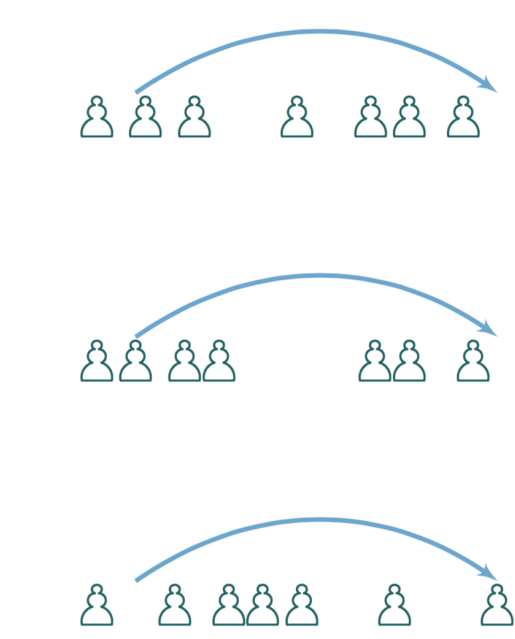


PSI (Ψ): a Private data Sharing Interface

Marco Gaboardi, James Honaker, Gary King, Kobbi Nissim, Jonathan Ullman, Salil Vadhan

<http://privacytools.seas.harvard.edu/psi>



Objectives

Differential privacy offers an attractive approach to enabling data sharing among social science researchers. Our desiderata are:

- **Accessibility by non-experts:** researchers in the social sciences should be able to use the system to share and explore data with no involvement from experts in data privacy, computer science, or statistics.
- **Generality:** the system should be applicable and effective on a wide variety of heterogeneous datasets hosted in a repository such as the Harvard Dataverse.
- **Workflow-compatibility:** the system should fit naturally in the workflow of its users (e.g. researchers in the social sciences), and be positioned to offer clear benefits (e.g. more access to sensitive data or less risk of an embarrassing privacy violation) rather than being an impediment.

Motivation

Researchers in all experimental and empirical fields are increasingly expected to widely share the data behind their published research, to enable other researchers to verify, replicate, and extend their work. Indeed, data-sharing is now often mandated by funding agencies.

However, many of the datasets in the social and health sciences contain sensitive personal information about human subjects.

- 1 Numerous data sets, such as surveys, that have been "deidentified" via traditional means are increasingly being deposited in publicly accessible data repositories at risk of reidentification.
- 2 Numerous other data sets are not made available at all, or only with highly restrictive and time-consuming provisions.

An important problem is to develop methods that can be used to offer greater privacy protections for datasets of the first type, and enable the safe sharing of datasets of the second type.

Actors and Workflow

We have three different kinds of actors:

- 1 **Data depositors:** Users that come to deposit their privacy-sensitive dataset in a data repository, and wish to make differentially private access to their dataset available.
- 2 **The data curators:** Data-repository managers that maintain the hardware and software on which PSI runs and the data repository infrastructure. They are trusted, and may also have legal obligations to maintain privacy.
- 3 **Data analysts:** Researchers that come to access sensitive datasets in the repository, often with the goal of data exploration. Have access to all of the differentially private statistics released by the data depositor, as well as ability to make their own differentially private queries (subject to the overall privacy budget).

Budgeting Tool

Variable	Type	Statistic	Upper Bound	Lower Bound	Granularity	Number of Rows	Epsilon	Accuracy	Hold
X	age	Numerical	Mean	100	0	na	na	0.0316	0.0474
X	educ	Numerical	Histogram	na	na	na	20	0.0316	0.01
X	sex	Categorical	Histogram	na	na	na	2	0.0316	0.0948
X	income	Numerical	Quantile	1000000	0	1000	na	0.021	0.0331
X	income	Numerical	Mean	1000000	0	na	na	0.0509	0.0200
X	white	Boolean	Histogram	na	na	na	2	0.0316	0.0948
X	black	Boolean	Histogram	na	na	na	2	0.0316	0.0948

Figure 1: Interactive Interface for Budgeting DP Releases

We have developed a privacy budgeting tool that exposes the privacy-accuracy tradeoff to the user.

To ensure that we get the most utility out of the global privacy budget, we use the "approximate optimal composition theorem" which in fact was developed for the purpose of our privacy budget tool.

