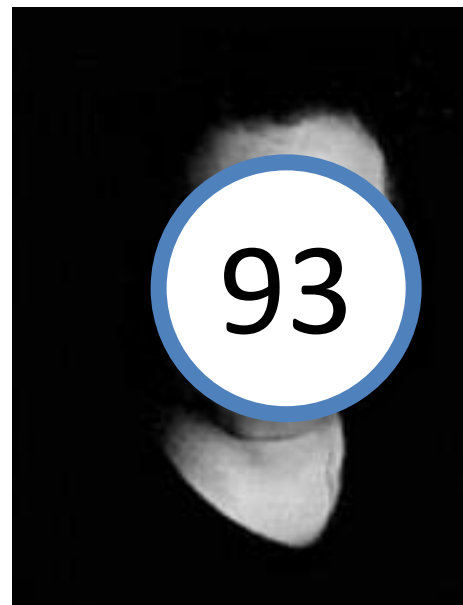# De-Identification and the Health Insurance Portability and Accountability Act (HIPAA)

## Overview and framing of current issues

93

**Simson L. Garfinkel, Ph.D.**

**Information Technology Laboratory**

**National Institute of Standards and Technology**

Subcommittee on Privacy, Confidentiality & Security
National Committee on Vital and Health Statistics
May 24, 2016

NIST
National Institute of
Standards and Technology
U.S. Department of Commerce

# National Institute of Standards and Technology



*Founded in 1901*

Non-regulatory federal laboratory.

Mission:

"To promote US innovation and industrial competitiveness by advancing measurement science, standards, and technology in ways that enhance economic security and improve our quality of life."

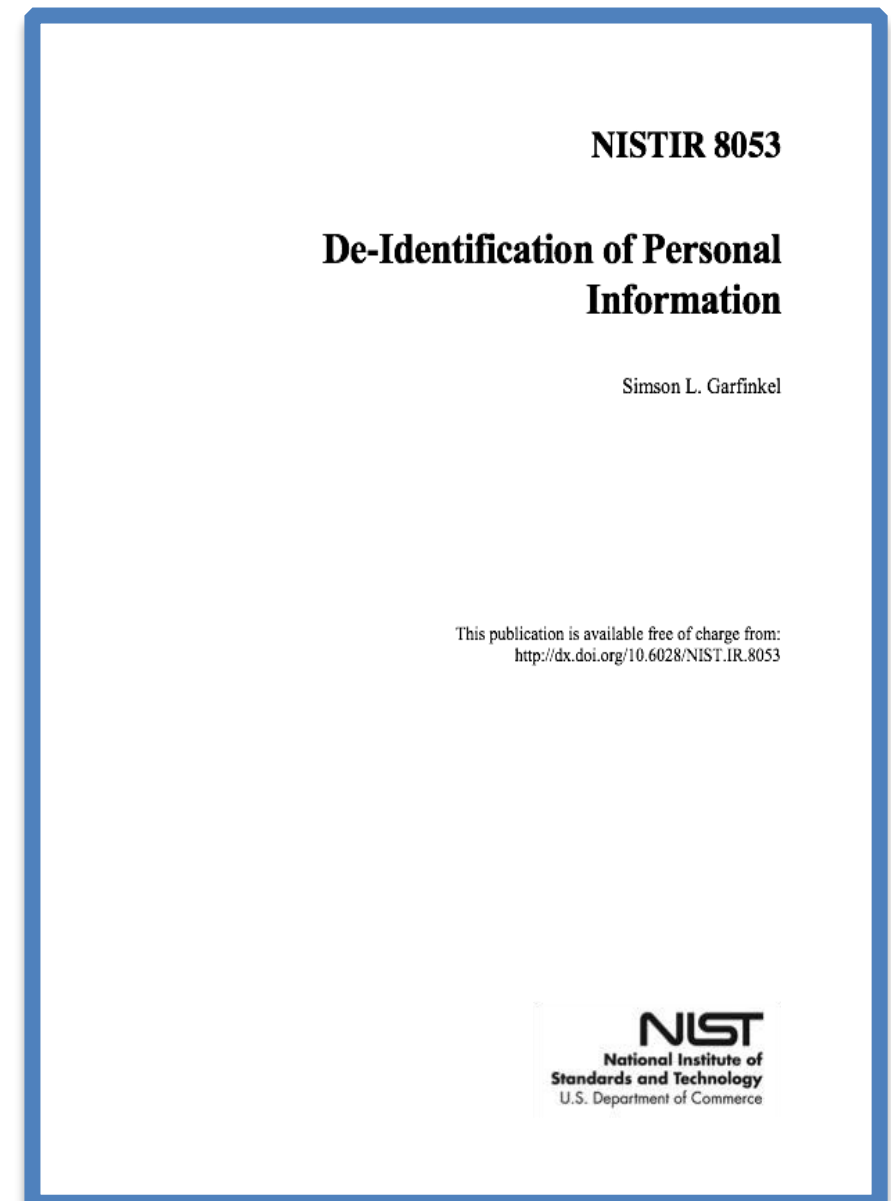# NISTIR 8053:
# De-Identification of Personal Information

