. Regardless, the data collection procedure records as Hubble outputs all display ads observed while visiting web pages. To detect ads on arbitrary pages, we will modify AdBlock to report any identified ad but does not disable it [6]. In addition to collecting display ads, the data collection also records all trackers detected on each site for use in future experiments in the workflow. To detect trackers, AdObservatory will leverage the CollectionObservatory tool we will develop as part of this project. AdObservatory will use Hubble’s statistical correlation and causation methods (§4.1.2) to identify which output display ads target which input sites.

E_2 aims to validate targeting hypotheses from E_1 in a more rigorous and controlled fashion. E_2 creates a group of n profiles, half of which get assigned the targeted website. In addition, all of the profiles get assigned the websites on which the ad appeared; excluding the targeted website if the ad appeared there too. Since, our input is the presence of the targeted website in a profile, E_2 is restricted to ads that appear on at least one website other than the targeted. The ads and their respective groups of site validated as targeted in E_2 will be used in E_3 .

E_3 is similar to the first two experiments and uses the same groups of sites as E_2 but uses trackers collected in E_1 as inputs rather than sites. Using the standard Hubble assignment mechanism each tracker is randomly assigned to half of the accounts. The data collection worker used in E_3 drives the profile to all sites in the assigned group blocking trackers all trackers no allocated to that profile.

4.2.2 DiscriminationObservatory: Revealing Online Price Discrimination

A second tool that we propose building is *DiscriminationObservatory*, a tool that leverages Hubble to reveal to the users how arbitrary websites are personalizing their content based on their personal data. A specific use case we are aiming to support is to reveal price or offer differentiation on eCommerce, mortgage, and loan websites based on users’ personal information, such as Facebook or Google+ profile information, or web histories purchased from trackers. Many websites today

¹Note that this last step should not be performed too naively, as a user should not reveal her inputs while asking for certain transparency results, but there are multiple efficient ways to do that.

leverage the single-signon capabilities of a handful of giant-scale services, such as Facebook Connect and Google OAuth, to authenticate their users. Upon authenticating a user through Facebook or Google, the websites obtain access to various aspects of a user’s profile on these services, and the level of access depends upon the permission level that the website asks for. These pieces of information are often used by the websites to personalize content. For example, Pinterest leverages Facebook Connect to authenticate its users; it requests access to the friend list of a user (among other things) and uses that list to personalize the content it recommends its users. Many other websites do this, and there has been speculation in the media recently that mortgage, loan, and other eCommerce websites might soon start using Facebook Likes and other social information to present the users with differentiated quotes on their websites [2]. At present, no one knows whether any such websites apply such differential treatment, and (worse) no one can find out, because there are no scalable, robust, and generic tools that can identify this kind of behavior in the wild.

Our goal in DiscriminationObservatory is to *detect* such differential treatment on arbitrary websites and *surface* sufficient information to the end-users and privacy watchdogs. End users may leverage this information to inform their decisions about the offers they receive. Privacy watchdogs can use DiscriminationObservatory to search for websites on the web that discriminate against protected user categories, such as specific races, ethnic groups, or genders. To first order, DiscriminationObservatory will leverage Hubble to obtain and analyze the contents of websites of interest (the Document Object Model, or DOM) from the vantage points of users with differentiated profiles. It will compare the contents of the websites at DOM tree level to identify differences that are consistent with differences in the user profiles. Finally, it will highlight visually any DOM portions that receive differential treatment based on various aspects available in social profiles (e.g., gender, relationship, Likes, friend list, etc.). We will leverage known algorithms for discerning differences between DOM trees (e.g., [1, 42, 91]).

4.2.3 CollectionObservatory: Revealing Third-Party Content and Tracking

A third tool we propose to build is *CollectionObservatory*, a comprehensive tool to detect data collection and web tracking. CollectionObservatory is valuable both as a user awareness tool itself (to help inform users about third-party content on the webpages they visit collecting information about their browsing behaviors) as well as to support our other user awareness tools, such as AdObservatory and DiscriminationObservatory. For example, AdObservatory will leverage CollectionObservatory to identify those trackers that collect and use users’ information to target ads at them. In prior work we developed a more limited browser-based web tracking detection and measurement platform which detects only cookie-based tracking [78, 81]. CollectionObservatory will build upon our experience but move significantly beyond it to (1) detect web tracking behaviors of much more diverse and subtle types and (2) provide effective user-facing visualizations of the observed behaviors. The key contributions in CollectionObservatory will thus be: (1)

1. *The most comprehensive web tracking detection:* CollectionObservatory will automatically detect a wide range of web tracking behaviors, including not only well-understood trackers based on browser cookies but also stateless fingerprint-based trackers, which use browser and machine fingerprints to re-identify users, and more esoteric tracking mechanisms (e.g., cache-based, Flash cookies, etc. [59]).
2. *Effective user-facing visualization and awareness:* CollectionObservatory will surface and visualize third-party content and data collection for users in a way that is effective and actionable, helping users take control of the collection and use of their web browsing behaviors.

3. *Research platform:* CollectionObservatory is intended as a platform, allowing other researchers to adapt and build upon our tracking detection capabilities, visualization primitives, and user-facing actions to construct new user awareness tools. Indeed, our own tools (e.g., AdObservatory, DiscriminationObservatory), will leverage CollectionObservatory.

Detecting fingerprint-based trackers. Our previous work, TrackingObserver [81], detects primarily cookie-based tracking that explicitly store state in the user’s browser. In CollectionObservatory, we will extend the scope of this automatic detection to include additional tracking behaviors, primarily *fingerprint-based trackers*. Fingerprinting-based trackers re-identify users based on unique combinations of attributes such as IP address, user agent, installed fonts and plugins, etc [35]. While researchers have explored how fingerprinting works and conducted limited measurement studies of specific fingerprinting techniques or known fingerprinting libraries (e.g., [4, 5, 69, 92]), there has been no extensive non-blacklist-based study of fingerprinting in the wild nor a user-facing tool to detect these behaviors. Implementing fingerprint-based tracking detection in CollectionObservatory, e.g., via hooks on the JavaScript APIs commonly used to generate fingerprints, would allow us to perform a similar study for these trackers. We will conduct a measurement study of tracking on a large number of popular and less popular websites, including from different vantage points (e.g., from different geographic locations). Ultimately, these findings will inform a user awareness tool for web tracking, described below.

Revealing third-party web content. Several tools exist to reveal which third-party trackers are loaded on a given web page, but (as described above) none of these tools localize those trackers on the page. That is, a user can learn that `doubleclick.net` was contacted as the page was loaded, but not which, if any, ads on the page were served by `doubleclick.net`. Similarly, a user cannot easily answer “where did this ad come from?” for a given ad, since even ads loaded from a particular domain may have been placed there by a different third-party (typically an advertising network) [78]. Indeed, some ads might not even have been intended by the web page developer, such as those injected by malicious browser extensions [8]. We propose a tool to identify third-party content on a page and attribute it to its source; achieving this requires addressing a number of technical challenges, including identifying content modifications on the first-party page that result from third-party scripts. We plan to integrate this tool with CollectionObservatory, and envision that it can be used to bootstrap both a user study of attitudes towards and expectations surrounding web tracking (see Section 4.3) as well as a measurement study of third-party content on the web.

Full-fledged web tracking transparency tool. Building on the above and on other aspects of the Hubble infrastructure, and informed by the user studies we describe in Section 4.3, we will ultimately extend CollectionObservatory into a full-fledged web tracking transparency tool for end users. In addition to providing useful visualizations to users about how their information is collected and used as they browse the web, this tool will provide useful, actionable, and verifiable changes that users can make to improve their privacy. We will release this final version of CollectionObservatory as open source, and we will deploy the tool publicly, ideally as part of an existing tool (e.g., as part of the Electronic Frontier Foundation’s Privacy Badger tool [36]), as we have done with ShareMeNot [80]) in the past. This deployment will serve as a field study of the tool, which in turn will inform additional iteration on the tool itself.

