# OpenDP

# An open-source suite of tools for deploying differential privacy

**Machine Learning in Science & Engineering 2020**
December 14-15, Columbia University Data Science Institute

**Mercè Crosas, Ph.D., Harvard University**
scholar.harvard.edu/mercecrosas  @mercecrosas

IQSS
The Institute for Quantitative Social Science

HARVARD
UNIVERSITY

# What is Differential Privacy?

- A **differentially private (DP)** algorithm introduces a minimum amount of noise to released statistics to mathematically guarantee the privacy of any individual in a dataset

- Aims to [Dwork, McSherry, Nissim, Smith, '06]:

    - enable statistical analysis of datasets -  **utility**

    - while protecting individual level data – **privacy**

- In the last years, DP has moved from theory to practice and starts to be deployed in products

# What is OpenDP?

- OpenDP is a **community effort** to build a **trustworthy** and **open-source** suite of **differential privacy tools** to explore and analyze sensitive data

- Aims to:

  - Channel **collective advances** on the science and practice of DP

  - Enable wider **adoption** of DP

  - Address high-demand, compelling **use cases**

  - Provide the building blocks for deploying **custom DP solutions**

  - Identify important **research directions** for the field

# The OpenDP Team

OpenDP

## Executive Committee

**Salil Vadhan**
Faculty Co-Director

**Gary King**
Faculty Co-Director

**Annie Wu**
Program Director

**James Honaker**
Chief Privacy Engineer

**Mercè Crosas**
Chief Data Science & Technology Officer, IQSS

## Development & Staff

Ellen Kraffmiller
Technical Lead

Ethan Cowan
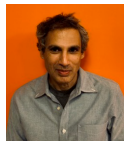Software Developer

Christian Covington
Research Fellow

Ira Globus-Harris
Software Engineer

Jayshree Sarathy
Graduate Student

Lindsay Froess
Project Coordinator

Raman Prasad
Technical Lead for Research Software

Audrey Haque
UX Designer & Researcher

Michael Shoemate
Senior Software Developer

Robert Treacy
Senior Application Architect

Danny Brooke
Program Manager, Product Development, IQSS

Gustavo Durand
Dataverse Technical Lead and Architect

Tania Schlatter
UX & UI Lead, Product Development

Michael Phelan
Senior Application Architect

Andy Vyrros
Senior Library Architect

# Advisory Board

OpenDP

**Gerome Miklau (co-chair)**
UMass, Amherst

**Adam Smith (co-chair)**
Boston University

**John Abowd**
US Census Bureau, Cornell University

**John Friedman**
Brown University

**Margaret Levenstein**
University of Michigan

**Jules Polonetsky**
Future of Privacy Forum

**Gilles Barthe**
Max Planck Institute, IMDEA Software Institute

**Jeff Gill**
American University

**Carlos Maltzahn**
UC Santa Cruz CROSS

**Aaron Roth**
University of Pennsylvania

**Barbara Bierer**
Brigham & Women's Hospital Harvard Medical School

**Daniel Goroff**
Alfred P. Sloan Foundation (ex officio)

**Kenneth Mandl**
Harvard Medical School Boston Children's Hospital

**Aleksandra Slavkovic**
Pennsylvania State University

**Sarah Bird**
Microsoft

**Frauke Kreuter**
University of Mannheim Institute for Employment Research, Germany

**Ilya Mironov**
Facebook

**Dawn Song**
University of California, Berkeley

**danah boyd**
Microsoft Research

**Orran Krieger**
PI Mass Open Cloud Boston University

**Helen Nissenbaum**
Cornell University

**Latanya Sweeney**
Harvard University

**Cynthia Dwork**
Harvard University Radcliffe Institute Microsoft Research