## Covers:

- Why de-identify?
- De-identification terminology
- Famous re-identification cases
- De-identifying and re-identifying *structured data*
  *(e.g. survey data, Census data, etc.)*
- Challenges with de-identifying *unstructured data*
  *(e.g. medical text, photographs, medical imagery, genetic information)*

http://nvlpubs.nist.gov/nistpubs/ir/2015/NIST.IR.8053.pdf

October 2015

vi+46 pages

**NISTIR 8053**

**De-Identification of Personal Information**

Simson L. Garfinkel

This publication is available free of charge from:
http://dx.doi.org/10.6028/NIST.IR.8053

**NIST**
National Institute of
Standards and Technology
U.S. Department of Commerce

# Today there is a significant and growing interest in de-identification.



Controlled Sharing



Open Science



Data Publishing

# Big-data is not a new science—it's the future of all science.



"… Qualified researchers from many organizations will, with appropriate protection of participant confidentiality, have access to the cohort's **de-identified data** for research and analysis."

Request for Information: NIH Precision Medicine Cohort
NOT-OD-15-096
https://grants.nih.gov/grants/guide/notice-files/NOT-OD-15-096.html

# Under the current HIPAA Privacy Rule, de-identified Protected Health Information can be distributed without restriction.



https://en.wikipedia.org/wiki/Medical_record

**Medical Records**

**De-identification**

✗ name
✗ address
✗ birthday
✗ medical record number
*etc.*



https://commons.wikimedia.org/wiki/File:Applications-internet.svg

**Public Internet**

# Interest in de-identification extends far beyond healthcare.


https://www.flickr.com/photos/usdagov/4423599680

Social Science Data


https://pixabay.com/en/credit-card-bill-bank-statement-1104961/

Consumer Financial Data



Website
"We will never share your personal information..."

National Institute of Standards and Technology / U.S. Department of Commerce

# De-identified data can be re-identified



Sometimes data are not properly de-identified.

# De-identified data can be re-identified



| 24 | Grover Cleveland |
|----|------------------|
| 25 | William McKinley |
| **26** | **Theodore Roosevelt** |
| 27 | William Howard Taft |
| 28 | Woodrow Wilson |

Sometimes de-identified data can be *linked* to another dataset

# Simple statistics can be identifying.

| Title | Age | Sex | Address | ICD-10 | Diagnosis |
|-------|-----|-----|---------|--------|-----------|
| … | | | ██████ | | … |
| Lab Tech | 35 | M | ██████ | K25.0 | Gastric Ulcer with hemorrhage |
| Lab Tech | 56 | F | ██████ | J00 | Acute nasopharyngitis [Common Cold] |
| Professor | 35 | M | ██████ | C64.1 | Malignant neoplasm of right kidney |
| Professor | 69 | F | ██████ | C64.1 | Malignant neoplasm of right kidney |
| Contracts Specialist | 52 | F | ██████ | L30.9 | Dermatitis, unspecified [Eczema] |
| University President | 56 | F | ██████ | C64.1 | Malignant neoplasm of right kidney |
| . . . | | | ██████ | | . . . |

Hypothetical dataset from university healthcare system

# Re-identified information can link with other data.

## Research Database:

Patient 234-334-11
Diagnostic Codes: A98.4, J00, L30.9
…

Patient 234-334-11
Age: 35
Genetic History. …

Patient 234-334-11
Psychological Records
…

Patient 234-334-11
Social Services History
…

…

**NEWS**

| Ebola Patients | | | ICD-10 | Diagnosis |
|---|---|---|---|---|
| Alice | 30 | F | A98.4 | Ebola |
| Bob | 35 | M | A98.4 | Ebola |
| Carol | 40 | F | A98.4 | Ebola |

# Techniques for limiting identity disclosure:

| Title | Age | Sex | Address | ICD-10 | Diagnosis |
|---|---|---|---|---|---|
| University President | 56 | F | ████████████ | C64.1 | Malignant neoplasm of right kidney |

**Generalization:** University President ⇒ Senior Administrator
Age: 56 ⇒ Age: 50-59

**Field Swapping:** Age: 52 ⇒ Age: 56
Age: 56 ⇒ Age: 52

**Noise Addition:** University President ⇒ VP Finance
Age: 56 ⇒ Age: 58 ±5

**Suppression:** University President ⇒ XXXXXXXXXXXXX
Age: 56 ⇒ Age: XXX

# Lowering identifiability lowers data quality.



Identified & high quality.