4.2.4 LocationObservatory: Revealing Privacy Implications of Location Tracking

[augustin]

xxx

4.3 Thrust 3: User Studies and Transparency Tool Measurements

[Augustin: Add IRB note.]

xxx

To maximize the effectiveness of the transparency infrastructure and the user awareness tools that we build, it is critical that we understand users themselves. To this end, our proposed work will include user studies of two types: (1) user studies to help us understand *users’ existing mental models and attitudes*, and (2) user studies to help us *evaluate the effects of our tools*. We will work with our institutions’ human subject review boards to obtain IRB approval before conducting any studies involving human subjects.

User Studies for Existing User Mental Models and Attitudes. Our transparency and user awareness tools aim to close the gap between users’ existing mental models and attitudes with respect to the privacy of their data and the reality of what today’s applications and services collect and use. To achieve this, we must first understand what users already know or believe about the collection and use of their private data. Prior work has studied users’ mental models and attitudes in contexts such as targeted advertising (e.g., [61, 64, 71, 87]); we propose to extend that work here, and to update the findings for current users and systems.

Example 1: Reactions to Ad Targeting. As one example, we detail a user study intended to help inform our transparency and user awareness tools for web tracking and targeted advertising. We ask: what are users’ mental models about ad targeting? How will they react upon learning that a particular ad is targeted at them? To explore this question, we will initially design a study in which we post ads (e.g., on Facebook or via Google ads) targeted at specific—possibly sensitive—keywords. The content of our ads will inform the person viewing them about the targeting, e.g., by revealing the keyword that was used to target that particular ad. Clicking on the ad will direct the participant to a page with additional information about targeted advertising and about our study, including several survey questions to help us evaluate the participants’ reactions to (1) learning about the targeting as well as to (2) the targeting itself. As a second part of this study, we will conduct a user study with our AdObservatory tool to study users’ reactions to real targeting that they encounter in the wild.

By evaluating and comparing participants’ reactions to different targeting keywords, our results can help motivate and inform our transparency tools, which may in turn motivate changes within targeting systems themselves. For example, if we find that users are comfortable with ads targeted at debt-related keywords but not cancer-related keywords, we might recommend that ad targeting companies stop targeting cancer, or offer an opt-in to such “sensitive” topics. More broadly, studies such as this one will help us understand the notion of “sensitivity” — how much does it depend on the user, what kinds of things are uniformly “sensitive,” etc.? These findings will ultimately inform our transparency and user awareness tools as well as others working in this space.

Example 2: Blah blah. [Something from Augustin somewhere around here?]

xxx

User Studies to Evaluate our Tools. In addition to user studies aimed to teach us about users in general, we must also evaluate the effectiveness of our tools – AdObservatory, DiscriminationObservatory, CollectionObservatory, and LocationObservatory– with real users. Our goals here include (1) validating that the information our tools surface to users is comprehensible, (2) studying immediate user attitudes and reactions to learning this information, and (3) observing and tracking sustained user behavior changes (or lack thereof) in response to learning this information.

These studies will take several forms throughout the design of each tool, beginning with limited usability studies of preliminary designs, followed by more in-depth studies to evaluate the

Component	Sub-tasks	Responsible PI(s)
Hubble infrastructure	1.1, 2.1, 3.1	Geambasu
Statistical correlation and causation	1.3, 1.4, 2.3, 2.4, 3.3, 3.4	Hsu
Privacy-preserving transparency	1.5, 2.5, 3.5	Chaintreau
CollectionObservatory	1.7, 2.7, 3.7	Roesner
AdObservatory	1.2, 2.2	Geambasu
DiscriminationObservatory	2.2, 3.2	Geambasu
LocationObservatory	1.6, 2.6, 3.6	Chaintreau
User studies	1.8, 2.8, 3.2, 3.6, 3.8	Roesner, Chaintreau, Geambasu
Integration, Evaluation on TA3	1.9, 2.9, 3.9	All PIs

Fig. 4: Team member responsibilities (research areas and subtasks).

effectiveness of our tools to improve user comprehension and to positively affect user behaviors, culminating in full-fledged beta-tests with real user populations. For example, co-PI Roesner has previously released a user-facing anti-web tracking tool (originally called ShareMeNot [80]) as part of the Electronic Frontier Foundation’s Privacy Badger tool [36]. We will use connections like these to iteratively beta-test our tools with large numbers of real users in real contexts.

5 Personnel and Management Plan

The team members are faculty at two institutions: Columbia University and University of Washington. Columbia University will be the Prime Contractor for the project, with University of Washington acting as a subcontractor; the formal agreements are already in place for this project. Roxana Geambasu will be the overall project PI, responsible for general technical direction, coordination and reporting (in addition to conducting a portion of the research). Each

Key Individual	2015	2016	2017	2018	2019
Geambasu	53 h	160 h	160 h	160 h	160 h
Chaintreau	53 h	160 h	160 h	160 h	160 h
Hsu	53 h	160 h	160 h	160 h	160 h
Roesner	53 h	160 h	160 h	160 h	160 h

Fig. 3: Team member commitments.

co-PI will be responsible for one or more component and associated sub-tasks (see Fig. 4). Each faculty member will be responsible for supervising Ph.D. Graduate Research Assistants (GRAs). Each faculty member will dedicate a significant amount of their time to this project (see Fig. 3).

The management structure is relatively flat, with Geambasu the lead PI and everyone else working with each other and under the general guidance of Geambasu. The PIs already have a history of collaboration with each other and are co-advising students. For example, Chaintreau and Geambasu co-authored the XRay paper [60] and are co-advising a Ph.D. student, the paper’s first author. Chaintreau, Geambasu, and Hsu have been working on follow-on technology and are now writing a joint paper for CCS’15 on a related topic. Geambasu and Roesner have already started a collaboration in the space of user awareness studies. The Columbia Co-PIs meet face-to-face almost on a daily basis. To facilitate collaboration with the UW Co-PI, we will have regular meetings over Skype or other technology. We will also organize two physical meetings per year, hosted on a rotating basis among the institutions and/or co-located with the program PI meetings. We will use a wiki and Github for coordination and record keeping. **All code will be made public on Github.**

5.1 Personnel

The PIs span a broad range of expertise: *systems* (Geambasu), *theory and social networks* (Chaintreau), and *machine learning and statistics* (Hsu), and *security and human factors* (Roesner). We will combine this broad expertise in a close collaboration to produce the first scalable infrastruc-

ture for transparency and the first valuable tools for end-user privacy awareness. Following are the biographies of each participant. Section 6 describes the team’s relevant expertise and joint projects.

Roxana Geambasu is an Assistant Professor of Computer Science at Columbia University. She has made research contributions in software systems across a broad range of areas, including operating systems, distributed systems, and security and privacy. One over-arching theme of her research relates to increasing privacy in today’s data-driven world by developing transparency, fairness, and data management tools for programmers, privacy watchdogs, and end-users. Her publications are available at: www.cs.columbia.edu/~roxana. Prof. Geambasu is a member of the Information Science and Technology (ISAT) focus group. For her work in privacy, Prof. Geambasu received a Microsoft Research Faculty, a Popular Science “Brilliant 10” listing, an Honorable Mention for the inaugural Dennis M. Ritchie Doctoral Dissertation Award, a William Chan Dissertation Award, two best paper awards at top systems conferences (USENIX Security and EuroSys), and the first Google Ph.D. Fellowship in Cloud Computing. Geambasu’s research was featured in multiple articles in New York Times, The Economist, NPR.