**David Lazer**
Northeastern University

**Kobbi Nissim**
Georgetown University

**Omer Tene**
International Association of Privacy Professionals

**Úlfar Erlingsson**
Apple

**Katrina Ligett**
Hebrew University

**Dina N. Paltoo**
National Heart, Lung and Blood Institute

**Stefaan G Verhulst**
New York University

# The OpenDP initiative launched in 2020

## OpenDP Community Meeting 2020

### Plenary Session - Wednesday, May 13th

▶ 11:00 - 11:30: Overview of OpenDP and Goals for Meeting

▶ 11:35 - 11:50: Use Cases

▶ 11:55 - 12:10: Differential Privacy for COVID-19 Data

▶ 12:15 - 12:45: A Programming Framework for the DP Library

▶ 12:50 - 1:20: Break

▶ 1:20 - 1:40: Statistical Functionality

▶ 1:45 - 2:05: System Integrations

▶ 2:10 - 2:30: Governance and Licensing

▶ 2:35 - 2:55: Collaborations

- **OpenDP Community Meeting in May 2020**

- **Key elements in whitepaper:** use cases, governance, programming framework, statistical functionality, system integrations, collaboration & community

- Previous work with the **Harvard Privacy Tools Project**, funded by NSF, US Census Bureau, the Sloan Foundation, and Google

- Current grants from the **Sloan Foundation**

- Current collaboration with **Microsoft**

Alfred P. Sloan FOUNDATION

Microsoft

NSF

Google

# Open-source Governance and Components

# Open-source Governance and Components

**OpenDP Commons**: DP tools and packages built and used by the community

**Vetted Core**

OpenDP Library

Tool 1

Tool 2

Package 1

**Product Solutions powered by OpenDP**

**SmartNoise 0.x**

Product suite for deploying OpenDP, developed in partnership with Microsoft

System 2

Library

Tool 3

System 3

Tool 1

Package 2

Tool 4

Library

Contribute new components to OpenDP Commons

# SmartNoise 0.x: First end-to-end OpenDP system



**Use Data.
Preserve Privacy.**

A differential privacy toolkit for analytics
and machine learning

- **Released in June 2020**

- A collaboration with Microsoft: **Sarah Bird, John Kahan, Joshua Allen, Eddie de Leon, Kevin White**

- An OpenDP end-to-end system **proof-of-concept**

- **SmartNoise Core v0.2.0:**
  - Rust Library and Python bindings
  - Statistics: mean, variance, count, histogram, general quantiles, sum, covariance, simple linear regression
  - Mechanisms: Laplace, Gaussian, snapping, geometric, exponential

- **SmartNoise Tools/SDK v0.1.0**



Microsoft

IQSS
The Institute for Quantitative Social Science

Harvard John A. Paulson
**School of Engineering**
and Applied Sciences

# SmartNoise Core 0.2.0

## Basic PUMS Analysis with SmartNoise

```
In [8]:   # attempt 4 - succeeds!
          with sn.Analysis() as analysis:
              # load data
              data = sn.Dataset(path = data_path, column_names = var_names)

              ''' get mean age '''
              # establish data
              age_dt = sn.to_float(data['age'])

              # clamp data to range and impute missing values
              age_dt = sn.clamp(data = age_dt, lower = age_range[0], upper = age_range[1])
              age_dt = sn.impute(data = age_dt, distribution = 'Gaussian',
                                              lower = age_range[0], upper = age_range[1],
                                              shift = 45., scale = 10.)
```

Run a DP statistic on the 'PUMS_california_demographics_1000' dataset

```
Post-Release

DP mean of age: 44.70965828557311
Privacy usage: approximate {
  epsilon: 0.65
}
```