Poor privacy protection.

Naïve de-identification

All quality gone.
Bad Science,
Bad Decision

Data Quality

Identifiability

National Institute of Standards and Technology / U.S. Department of Commerce

# HIPAA Privacy Rule "Safe Harbor" Provision:
## Medical records are de-identified if 18 data elements are removed

## Direct Identifiers:
— *Names*

— *Individual numbers: phone, fax, SSN, medical record, account #s, etc.*

— *Email addresses, IP address, URLs*

— *Biometrics: fingerprints, voiceprints, photographs, etc.*

— *Any other uniquely identifying number, characteristic or code.*

## Indirect Identifiers:
— *Geographic subdivisions smaller than a state, except first 3 digits of ZIP, provided the combined ZIP codes contain more than 20,000 people.*

— *Dates directly related to an individual (except for "age 90 or older")*

# Geographic information requires special attention

| Indirect identifiers | | | Direct identifier | | |
|---|---|---|---|---|---|
| **Title** | **Age** | **Sex** | **Address** | **ICD-10** | **Diagnosis** |
| Lab Tech | 35 | M | 100 Utah St. Anytown, 20124 | K25.0 | Gastric Ulcer with hemorrhage |
| Lab Tech | 56 | F | 653 Pleasant St. Uptown, 20321 | J00 | Acute nasopharyngitis [Common Cold] |
| Professor | 35 | M | 564 Main St. Nassis, 25312 | T25.332S | Burn of third degree of left toe |
| Professor | 69 | F | 202 Sky Lane Katap, 20134 | C64.1 | Malignant neoplasm of right kidney |
| Contracts Specialist | 52 | F | 956 Diablo Rd. Quirky, 23990 | L30.9 | Dermatitis, unspecified [Eczema] |
| University President | 56 | F | 451 Termo Dr. Boltz, 25333 | C64.1 | Malignant neoplasm of right kidney |

Hypothetical dataset from university healthcare system

# Safe Harbor allows ZIP3
## (assuming there are 20,000 people living in the area)

| Indirect identifiers | | | Direct identifier | | |
|---|---|---|---|---|---|
| **Title** | **Age** | **Sex** | **Address** | **ICD-10** | **Diagnosis** |
| Lab Tech | 35 | M | 201XX | K25.0 | Gastric Ulcer with hemorrhage |
| Lab Tech | 56 | F | 203XX | J00 | Acute nasopharyngitis [Common Cold] |
| Professor | 35 | M | 253XX | T25.332S | Burn of third degree of left toe |
| Professor | 69 | F | 201XX | C64.1 | Malignant neoplasm of right kidney |
| Contracts Specialist | 52 | F | 239XX | L30.9 | Dermatitis, unspecified [Eczema] |
| University President | 56 | F | 253XX | C64.1 | Malignant neoplasm of right kidney |

Hypothetical dataset from university healthcare system

# Results of the 2010 Office of the National Coordinator for Health Information Technology Safe Harbor Re-Identification Test:



**15,000 Hispanic Patients**

216 distinct by Sex, ZIP3 & age

**30,000 Records from InfoUSA**

84 distinct by sex, ZIP3 & age

20 match on sex, ZIP3 & age

2 actual matches on last name, street address, and phone

Data from 2004-2009

# K-anonymity: assure at least "k" records have the same set of indirect identifiers.

| Indirect identifiers | | | Direct identifier | | |
|---|---|---|---|---|---|
| **Title** | **Age** | **Sex** | **Address** | **ICD-10** | **Diagnosis** |
| Lab Tech | 35 | M | 201XX | K25.0 | Gastric Ulcer with hemorrhage |
| Lab Tech | 56 | F | 203XX | J00 | Acute nasopharyngitis [Common Cold] |
| Professor | 35 | M | 253XX | T25.332S | Burn of third degree of left toe |
| Professor | 69 | F | 201XX | C64.1 | Malignant neoplasm of right kidney |
| Contracts Specialist | 52 | F | 239XX | L30.9 | Dermatitis, unspecified [Eczema] |
| University President | 56 | F | 253XX | C64.1 | Malignant neoplasm of right kidney |

*Color background indicates values modified for k=2 k-anonymity*

# "Tiger Teams" are another way to test re-identification.



Re-identification

15,000
Hispanic
Patients

216 distinct by
Sex, ZIP3 & age

20 match
on sex,
ZIP3 & age

30,000
Records from
InfoUSA

84 distinct by
sex, ZIP3 & age

2 actual matches on last name,
street address, and phone

Data from 2004-2009

Estimated Re-identification rate:
<u>No</u> verification: 20 in 15,000
Verification: 2 in 15,000

National Institute of Standards and Technology / U.S. Department of Commerce

# Re-identification tests assume data available to match. As more data become available, re-identification gets easier.



300M records from provider "C" (Social Media?)