Augustin Chaintreau is an Assistant Professor of Computer Science at Columbia University since 2010, where he directs the Mobile Social Lab. He designs algorithms leveraging social mobile behaviors or incentives, to reconcile the risk and value of personal data networking. His research addressing transparency in personalization, fairness in personal data markets, efficiency in social information sharing, cross domain data linkage, and human mobility lead to 25 papers in tier-1 conferences (five receiving best paper awards at ACM CoNEXT, SIGMETRICS, USENIX IMC, IEEE MASS, Algotel, some with media coverage such as the NYT blog). An ex student of the Ecole Normale Supérieure in Paris, he earned a Ph.D in mathematics and computer science in 2006, the NSF CAREER Award in 2013 and the ACM SIGMETRICS Rising star award in 2013. He has been an active member of the network research community, serving in the program committees of ACM SIGMETRICS, SIGCOMM, WWW, CoNEXT (as chair), MobiCom, IMC, WSDM, COSN, AAAI ICWSM, and IEEE Infocom, and is area editor for IEEE TMC and ACM SIGCOMM CCR.

Daniel Hsu is an Assistant Professor of Computer Science at Columbia University, and is a member of the Data Science Institute, also at Columbia. His research interests are in algorithmic statistics, machine learning, and privacy-preserving data analysis. His work on interactive and unsupervised learning have yielded the first computationally efficient and statistically consistent algorithms for a number of core estimation and learning problems that were only previously tackled using heuristics or suboptimal methods. Much of his current research focuses on developing scalable and statistically sound learning algorithms for discovering hidden structure in massive data, as well as on the interaction between statistical inference and privacy. He was the organizer of several workshops and tutorials on algorithms for learning hidden variable models at premier machine learning venues (ICML, NIPS); he received a Yahoo Academic Career Enhancement Award in 2014 and the UC San Diego Departmental Dissertation Award in 2010.

Franziska Roesner is an Assistant Professor of Computer Science and Engineering at the University of Washington. She has made research contributions in computer security and privacy, spanning broadly from systems to human factors. Her work involves designing and building systems that address security and privacy challenges faced by end users of existing and emerging technologies. For example, she has made contributions in computer security and privacy in the contexts of third-party web tracking, permission granting in modern operating systems (such as smartphones), secure embedded user interfaces, and emerging augmented reality platforms. A list of her publications is available at: <http://www.franziroesner.com>. Her work on web

privacy included the development of ShareMeNot, a defense for one type of web tracker, which was incorporated into the Electronic Frontier Foundation’s Privacy Badger tool in 2014. For her work in security and privacy, Prof. Roesner received the William Chan Memorial Dissertation Award in 2014, the IEEE Symposium on Security and Privacy Best Practical Paper Award in 2012, a NSF Graduate Research Fellowship, and a Microsoft Research PhD Fellowship.

5.2 Integration and Evaluation

Although each component is led by a particular team member, the PIs will work together as part of a unified team to integrate their components in a coherent system and a useful suite of tools. The end product will be a robust infrastructure that can be leveraged by other researchers and developers to reveal other aspects of personal data treatment on the web. The infrastructure can be used to audit TA3 systems that rely on personal information; these systems will expose the types of personal data that are used as input (e.g., user browsing data, e-mails), and what kinds of measurable outputs the system will produce (e.g., ad placement decisions). Our tools will provide a means to inform users of the privacy implications of using or interacting with such systems.

We will conduct a thorough evaluation of our infrastructure prototype and tools using well-established evaluation metrics and methods. We will assess the scalability of our system as the number of possible targeting inputs and outputs grows. We will also assess the precision and recall of our system in detecting known targeting behavior in systems where ground-truth has been established. As already discussed in §4.3, the user-facing transparency tools will be thoroughly evaluated in comprehensive user studies.

6 Capabilities

PI Geambasu has been working on increasing privacy and transparency in computer systems for multiple years. As part of a DARPA MRC project (MEERKATS), she has CleanOS, a mobile operating system designed with privacy and transparency in mind [82, 84]. Unlike existing mobile OSes, CleanOS manages users’ data carefully so as to (1) minimize exposure of users’ personal data at any time in anticipation of attack and (2) provides visibility post-attack into what specific data might have been compromised.

Geambasu and Chaintreau have recently developed *XRay*, a preliminary transparency infrastructure that reveals data targeting in Web services [60]. The system, which is our preliminary foray into the topic of Web transparency and our inspiration for Hubble, is the very first to accurately reverse targeting in multiple services, including Gmail, Amazon, Youtube, and Google Search. The system, however, is limited in scale, applicability, features, and our personal experience with it. We will address these and more limitations in Hubble.

Geambasu, whose core expertise lies in building scalable, extensible, and robust distributed systems [43–45, 88], has a track record of transition into practice of the systems she builds. For example, Synapse [88], a scalable, heterogeneous-database replication system, has been deployed at Crowdmap, a data-driven marketing startup in NYC, which has been running it in production for about a year with great success. As another example, Geambasu deployed the first security measures in a commercial, giant-scale distributed hash table (DHT) with millions of users. Her defenses alleviated the potential for certain Sybil attacks on that DHT [43].

PI Chaintreau produced research addressing the need for better fairness and transparency in personalization. In addition to the *XRay* system already mentioned above, he has also explored how to design incentive, based on fluid limit of cooperative game theory, proving stability of simple

rewarding schemes in peer-assisted services [66]. More recently, with Telefonica and AT&T, he proposed a solution based on pseudonyms and auctions to redistribute economic value of web-browsing to the users [73], and reveal the revenue made available to online advertising through tracking [48]. Multiple efforts to experiments to those designs with real users passed IRB approval and are ongoing. These results, along with new techniques to increase the transparency of online services, proposes concrete steps to reconcile privacy and the deployment of big data.

In parallel, following a Ph.D. on models of interaction processes within TCP/IP networks [17], his research repeatedly demonstrated the importance of *social properties* and algorithms leveraging those to design large scale systems. (1) how human mobility exhibits statistical [20] and topological [22] properties that characterize the convergence of opportunistic routing protocols; (2) how small world navigation emerges from simple local dynamics of mobility and aging [18], which distributed gossip algorithms can exploit [19]; (3) how populations of mobile users follow spatial models, informing capacity planning [21]; (4) how distributed algorithms can exploit arbitrary processes of contacts between nodes for update [56], caching [57, 72], routing [40, 41], and personalization [58]. These results were collaborations with industry at IBM, Alcatel-Lucent, Intel and Technicolor, and led to 8 filed patents.

PI Hsu works on algorithmic statistics, machine learning, and privacy-preserving data analysis. He has developed several foundational algorithms in the areas of active learning and unsupervised learning. His work on noise-tolerant and statistically-consistent active learning [14, 27, 28] provides an algorithmic basis for adaptive experimental design in classification problems, which we hope to employ in Hubble for scalability. His work on efficient algorithms for learning latent variable models [10–12, 51, 53, 54] provides techniques for capturing hidden structure in data, which can be used to improve the statistical power for finding significant correlations and causal relationships.

In addition to work on active and unsupervised learning, Hsu has studied implications of privacy constraints on statistical machine learning and data analysis, and he has experience in developing scalable learning algorithms for complex regression problems. His work has revealed practically-relevant limitations of requiring differential privacy guarantees on learning algorithms and statistical estimators [23, 24], and also developed new methods that exploit conditions that are favorable for learning when it is available [25]. He also developed highly scalable algorithms for discovering complex variable interactions that are useful in regression [7], as well as methods for exploiting sparse linear regression in the context of complex output prediction [52].