Architecture

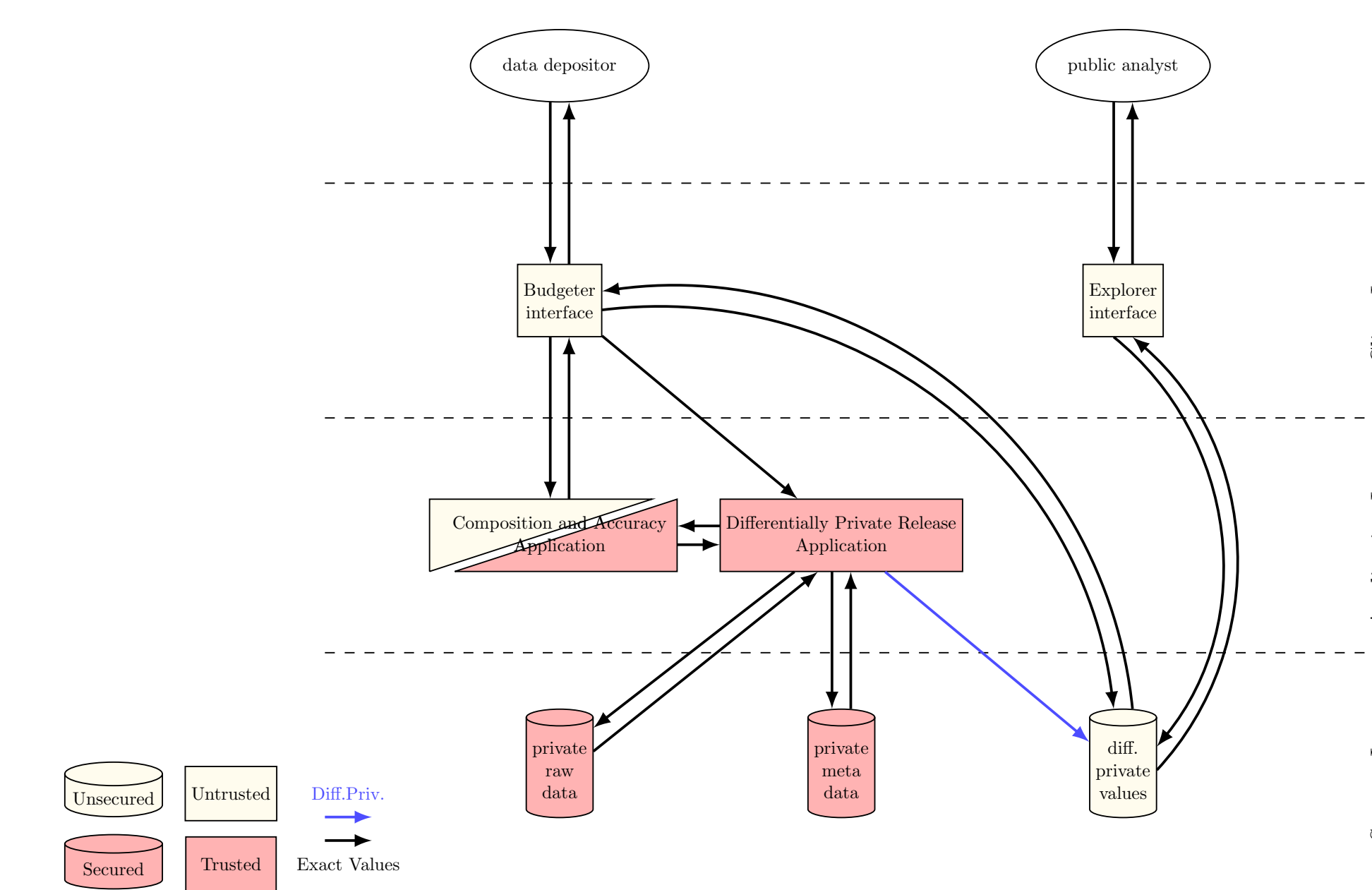


Figure 2: Architecture Diagram for Components of System

Conclusion

The goal of PSI is to provide increased data sharing and exploration of private data in repositories. We continue to add new algorithms to further aid researchers in exploratory analysis, while judging the utility of the system by experimental studies, replication analyses of published work using our tools, and usability experiments, to bring differential privacy as a practical tool to social science research.

Acknowledgements

This work is part of the "Privacy Tools for Sharing Research Data" project at Harvard, supported by NSF grant CNS-1237235 and a grant from the Sloan Foundation. The accompanying working paper describes a vision for work that is still in progress, and is therefore authored by the leadership of the efforts. The software project includes contributions from Andreea Antueta, Connor Bain, Victor Balcer, Jessica Bu, Mark Bun, Stephen Chong, Vito D'Orazio, Anna Gavrilman, Caper Gooden, Paul Handorff, Raquel Hill, Allyson Kaminsky, Murat Kuntarcioglu, Vishesh Karwa, George Kellaris, Hyun Woo Lim, Nathan Manohar, Dan Muise, Jack Murtagh, Sofya Raskhodnikova, Ryan Rogers, Or Sheffet, Adam D. Smith, Thomas Steinke, David Xiao, Haoqing Wang, and Joy Zheng

Workflow



The System

Unique features of PSI include:

- None of its users, are expected to have expertise in privacy, computer science, or statistics.
- It is designed to be integrated with existing and widely used data repository infrastructures.
- Its initial set of differentially private algorithms have wide use in the social sciences, and are integrated with existing statistical software.
- We have included pedagogical materials explaining differential privacy in an intuitive but accurate manner, with a minimum of technical terminology and notation.

Statistics

The initial set of dp-algorithms were chosen to give immediate utility for social science research:

- Univariate descriptive statistics, such as means, quantiles, histograms, and approximate cdfs.
- Basic statistical estimators, such as difference-of-means testing for causal inference, hypothesis tests for the independence of categorical variables, and low-dimensional covariance matrices for least-squares regressors, and principal components.
- Transformations for creating new features (variables) out of combinations of already existing ones.

Contact Information

- Web: <http://privacytools.seas.harvard.edu>
- Email: privacytools-info@seas.harvard.edu