15,000 Hispanic Patients

30,000 Records from InfoUSA

1M records From provider "B"

**Hypothetical Example**

# A constellation of diseases can be an identifier



**De-identified medical records from provider "N"**

Smallpox

Concussion at age 9

Malaria at age 21

Depression

Fractured jaw

Linked records

**Hypothetical Example**

# Atreya, Smith, McCoy, Malin & Miller (2013) "Reducing patient re-identification risk for laboratory results within research datasets."

Medical tests



Database with linked test results

A single identified blood test can be the link to dozens of de-identified records

# Blood tests can be de-identified by adding noise

## Example Lab Report

### Patient copy

1 University Medical Center, Dept. of Pathology
123 University Way, City, ST 12345    02/14/2008  16:13    2

Doe, Mr. John Q.  3

Patient ID No. 987654321  3    D.O.B. 01/01/1941    67Y
Ordering MD: Smith, Peter MD  4    Physician Copy for Dr: Smith, Jane
PT medications: multivitamins  5

Specimen(s) Collected: 2/10/08 14:30    Lab Acc'n No. 2234
Specimen: Serum    Date Reported: 2/10/08 16:
Comments: Specimen is non-fasting ; sl. hemolysis

| Test Name | Patient's Results | Ref. Range | Units |
|---|---|---|---|
| BMP | | | |
| Na | L124 | 136-145 | mEq/L |
| K | H5.8 | 3.5-5.1 | mEq/L |
| CO2 | 25 | 23-29 | mEq/L |
| Cl | 101 | 98-107 | mEq/L |
| Glucose | H107 | 74-100 | mg/dL |
| Ca | 10.1 | 8.6-10.2 | mg/dL |
| BUN | 17 | 8-23 | mg/dL |
| Creatinine | 0.9 | 0.8-1.3 | mg/dL |

Key: L=Abnormal Low, H=Abnormal High, WNL=Within Normal Limits, *=critical value

Na:        124        $\Rightarrow$ 126
K:         5.8        $\Rightarrow$ 5.9
$CO_2$:    25         $\Rightarrow$ 24
Cl:        101        $\Rightarrow$ 104
Glucose: 107    $\Rightarrow$ 110
Ca:        10.1       $\Rightarrow$ 9.9
BUN:    17         $\Rightarrow$ 17
Creatinine: 0.9 $\Rightarrow$ 1.0

(values for demonstration only)