PI Roesner has worked on web privacy topics for several years, focusing on (1) studying and measuring the existing state of web privacy, (2) building tools to enable measurement and other follow-on work, and (3) providing users with visibility into and control over their privacy on the web. Her 2012 taxonomy and measurement study of third-party web tracking in the wild [78] was among the first efforts to deeply understand the web tracking space. As part of this work, Roesner developed *ShareMeNot*, a defense for social media web trackers (such as the Facebook “Like” button). *ShareMeNot*’s techniques were adopted by Ghostery [46], a popular anti-tracking browser add-on, and *ShareMeNot*’s code itself was incorporated into the Electronic Frontier Foundation’s Privacy Badger [36] web privacy tool in 2014. Roesner’s work has also focused on ensuring that the security and privacy properties of systems match users’ expectations in other contexts. For example, she developed *user-driven access control* [77] as a new approach for permission granting in modern operating systems (such as smartphones), by which the operating system is able to extract a user’s permission granting intent from the way he or she naturally interacts with any application. Roesner implemented user-driven access control in *LayerCake*, a modified version of

Android that provides security for embedded user interfaces [74,75]. Her work has also focused on emerging security and privacy challenges in emerging augmented reality and continuous sensing platforms [76,79].

7 Statement of Work

Our effort is composed of one overall task, aimed at developing a complete and demonstrable Hubble prototype and tools. We define a number of subtasks that partition the effort into smaller, independently developed components, which are integrated increasingly with each program phase and evaluated against TA3 systems.

TASK: Objective: Investigate, develop, and experimentally evaluate a Hubble prototype; develop and evaluate user awareness tools built upon its primitives.
General Description: This is our main goal and high-level task, around which a number of smaller tasks (broken down by phase) are organized. We will develop and integrate the individual components, and evaluate the integrated architecture across the full duration of the project.
Responsible Organization and Location: Columbia University (NYC), University of Washington (Seattle).
Exit Criteria: An extensible, scalable, and robust infrastructure system for building transparency tools to increase users' awareness of how their data is being collected, used, and exchanged by online services. A scientific understanding of how such tools can help change user perceptions of the risks involved and improve their mental models of protection techniques that exist or are being developed as part of the Brandeis program. Evaluation in terms of accuracy, scale, performance, and lightweightness are successful on real systems and TA3 systems.
Deliverables: Prototype implementation of Hubble and tools, including documentation and the final project report, quarterly technical progress reports, slide presentations, evaluation data, and other reports per requirements. All source code for Hubble and tools will be released publicly on Github.

7.1 Phase 1 (Months 1-18)

TASK 1.1: Objective: Design and implement basic Hubble infrastructure and tool development API.
General Description: Design early version of Hubble's architecture and developer APIs. The architecture will support single-stage experiments (no validations / refinements). Implement this architecture with basic statistical correlation engine available, stub other unavailable components. Focus on controlled-input use cases.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Infrastructure that reveals input/output targeting by measuring correlation on differentiated profiles. Supports 10s-100 inputs and has precision/recall for detecting targeting of 70-90%.
Deliverables: Early software prototype and design documents.
TASK 1.2: Objective: Design and implement basic AdObserver tool.
General Description: Implement a basic version of the AdObserver tool to exercise Hubble's architecture and APIs. Use AdBlocker to identify ads on arbitrary pages.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Tool that can reveal ads targeted on previously visited websites or other data.
Deliverables: Software prototype and design documents.
TASK 1.3: Objective: Develop basic statistical methodology for testing targeting hypotheses.
General Description: Developing a formal specification for targeting hypotheses as generated by Hubble, together with a methodology for reliable testing of the hypotheses.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Concrete specification of targeting hypotheses, and software tool that computes valid statistical tests at any specified level, incorporated into Hubble.
Deliverables: Software prototype and design documents.

TASK 1.4: Objective: Apply scalable sparse linear regression methods to generation of targeting hypotheses.
General Description: Develop sparse linear regression approach to infer putative targeting hypotheses from data collected by Hubble. Evaluate scalability using simulated targeting mechanisms and real data collected by Hubble.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Scalable and empirically-validated implementation, incorporated into Hubble pipeline.
Deliverables: Software prototype and design documents.
TASK 1.5: Objective: Design and implement basic privacy-preserving transparency protocol.
General Description: XXX AUGUSTIN. Relies on central collection service.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: XXX.
Deliverables: XXX.
TASK 1.6: Objective: Design and implement basic LocationObserver to reveal information that can be inferred from location.
General Description: XXX AUGUSTIN. Evaluate privacy-preserving protocol against alternative designs.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: XXX.
Deliverables: XXX.
TASK 1.7: Objective: Implement fingerprint tracking detection infrastructure in CollectionObservatory.
General Description: Build on TrackingObserver, our prior web tracking detection and measurement platform, to begin developing CollectionObservatory. First, implement detection of fingerprint-based web trackers that use browser and machine fingerprinting techniques to re-identify users. Rather than using a known list of fingerprinting scripts, detect fingerprinting behavior using a measurement of entropy extracted by a potential tracker's JavaScript API accesses.
Responsible Organization and Location: University of Washington (Seattle, WA)
Exit Criteria: Initial version of CollectionObservatory that successfully detects a large fraction of fingerprint-based trackers, evaluated by a comparison with blacklist-based tracking detection tools.
Deliverables: Initial version of CollectionObservatory that detects fingerprint-based trackers.
TASK 1.8: Objective: Conduct user study of attitudes towards targeting.
General Description: Conduct a user study to better understand users' attitudes towards targeted advertising. Target ads using a variety of keywords (including sensitive keywords) and inform users about the targeting in the content of the ads. For participants who click on the ad, debrief them about the study and ask addition survey questions.
Responsible Organization and Location: University of Washington (Seattle, WA)
Exit Criteria: Sufficient participation in the user study to draw statistically significant conclusions.
Deliverables: Conclusions drawn from user study results.
TASK 1.9: Objective: Demonstrate our TA2 technology on a TA3 Research System.
General Description: Integrate basic Hubble and transparency tools prototypes with the Research System(s) developed by TA3 performers. Our tools should be able to identify data uses or privacy attacks built within those systems.
Responsible Organization and Location: Columbia University (New York), University of Washington (Seattle)
Exit Criteria: Successful detection of data use in TA3 Research System.
Deliverables: Software prototypes, design documents, and results from evaluation.

7.2 Phase 2 (Months 19-36)

TASK 2.1: Objective: Extend Hubble and APIs for multi-stage transparency tool designs.
General Description: Incorporate support for multi-stage transparency tools. Support validation and refinement as abstractions for multi-stage tools. Develop API for such tools. Incorporate causal inference building block into Hubble as part of the validation phase. Continue to focus on controlled-input use cases.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: A system capable of validating and explaining its own causal targeting hypotheses. Its evaluated scale will be in the range of 100s-1000s inputs, but we expect its recall/precision to grow thanks to validations.
Deliverables: Software prototype and design documents.

TASK 2.2: Objective: Extend AdObserver and DiscriminationObserver to leverage Hubble’s multi-stage architecture.
General Description: Design and implement using Hubble’s APIs validation and refinement stages for each tool. Run experiments to test and evaluate.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Tools that both scale and validate/explain their own assessments to the users.
Deliverables: Software prototype and design documents.
TASK 2.3: Objective: Develop and evaluate methodology for generating and testing targeting hypotheses from observational data.
General Description: Explore and evaluate techniques for estimating causal effects from observational based on an assumed casual model. Also develop and evaluate correlation hypotheses that do not assert causal implications.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Software tool for hypothesis generation and testing using observational data.
Deliverables: Software prototype and design documents.
TASK 2.4: Objective: Extend sparse linear regression methodology to support complex targeting hypotheses.
General Description: Develop multi-stage methodology to support testing of complex targeting hypotheses with higher-order input interactions. Evaluate this strategy using simulated data and real data collected by Hubble.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Scalable and empirically-validated implementation linear regression approach using higher-order inputs, incorporated into Hubble pipeline.
Deliverables: Software prototype and design documents.
TASK 2.5: Objective: Extend privacy-preserving transparency to avoid trust in a central point.
General Description: XXX AUGUSTIN.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: XXX.
Deliverables: XXX.
TASK 2.6: Objective: Extend LocationObserver to integrate privacy-preserving techniques.
General Description: XXX AUGUSTIN.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: XXX.
Deliverables: XXX.
TASK 2.7: Objective: Measurement study with CollectionObservatory.
General Description: Using fingerprint-based tracking detector in CollectionObservatory (as well as existing capabilities in TrackingObserver from prior work), conduct large-scale measurement study of tracking on the web.
Responsible Organization and Location: University of Washington (Seattle, WA)
Exit Criteria: Conduct a measurement study of tracking on a large number of popular and less popular websites, including from different vantage points (e.g., from different geographic locations).
Deliverables: Measurement study results, including the prevalence and effectiveness of fingerprint-based trackers, a comparison with previous measurement results, etc.
TASK 2.8: Objective: Small-scope user awareness tool that visualizes third-party content.
General Description: Develop an initial user awareness tool for web tracking that identifies third-party content on a webpage and visualizes it for the user. This tool, combined with CollectionObservatory, will serve as a building block for our later, more full-fledged web tracking user awareness tool.
Responsible Organization and Location: University of Washington (Seattle, WA)
Exit Criteria: Software prototype that identifies and visualized third-party content on a webpage.
Deliverables: Software prototype and design documents.