Release DP mean of a variable, with privacy loss information

# Open-source Governance and Components

**OpenDP Commons:** DP tools and packages built and used by the community

OpenDP Library

**Product Solutions powered by OpenDP**

**SmartNoise 0.x**

Product suite for deploying OpenDP, developed in partnership with Microsoft

System 2

Library

Tool 3

System 3

Tool 1

Package 2

Tool 4

Library

Contribute new components to OpenDP Commons

# OpenDP Library

- A **novel programming framework** [Marco Gaboardi, Michael Hay, Salil Vadhan]
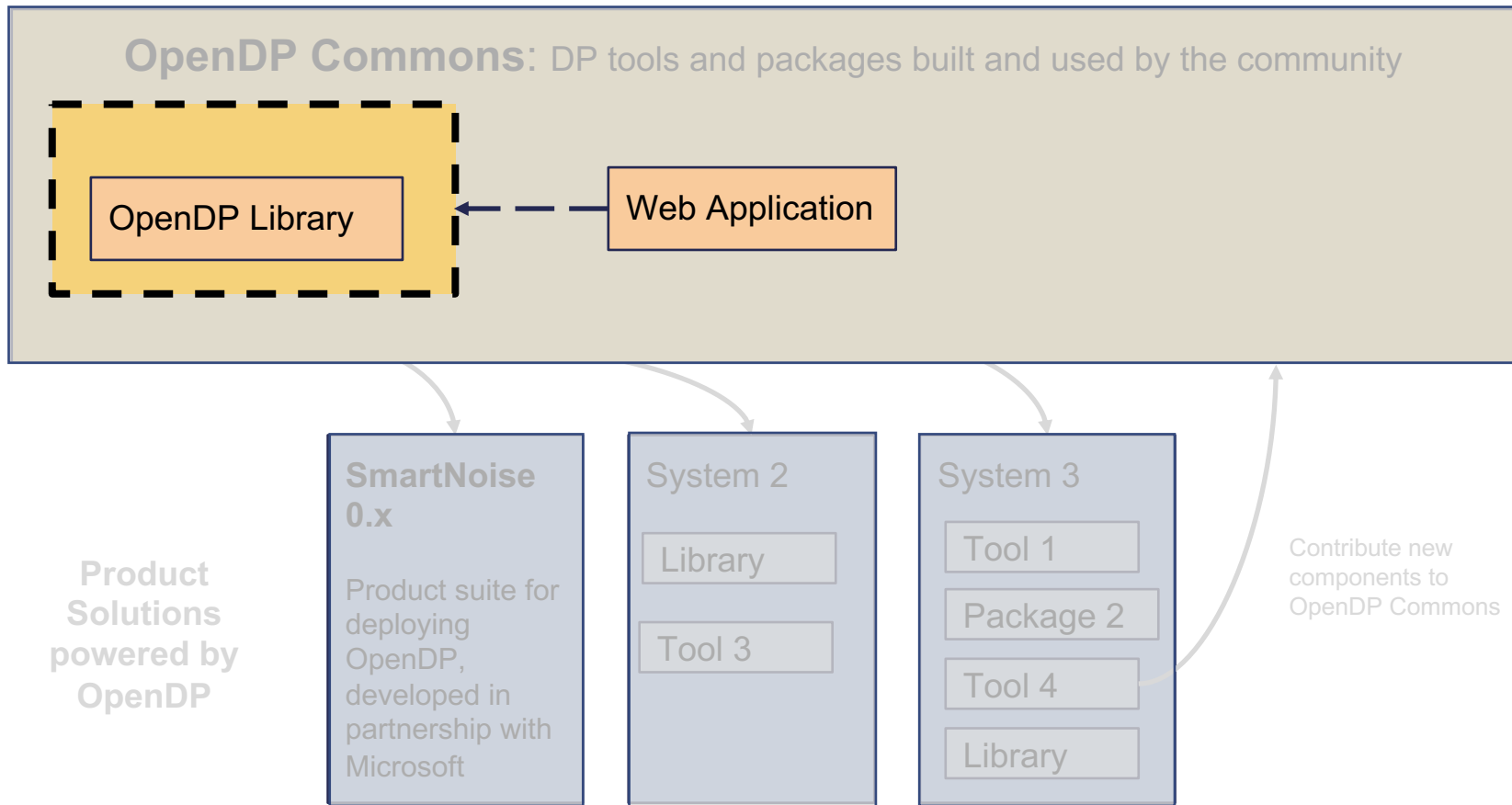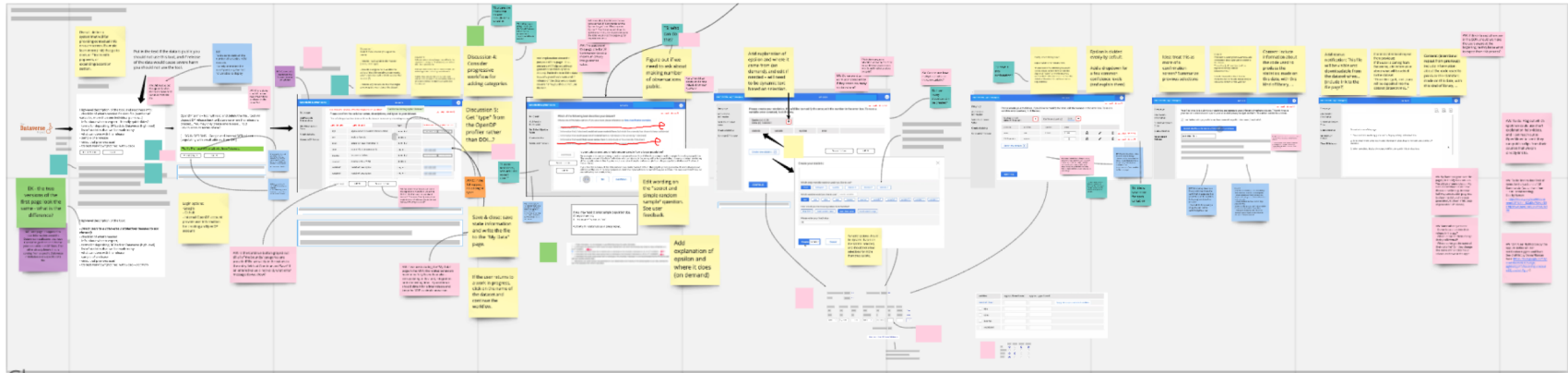
- Goes beyond simply solving a fixed set of DP models