### Research database

ology / U.S. Department of Commerce

# "Differential Privacy" adds systematic noise to query results



**Data Enclave**

Real Data
+
Computation

**Query**

**Result**

*Synthetic Data*

The Common Data Project
Private Map Maker

*Key concepts: Privacy Budget & Noise*

# The Census Bureau distributes synthetic data to protect privacy while preserving some data quality.

# Can synthetic datasets designed to enable research also be used to promote accountability and transparency?


Animated encounter data


Synthetic tabular data


Body-worn camera video with replaced faces

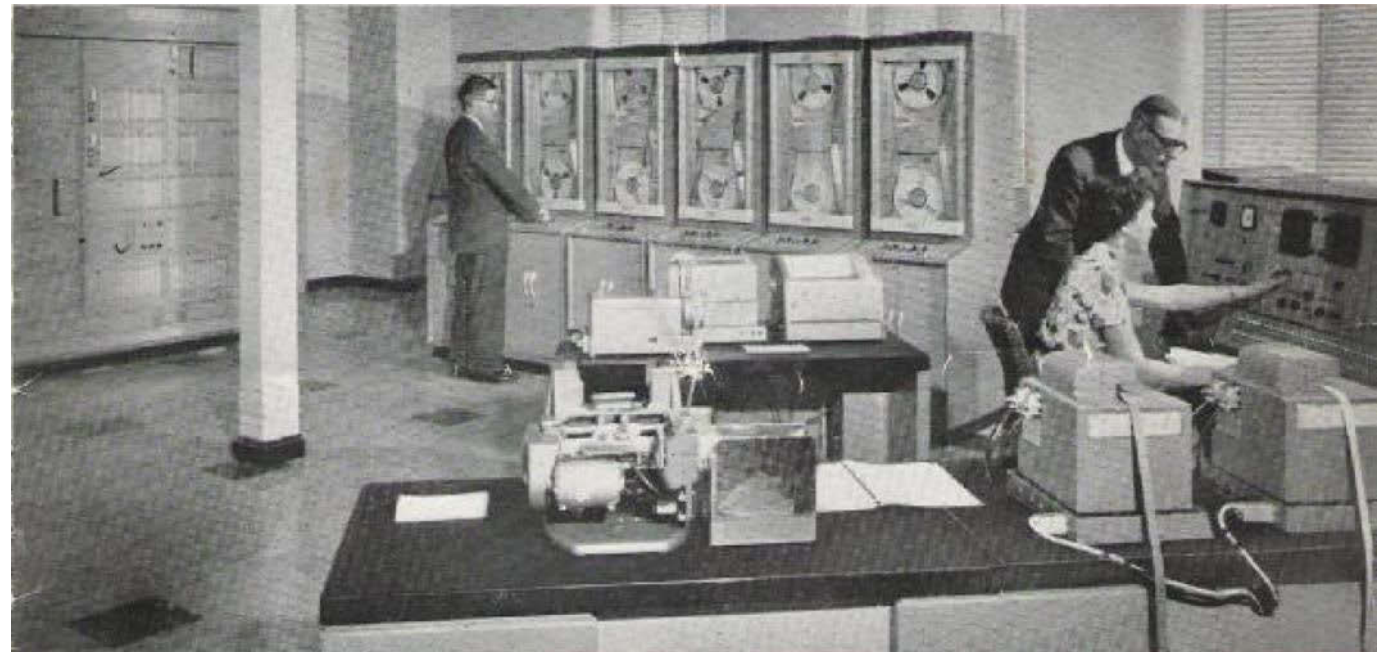# De-identification strategies should be formally evaluated.

Do they meet the stated policy goals?

Does the software faithfully implement the stated algorithm?

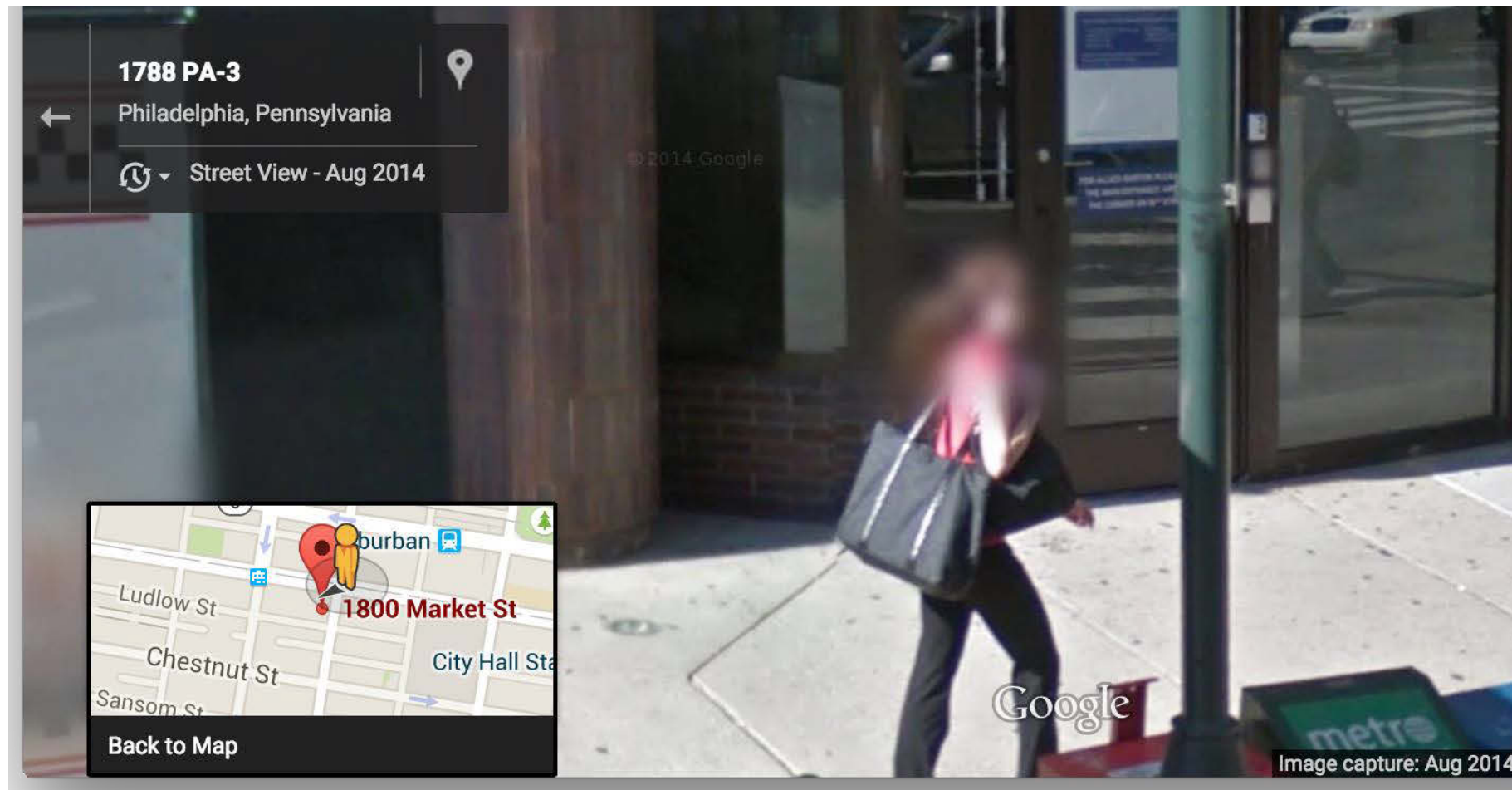Are the statistical privacy guarantees actually met?