TASK 2.9: Objective: Demonstrate our enhanced TA2 technology on a TA3 Research System. Initial trial of demonstration on a TA3 Existing System.
General Description: Integrate enhanced implementation of Hubble and transparency tools with the Research System(s) implemented by TA3 researchers. Begin integration of our tools with Ta3 Existing System(s), as well as TA1 and TA2 protection-oriented technologies to enable auditing of the effectiveness of their protection.
Responsible Organization and Location: Columbia University (New York), University of Washington (Seattle)
Exit Criteria: Successful detection of data use in TA3 Research System. Our tools should detect data uses in TA3 Research systems.
Deliverables: Software prototypes, design documents, and results from evaluation.

7.3 Phase 3 (Months 37-54)

TASK 3.1: Objective: Extend Hubble to support collaborative transparency scenarios.
General Description: Incorporate statistical correlation building block for uncontrolled inputs to support end-user scenarios. Also incorporate privacy-preserving protocols to limit the need for users to trust Hubble. Run experiments with simulated users to evaluate.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: A privacy-preserving collaborative transparency system where users can submit their inputs/outputs partially and retrieve targeting assessments.
Deliverables: Software prototype and design documents.
TASK 3.2: Objective: Extend AdObserver, DiscriminationObserver to the collaborative use case.
General Description: Port the tools to the collaborative version of Hubble and re-run measurements in a simulated collaborative scenario for evaluation.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Transparency tools now run on observational data of the end users without the need to trust Hubble.
Deliverables: Software prototype and design documents.
TASK 3.3: Objective: Develop and evaluate statistical testing methodology for stratification structure.
General Description: Develop methods for discovering latent population stratification (clustering), together with hypothesis tests that leverage this stratification structure to increase the statistical power to detect targeting.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Software tool for computation of statistical tests.
Deliverables: Software prototype and design documents.
TASK 3.4: Objective: Extend sparse linear regression techniques to use adaptive multi-stage experimental designs, and incorporate statistical testing methods to generate higher-order targeting hypotheses.
General Description: Develop multi-stage methodology for exploiting groups of related targeting inputs and outputs. The group structures are inferred in a first experimental stage, and the subsequently exploited in a second stage using group-sparse linear regression methods to discover group-level targeting hypotheses.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: Scalable and empirically-validated implementation of multi-stage group-sparse linear regression approach, incorporated into Hubble pipeline.
Deliverables: Software prototype and design documents.
TASK 3.5: Objective: Finalize privacy-preserving, collaborative transparency building blocks and integrate into Hubble.
General Description: AUGUSTIN
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: XXX.
Deliverables: XXX.
TASK 3.6: Objective: Finalize LocationObserver tool and run studies of impact of transparency on user actions.
General Description: XXX AUGUSTIN.
Responsible Organization and Location: Columbia University (NYC)
Exit Criteria: XXX.
Deliverables: XXX.

TASK 3.7: Objective: User study of third-party content visualization tool.
General Description: We will conduct a usability study of the previously developed third-party content visualization tool, to understand whether and how the tool is effective with real users: does it effectively convey information to users? Do users take useful actions in response to this information? etc.
Responsible Organization and Location: University of Washington (Seattle, WA)
Exit Criteria: Conduct user study with a sufficient number of participants to generate statistically significant results and inform interactive improvements to the tool.
Deliverables: User study results and iterative improvements to the software prototype.
TASK 3.8: Objective: Larger-scope web privacy user awareness tool.
General Description: Develop a more full-fledged web tracking user awareness tool, integrating functionality from the previously developed third-party content visualization tool and from CollectionObservatory. This tool will be informed by the user studies and measurements and will build on other infrastructure developed in the project.
Responsible Organization and Location: University of Washington (Seattle, WA)
Exit Criteria: Develop a more full-fledged web tracking user awareness tool informed by and building on other aspects of the project.
Deliverables: Software prototype and design documents.
TASK 3.9: Objective: Demonstrate our final TA2 technology on TA3 Research and Existing Systems.
General Description: Integrate final implementation of Hubble and transparency tools with the Research and Existing Systems implemented by TA3 researchers. Finalize integration of our tools with some TA1 and TA2 protection-oriented technologies to enable auditing of their effectiveness.
Responsible Organization and Location: Columbia University (New York), University of Washington (Seattle)
Exit Criteria: Successful detection of data use in TA3 Research and Existing Systems. Successful auditing of effectiveness of other TA1, TA2 technologies with which we integrate, as applied to the same TA3 systems.
Deliverables: Software prototypes, design documents, and results from evaluation.

8 Schedule and Milestones

The Gantt chart below provides a graphic representation of the project schedule at the level of sub-tasks, all of which fall with the one overall task of Hubble, aimed at developing a complete and demonstrable Hubble prototype and tools. The performing organization is indicated via color: blue tasks correspond to Columbia University, green tasks correspond to University of Washington. Program milestones are indicated via bullets, and the duration of each sub-task is provided in the final column of the graphic.

[Someone, please can you generate this gantt chart? I don't know how to make it nice, I only use OpenOffice and it's very primitive. Look at the MEERKATS proposal I sent for guidance. Table ?? contains the timeline data for us.] xxx

Task	Period	PI(s)
Phase 1		
Task 1.1	Months 1-18	Geambasu
Task 1.2	Months 1-18	Geambasu
Task 1.3	Months 1-18	Hsu
Task 1.4	Months 1-18	Hsu
Task 1.5	Months 1-18	Chaintreau
Task 1.6	Months 1-18	Chaintreau
Task 1.7	Months 1-18	Roesner
Task 1.8	Months 1-18	Roesner
Task 1.9	Months 15-18	All
Phase 2		
Task 2.1	Months 19-36	Geambasu
Task 2.2	Months 19-36	Geambasu
Task 2.3	Months 19-36	Hsu
Task 2.4	Months 19-36	Hsu
Task 2.5	Months 19-36	Chaintreau
Task 2.6	Months 19-36	Chaintreau
Task 2.7	Months 19-36	Roesner
Task 2.8	Months 19-36	Roesner
Task 2.9	Months 33-36	All
Phase 3		
Task 3.1	Months 37-54	Geambasu
Task 3.2	Months 37-54	Geambasu
Task 3.3	Months 37-54	Hsu
Task 3.4	Months 37-54	Hsu
Task 3.5	Months 37-54	Chaintreau
Task 3.6	Months 37-54	Chaintreau
Task 3.7	Months 37-54	Roesner
Task 3.8	Months 37-54	Roesner
Task 3.9	Months 51-54	All

9 Cost Summary

Entire Performance Period (Total: \$3,960,419)