- **Abstracts** the treatment of privacy mechanisms

- Allows building of **new, complex privacy** systems

- Facilitates and verifies **contributions**

- **First release to the community: Early 2021**

# Open-source Governance and Components

**OpenDP Commons:** DP tools and packages built and used by the community

OpenDP Library

Web Application

**Product Solutions powered by OpenDP**

**SmartNoise 0.x**

Product suite for deploying OpenDP, developed in partnership with Microsoft

System 2

Library

Tool 3

System 3

Tool 1

Package 2

Tool 4

Library

Contribute new components to OpenDP Commons

# OpenDP Web Application



- Uses the OpenDP library and focuses on UI/UX for **user-friendly, wide adoption**
- First **use case** aims to support archival **data repositories**:
  - Enables secondary reuse of data and reproducibility of published results
  - Integrates with > 60 **Dataverse repositories** world-wide
- DP released statistics must expose measures of **utility and uncertainty**
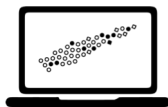- **First release to the community: Early 2021**

# OpenDP + Dataverse

## Public Data Repository
Sensitive datasets discoverable via repository (only metadata is open)

### Data User
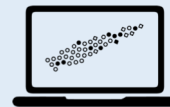**DP statistics release + Privacy loss + Measure of utility or uncertainty**

Launch OpenDP for a dataset found in Dataverse

Deposit DP release back to Dataverse

## Data owner or Data analyst
**Select DP Statistics and privacy loss**

Trusted Remote Storage Agents (TRSA) or data enclaves

Impact: Infrastructure for Privacy-Assured Computation
https://cyberimpact.us/

# Open-source Governance and Components

# Spring 2021

- **Feedback** on OpenDP Library
- Review and vetting process for **new contributions**
- OpenDP Community Meeting
- New OpenDP fellows program
- Launch working groups
- Build the community to help deploy **new systems**



OpenDP

# Get Involved

https://opendp.org

Mailing list: https://opendp.org/join

https://github.com/opendifferentialprivacy/

@opendp_io