Is the necessary training in place?

Will there be monitoring and auditing?



https://en.wikipedia.org/wiki/EMIDEC_1100

# De-identification of non-tabular data poses special problems.



1788 PA-3
Philadelphia, Pennsylvania

Street View - Aug 2014

Ludlow St
Chestnut St
Sansom St

1800 Market St
City Hall Sta

Back to Map

Image capture: Aug 2014



AYAKTA

http://www.randomhistory.com/photos/2014/scoliosis-xray.jpg

Google claims 90% of faces and 95% of license plates removed through automated processing.
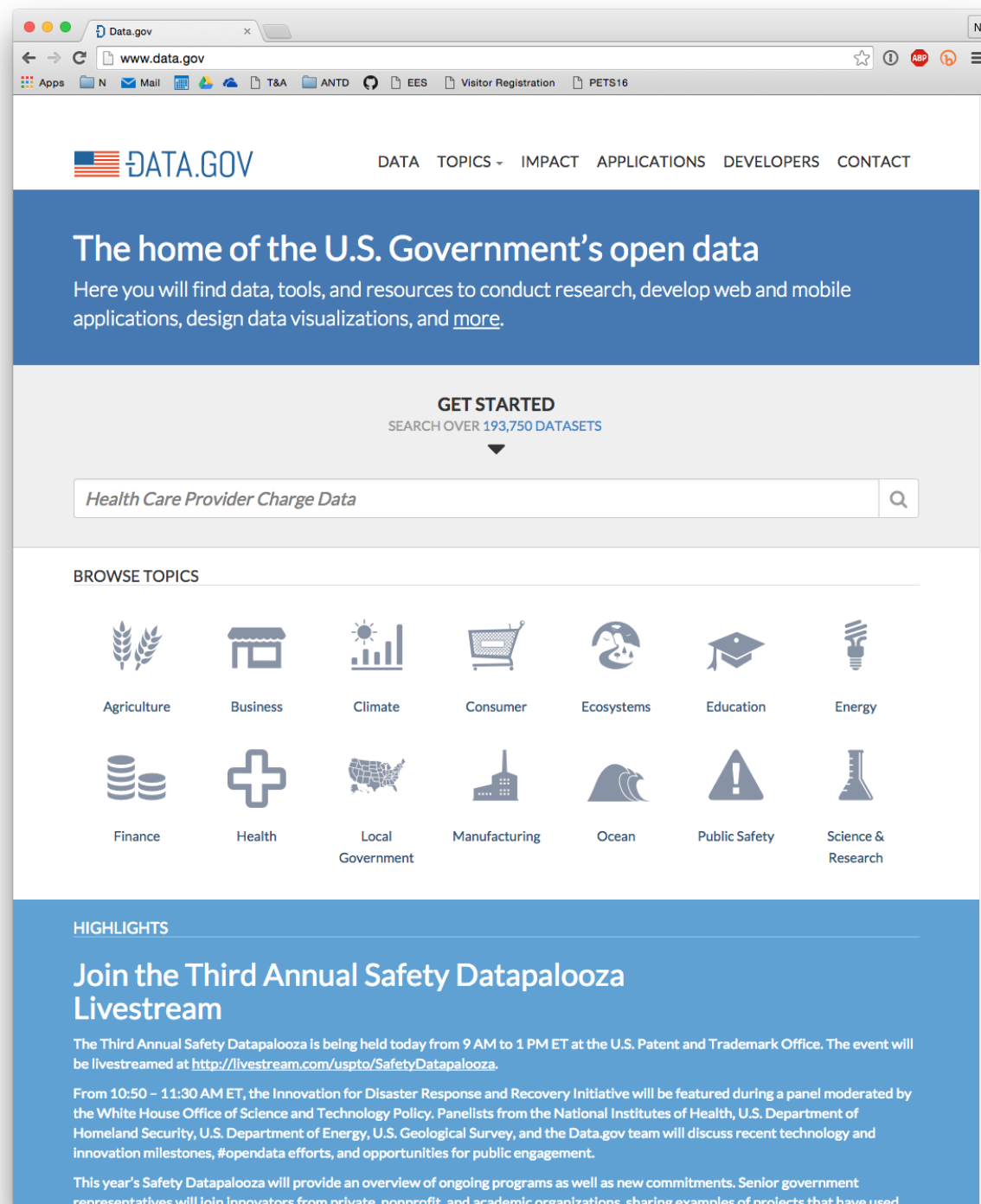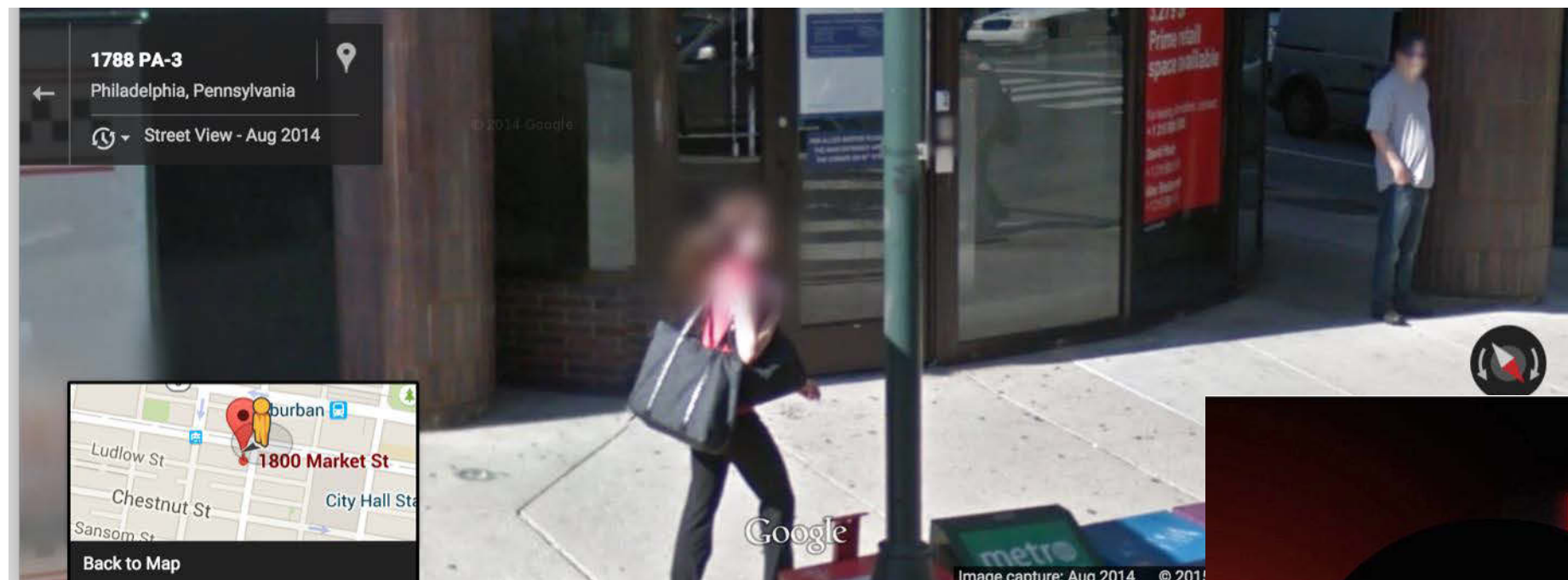
**Medical imagery can be highly identifying.**