	Columbia University (CU) (prime)	University of Washington (UW) (sub)	Category Total
Direct Labor	1,270,820	402,258	1,673,078
Materials ODC	783,786	324,082	1,107,868
Indirect Costs	910,322	269,151	1,179,473
Member Totals	2,964,928	995,491	3,960,419

GFY 15 (Total: \$104,438)

	CU (prime)	UW (sub)	Total
Direct Labor	23,921	7,844	31,765
Materials	1,517	1,196	2,713
ODC	18,594	12,725	31,319
Indirect Costs	32,962	5,679	38,641
Member Totals	76,994	27,444	104,438

GFY 16 (Total: \$885,835)

	CU (prime)	UW (sub)	Total
Direct Labor	297,808	94,277	382,085
Materials	39,200	12,855	52,055
ODC	142,269	42,071	184,340
Indirect Costs	205,205	62,148	267,353
Member Totals	674,481	211,351	885,832

GFY 17 (Total: \$923,294)

	CU (prime)	UW (sub)	Total
Direct Labor	296,940	96,167	393,107
Materials	39,200	12,930	52,130
ODC	146,087	58,067	204,154
Indirect Costs	210,684	63,219	273,903
Member Totals	692,911	230,383	923,294

GFY 18 (Total: \$949,555)

	CU (prime)	UW (sub)	Total
Direct Labor	306,367	97,493	403,860
Materials	39,200	12,984	52,184
ODC	150,019	63,182	213,201
Indirect Costs	216,340	63,970	280,310
Member Totals	711,926	237,629	949,555

GFY 19 (Total: \$898,226)

	CU (prime)	UW (sub)	Total
Direct Labor	294,556	86,783	381,339
Materials	31,683	12,209	43,892
ODC	141,732	68,808	210,540
Indirect Costs	204,745	57,711	262,455
Member Totals	672,716	225,510	898,226

GFY 20 (Total: \$199,074)

	CU (prime)	UW (sub)	Total
Direct Labor	61,228	19,694	80,922
Materials	0	4,926	4,926
ODC	34,285	22,130	56,415
Indirect Costs	40,387	16,424	56,811
Member Totals	135,900	63,174	199,074

References

- [1] Facebook react. <http://facebook.github.io/react/docs/reconciliation.html>.
- [2] Time - lendup.com. <http://business.time.com/2012/11/16/can-a-payday-lending-start-up-use-facebook-to-create-a-modern-community-bank/>.
- [3] Personal Data: The Emergence of a New Asset Class. *World Economic Forum Report*, pages 1–40, Jan. 2011.
- [4] G. Acar, C. Eubank, S. Englehardt, M. Juarez, A. Narayanan, and C. Diaz. The web never forgets: Persistent tracking mechanisms in the wild. In *21st ACM Conference on Computer and Communications Security*, 2014.
- [5] G. Acar, M. Juarez, N. Nikiforakis, C. Diaz, S. Gürses, F. Piessens, and B. Preneel. FPDetective: Dusting the web for fingerprints. In *20th ACM Conference on Computer and Communications Security*. ACM, 2013.
- [6] Adblock plus. <https://adblockplus.org>.
- [7] A. Agarwal, A. Beygelzimer, D. Hsu, J. Langford, and M. Telgarsky. Scalable nonlinear learning with adaptive polynomial expansions. In *Advances in Neural Information Processing Systems* 27, 2014.
- [8] R. Amadeo. Adware vendors buy chrome extensions to send ad- and malware-filled updates. *Ars Technica*. <http://arstechnica.com/security/2014/01/malware-vendors-buy-chrome-extensions-to-send-adware-filled-updates/>.
- [9] A. Anandkumar, K. Chaudhuri, D. Hsu, S. M. K. akade, L. Song, and T. Zhang. Spectral methods for learning multivariate latent tree structure. In *Advances in Neural Information Processing Systems* 24, 2011.
- [10] A. Anandkumar, D. P. Foster, D. Hsu, S. M. Kakade, and Y.-K. Liu. A spectral algorithm for latent Dirichlet allocation. *Algorithmica*, 72(1):193–214, 2015.
- [11] A. Anandkumar, R. Ge, D. Hsu, and S. M. Kakade. A tensor approach to learning mixed membership community models. *Journal of Machine Learning Research*, 15(Jun):2239–2312, 2014.
- [12] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15(Aug):2773–2831, 2014.
- [13] P. Barford, I. Canadi, D. Krushevskaja, Q. Ma, and S. Muthukrishnan. Adscape: Harvesting and Analyzing Online Display Ads. *WWW '14: Proceedings of the 23rd international conference on World Wide Web*, Apr. 2014.
- [14] A. Beygelzimer, D. Hsu, J. Langford, and T. Zhang. Agnostic active learning without constraints. In *Advances in Neural Information Processing Systems* 23, 2010.
- [15] P. J. Bickel, Y. Ritov, and A. B. Tsybakov. Simultaneous analysis of lasso and dantzig selector. *Ann. Statist.*, 37(4):1705–1732, 08 2009.
- [16] E. J. Candès, J. K. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- [17] A. Chaintreau. *Processes of interaction in data networks*. PhD thesis, Ecole Normale Supérieure, 2006.

- [18] A. Chaintreau, P. Fraigniaud, and E. Lebhar. Networks Become Navigable as Nodes Move and Forget. In *ICALP '08: Proceedings of the 35th international colloquium on Automata, Languages and Programming, Part I*. Springer-Verlag, July 2008.
- [19] A. Chaintreau, P. Fraigniaud, and E. Lebhar. Opportunistic spatial gossip over mobile social networks. In *WOSN '08: Proceedings of the first workshop on Online social networks*. ACM, Aug. 2008.
- [20] A. Chaintreau, P. Hui, J. Crowcroft, C. Diot, R. Gass, and J. Scott. Impact of Human Mobility on Opportunistic Forwarding Algorithms. *IEEE Transactions on Mobile Computing*, 6(6):606–620, 2007.
- [21] A. Chaintreau, J.-Y. Le Boudec, and N. Ristanovic. The age of gossip: spatial mean field regime. In *SIGMETRICS '09: Proceedings of the eleventh international joint conference on Measurement and modeling of computer systems*. ACM Request Permissions, June 2009.
- [22] A. Chaintreau, A. Mtibaa, L. Massoulié, and C. Diot. The diameter of opportunistic mobile networks. In *CoNEXT '07: Proceedings of the 2007 ACM CoNEXT conference*. ACM, Dec. 2007.
- [23] K. Chaudhuri and D. Hsu. Sample complexity bounds for differentially private learning. In *Twenty-Fourth Annual Conference on Learning Theory*, 2011.
- [24] K. Chaudhuri and D. Hsu. Convergence rates for differentially private statistical estimation. In *Twenty-Ninth International Conference on Machine Learning*, 2012.
- [25] K. Chaudhuri, D. Hsu, and S. Song. The large margin mechanism for differentially private maximization. In *Advances in Neural Information Processing Systems 27*, 2014.
- [26] W. Cheng, Q. Zhao, B. Yu, and S. Hiroshige. Tainttrace: Efficient flow tracing with dynamic binary rewriting. In *Proceedings of the 11th IEEE Symposium on Computers and Communications*, 2006.
- [27] S. Dasgupta and D. Hsu. Hierarchical sampling for active learning. In *Twenty-Fifth International Conference on Machine Learning*, 2008.
- [28] S. Dasgupta, D. Hsu, and C. Monteleoni. A general agnostic active learning algorithm. In *Advances in Neural Information Processing Systems 20*, 2007.
- [29] A. Datta, M. C. Tschantz, and A. Datta. Automated Experiments on Ad Privacy Settings: A Tale of Opacity, Choice, and Discrimination. *arXiv.org*, Aug. 2014.
- [30] A. Datta, M. C. Tschantz, and A. Datta. Automated Experiments on Ad Privacy Settings: A Tale of Opacity, Choice, and Discrimination. In *Proceedings of the 15th Privacy Enhancing Technologies Symposium (PETS)*, 2015.
- [31] Y. A. de Montjoye, E. Shmueli, S. S. Wang, and A. S. Pentland. openPDS: Protecting the Privacy of Metadata through SafeAnswers. *PLoS One*, 2014.
- [32] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.
- [33] Y. Duan and J. Canny. Practical private computation and zero-knowledge tools for privacy-preserving distributed data mining. *Proceedings of the SIAM International Conference on Data Mining (SDM 2008), Atlanta, Georgia, USA*, pages 265–276, 2008.
- [34] E. Dwoskin. Why You Can't Trust You're Getting the Best Deal Online. *online.wsj.com*, Oct. 2014.
- [35] P. Eckersley. How unique is your web browser? In *Proceedings of the International Conference on Privacy Enhancing Technologies*, 2010.
- [36] Electronic Frontier Foundation. Privacy Badger, July 2014. <https://www EFF.org/privacybadger>.
- [37] W. Enck, P. Gilbert, B. gon Chun, L. P. Cox, J. Jung, P. McDaniel, and A. N. Sheth. TaintDroid: An information-flow tracking system for realtime privacy monitoring on smartphones. In *Proc. of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2010.
- [38] S. Englehardt, C. Eubank, P. Zimmerman, D. Reisman, and A. Narayanan. Web Privacy Measurement: Scientific principles, engineering platform, and new results. *Princeton University*, June 2014.
- [39] Ú. Erlingsson, V. Pihur, and A. Korolova. RAPPOR: Randomized Aggregatable Privacy-Preserving Ordinal Response. *arXiv.org*, July 2014.
- [40] V. Erramilli, A. Chaintreau, M. Crovella, and C. Diot. Diversity of forwarding paths in pocket switched networks. In *IMC '07: Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*. ACM Request Permissions, Oct. 2007.
- [41] V. Erramilli, M. Crovella, A. Chaintreau, and C. Diot. Delegation forwarding. In *MobiHoc '08: Proceedings of the 9th ACM international symposium on Mobile ad hoc networking and computing*. ACM Request Permissions, May 2008.

- [42] J. P. Finis, M. Raiber, N. Augsten, R. Brunel, A. Kemper, and F. Färber. Rws-diff: Flexible and efficient change detection in hierarchical data. In *Proceedings of the 22Nd ACM International Conference on Conference on Information & Knowledge Management, CIKM '13*, pages 339–348, New York, NY, USA, 2013. ACM.
- [43] R. Geambasu, T. Kohno, A. Krishnamurthy, A. Levy, H. M. Levy, P. Gardner, and V. Moscaritolo. New directions for self-destructing data systems. Technical Report UW-CSE-11-08-01, University of Washington, 2010.
- [44] R. Geambasu, T. Kohno, A. Levy, and H. M. Levy. Vanish: Increasing data privacy with self-destructing data. In *Proc. of USENIX Security*, 2009.
- [45] R. Geambasu, A. Levy, T. Kohno, A. Krishnamurthy, and H. M. Levy. Comet: An active distributed key/value store. In *Proc. of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2010.
- [46] Ghostery Enterprise. Ghostery. <https://www.ghostery.com/>.
- [47] D. B. Giffin, A. Levy, D. Stefan, D. Terei, D. Mazières, J. C. Mitchell, and A. Russo. Hails: protecting data privacy in untrusted web applications. In *OSDI'12: Proceedings of the 10th USENIX conference on Operating Systems Design and Implementation*. USENIX Association, Oct. 2012.
- [48] P. Gill, V. Erramilli, A. Chaintreau, B. Krishnamurthy, K. Papagiannaki, and P. Rodriguez. Follow the money: understanding economics of online aggregation and advertising. *IMC '13: Proceedings of the 2013 conference on Internet measurement conference*, Oct. 2013.
- [49] A. Hannak, P. Sapiezynski, A. M. Kakhki, B. Krishnamurthy, D. Lazer, A. Mislove, and C. Wilson. Measuring personalization of web search. In *WWW '13: Proceedings of the 22nd international conference on World Wide Web*. International World Wide Web Conferences Steering Committee, May 2013.
- [50] A. Hannak, G. Soeller, D. Lazer, A. Mislove, and C. Wilson. Measuring Price Discrimination and Steering on E-commerce Web Sites. *IMC '14: Proceedings of the 14th ACM SIGCOMM conference on Internet measurement*, 2014.
- [51] D. Hsu and S. M. Kakade. Learning mixtures of spherical Gaussians: moment methods and spectral decompositions. In *Fourth Innovations in Theoretical Computer Science*, 2013.
- [52] D. Hsu, S. M. Kakade, J. Langford, and T. Zhang. Multi-label prediction via compressed sensing. In *Advances in Neural Information Processing Systems* 22, 2009.
- [53] D. Hsu, S. M. Kakade, and P. Liang. Identifiability and unmixing of latent parse trees. In *Advances in Neural Information Processing Systems* 25, 2012.
- [54] D. Hsu, S. M. Kakade, and T. Zhang. A spectral algorithm for learning hidden Markov models. *Journal of Computer and System Sciences*, 78(5):1460–1480, 2012.
- [55] J. Huang and T. Zhang. The benefit of group sparsity. *Annals of Statistics*, 38:1978–2004, 2010.
- [56] S. Ioannidis, A. Chaintreau, and L. Massoulié. Optimal and Scalable Distribution of Content Updates over a Mobile Social Network. *INFOCOM 2009, IEEE*, pages 1422–1430, 2009.
- [57] S. Ioannidis, L. Massoulié, and A. Chaintreau. Distributed caching over heterogeneous mobile networks. In *SIGMETRICS '10: Proceedings of the ACM SIGMETRICS international conference on Measurement and modeling of computer systems*. ACM Request Permissions, June 2010.
- [58] S. Isaacman, S. Ioannidis, A. Chaintreau, and M. Martonosi. Distributed rating prediction in user generated content streams. In *RecSys '11: Proceedings of the fifth ACM conference on Recommender systems*. ACM Request Permissions, Oct. 2011.
- [59] S. Kamkar. Evercookie—virtually irrevocable persistent cookies. <http://samy.pl/evercookie/>.
- [60] M. Lecuyer, G. Ducoffe, F. Lan, A. Papancea, T. Petsios, R. Spahn, A. Chaintreau, and R. Geambasu. XRay: Enhancing the Web's Transparency with Differential Correlation. In *23rd USENIX Security Symposium (USENIX Security 14)*, San Diego, CA, 2014. USENIX Association.
- [61] P. G. Leon, B. Ur, Y. Wang, M. Sleeper, R. Balebako, R. Shay, L. Bauer, M. Christodorescu, and L. F. Cranor. What Matters to Users? Factors that Affect Users' Willingness to Share Information with Online Advertisers. In *Symposium on Usable Privacy and Security*, 2013.
- [62] B. Liu, A. Sheth, U. Weinsberg, J. Chandrashekar, and R. Govindan. Adreveal: Improving transparency into online targeted advertising. In *Proceedings of the Twelfth ACM Workshop on Hot Topics in Networks (HotNets)*, 2013.
- [63] B. Liu, A. Sheth, U. Weinsberg, J. Chandrashekar, and R. Govindan. AdReveal: improving transparency into online targeted advertising. In *HotNets-XII: Proceedings of the Twelfth ACM Workshop on Hot Topics in Networks*. ACM Request Permissions, Nov. 2013.
- [64] A. M. McDonald and L. F. Cranor. Americans' Attitudes about Internet Behavioral Advertising Practices. In *Proceedings of the Workshop on Privacy in the Electronic Society*, 2010.