# More research is needed to determine if systems can protect privacy and allow for unlimited use of data.
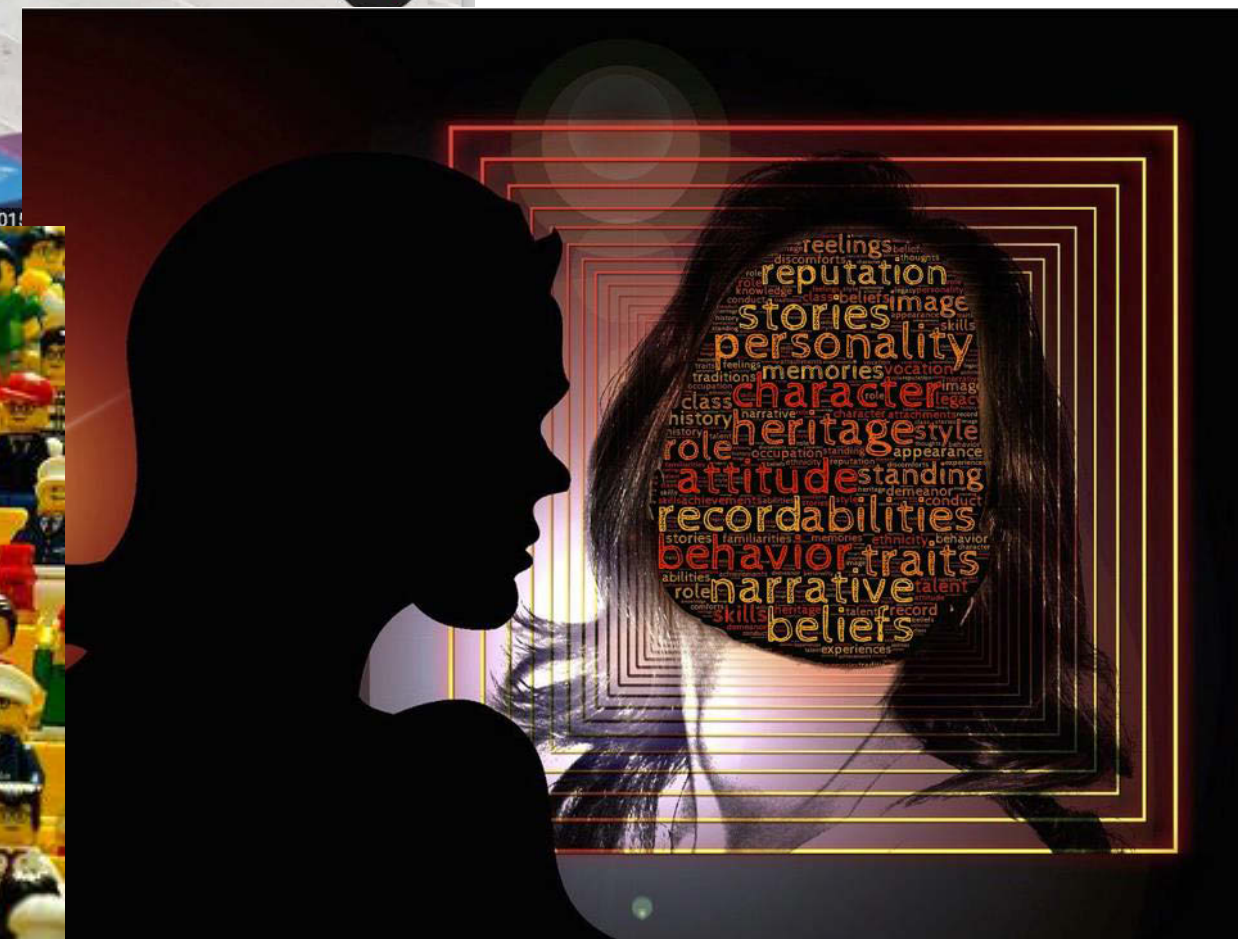
# Can raw data be transformed so completely that individuals cannot recognize their own data once they are in a crowd?



https://pixabay.com/en/lego-doll-the-per-amphitheatre-1044891/

https://pixabay.com/en/self-self-image-image-identity-792365/

National Institute of Standards and Technology / U.S. Department of Commerce

# In summary:
# We have learned a lot about de-identification in recent years.

The de-identification "toolkit" has several options
- suppression, generalization        *commonly used in healthcare*
- field swapping, noise addition      *commonly used in vital statistics*

K-anonymity and Differential Privacy are formal models for evaluating the quality of de-identification

We increasingly have the ability to:
- Modify data so that the data subjects' identity is rem
  leaving information that is somewhat useful.
- But the more useful it is,
  the more likely it can be re-identified

We need procedures for:
- Evaluating the effectiveness of de-identification
- Evaluating the usefulness of the data that remain.

We need these techniques for a wide range of
- Structured data, text, medical, video

**Lowering identifiability lowers data quality.**

National Institute of Standards and Technology / U.S. Department of Commerce