- [65] J. Mikians, L. Gyarmati, V. Erramilli, and N. Laoutaris. Detecting price and search discrimination on the internet. In *HotNets-XI: Proceedings of the 11th ACM Workshop on Hot Topics in Networks*. ACM Request Permissions, Oct. 2012.
- [66] V. Misra, S. Ioannidis, A. Chaintreau, and L. Massoulié. Incentivizing peer-assisted services: a fluid shapley value approach. In *SIGMETRICS '10: Proceedings of the ACM SIGMETRICS international conference on Measurement and modeling of computer systems*. ACM Request Permissions, June 2010.
- [67] Mozilla. Lightbeam. <https://www.mozilla.org/en-US/lightbeam/about/>.
- [68] M. Mun, S. Hao, N. Mishra, K. Shilton, J. Burke, D. Estrin, M. Hansen, and R. Govindan. Personal data vaults: a locus of control for personal data streams. *Proceedings of the 6th International Conference*, page 17, 2010.
- [69] N. Nikiforakis, A. Kapravelos, W. Joosen, C. Kruegel, F. Piessens, and G. Vigna. Cookieless Monster: Exploring the Ecosystem of Web-based Device Fingerprinting. In *Proceedings of the IEEE Symposium on Security and Privacy*, 2013.
- [70] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2009.
- [71] E. Rader. Awareness of behavioral tracking and information privacy concern in facebook and google. In *Symposium on Usable Privacy and Security*, 2014.
- [72] J. Reich and A. Chaintreau. The age of impatience: optimal replication schemes for opportunistic networks. In *CoNEXT '09: Proceedings of the 5th international conference on Emerging networking experiments and technologies*. ACM Request Permissions, Dec. 2009.
- [73] C. Riederer, V. Erramilli, A. Chaintreau, B. Krishnamurthy, and P. Rodriguez. For sale : your data: by : you. In *HotNets-X: Proceedings of the 10th ACM Workshop on Hot Topics in Networks*. ACM Request Permissions, Nov. 2011.
- [74] F. Roesner, J. Fogarty, and T. Kohno. User Interface Toolkit Mechanisms for Securing Interface Elements. In *Proceedings of the ACM Symposium on User Interface Software and Technology*, 2012.
- [75] F. Roesner and T. Kohno. Securing Embedded User Interfaces: Android and Beyond. In *Proceedings of the USENIX Security Symposium*, 2013.
- [76] F. Roesner, T. Kohno, and D. Molnar. Security and Privacy for Augmented Reality Systems. *Communications of the ACM*, 57:88–96, 2014.
- [77] F. Roesner, T. Kohno, A. Moshchuk, B. Parno, H. J. Wang, and C. Cowan. User-Driven Access Control: Re-thinking Permission Granting in Modern Operating Systems. In *Proceedings of the IEEE Symposium on Security and Privacy*, 2012.
- [78] F. Roesner, T. Kohno, and D. Wetherall. Detecting and defending against third-party tracking on the web. In *Proceedings of the 9th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2012.
- [79] F. Roesner, D. Molnar, A. Moshchuk, T. Kohno, and H. J. Wang. World-Driven Access Control for Continuous Sensing Applications. In *ACM Conference on Computer and Communications Security*, 2014.
- [80] F. Roesner, C. Rovillos, T. Kohno, and D. Wetherall. ShareMeNot: Balancing Privacy and Functionality of Third-Party Social Widgets. *USENIX ;login.*, 37, 2012. <https://sharemenot.cs.washington.edu/>.
- [81] F. Roesner, C. Rovillos, A. Saxena, and T. Kohno. Trackingobserver: A browser-based web tracking detection platform, 2013. <https://trackingobserver.cs.washington.edu/>.
- [82] R. Spahn, J. Bell, R. Geambasu, and G. Kaiser. Pebbles: Fine-grained data management abstractions for modern operating systems. In *Proc. of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2014.
- [83] L. Sweeney. Discrimination in online ad delivery. *Communications of the ACM*, 56(5), May 2013.
- [84] Y. Tang, P. Ames, S. Bhamidipati, A. Bijlani, R. Geambasu, and N. Sarda. CleanOS: Mobile OS abstractions for managing sensitive data. In *Proc. of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2012.
- [85] The Wall Street Journal. What they know, 2010–2012. <http://www.wsj.com/public/page/what-they-know-digital-privacy.html>.
- [86] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1994.
- [87] B. Ur, P. G. Leon, L. F. Cranor, R. Shay, and Y. Wang. Smart, useful, scary, creepy: perceptions of online behavioral advertising. In *8th Symposium on Usable Privacy and Security*, 2012.
- [88] N. Viennot, M. Lecuyer, J. Bell, R. Geambasu, and J. Nieh. Synapse: New data integration abstractions for agile web application development. In *Proc. of the ACM European Conference on Computer Systems (EuroSys)*, 2015.

- [89] X. Xing, W. Meng, D. Doozan, N. Feamster, and W. Lee. Exposing Inconsistent Web Search Results with Bobble. *cseweb.ucsd.edu*.
- [90] X. Xing, W. Meng, D. Doozan, N. Feamster, W. Lee, and A. C. Snoeren. Exposing Inconsistent Web Search Results with Bobble. *Passive and Active Measurements Conference*, 2014.
- [91] H. Xu, Q. Wu, H. Wang, G. Yang, and Y. Jia. Kf-diff+: Highly efficient change detection algorithm for xml documents. In R. Meersman and Z. Tari, editors, *On the Move to Meaningful Internet Systems 2002: CoopIS, DOA, and ODBASE*, volume 2519 of *Lecture Notes in Computer Science*, pages 1273–1286. Springer Berlin Heidelberg, 2002.
- [92] T.-F. Yen, Y. Xie, F. Yu, R. P. Yu, and M. Abadi. Host Fingerprinting and Tracking on the Web: Privacy and Security Implications. In *Proceedings of the Network and Distributed System Security Symposium*, 2012.
- [93] Y. Zhu, J. Jung, D. Song, T. Kohno, and D. Wetherall. Privacy scope: A precise information flow tracking system for finding application leaks. Technical Report UCB/EECS-2009-145, EECS Department, University of California, Berkeley, Oct 2009.

10 Appendix A

10.1 Team Member Identification

Individual Name	Role (Prime, Subcontractor or Consultant)	Organization	Non-US?		FFRDC or Govt?
			Org.	Ind.	
Geambasu	Prime	Columbia University	N/A	N/A	N/A
Chaintreau	Prime	Columbia University	N/A	N/A	N/A
Hsu	Prime	Columbia University	N/A	N/A	N/A
Roesner	Subcontractor	University of Washington	N/A	N/A	N/A

10.2 Government or FFRDC Team Member Proof of Eligibility to Propose

NONE

10.3 Government or FFRDC Team Member Statement of Unique Capability

NONE

10.4 Organizational Conflict of Interest Affirmations and Disclosure

NONE

10.5 Intellectual Property (IP)

The Offeror and subcontractors reserve the right to independently or jointly seek intellectual protection for the results of the work under this program. These rights will not compromise the values of the proposed work to the Government because it will have access to and use of the research and results of this work.

10.6 Human Subjects Research (HSR)

The proposed work includes user studies that will involve human subject research. The proposed studies will be designed and conducted according to procedures approved by the organizations' Institutional Review Boards (IRBs). Ample time will be allotted to complete the approval process for each study.

10.7 Animal Use

NONE

10.8 Representations Regarding Unpaid Delinquent Tax Liability or a Felony Conviction under Any Federal Law

(a) The proposer represents that it is [] is not [**X**] a corporation that has any unpaid Federal tax liability that has been assessed, for which all judicial and administrative remedies have been exhausted or have lapsed, and that is not being paid in a timely manner pursuant to an agreement with the authority responsible for collecting the tax liability.

(b) The proposer represents that it is [] is not [**X**] a corporation that was convicted of a felony criminal violation under a Federal law within the preceding 24 months.

10.9 Cost Accounting Standards (CAS) Notices and Certification

NONE