Team for Research in Ubiquitous Secure Technology

*Women's Institute in Summer Enrichment*

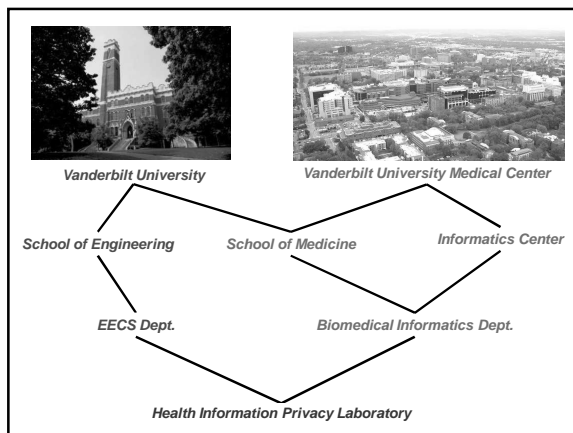## Building Systems to Support Health Information Privacy

Brad Malin, Ph.D.
Assistant Professor
Department of Biomedical Informatics, School of Medicine
Department of Computer Science, School of Engineering
Vanderbilt University

1

---

# Disclaimer

- **Privacy** is an overloaded word

- Today: **privacy** in the context of a specific domain

- Healthcare
  - Health Insurance Portability & Accountability Act (HIPAA)
  - NIH Data Sharing Policy
  - NIH Genome Wide Association Study Data Sharing Policy

2

---



Vanderbilt University

Vanderbilt University Medical Center

School of Engineering    School of Medicine    Informatics Center

EECS Dept.    Biomedical Informatics Dept.

Health Information Privacy Laboratory

---

# What's Going On?

- We study "privacy" in various operational realms
  - Primary Care
    - Clinical Information Systems Design
    - "Intelligent" Auditing

  - Secondary Uses
    - De-identification / Re-identification / Anonymization
    - Secure Data Integration and Analysis

---

# Privacy Everywhere

- We do not always control who gets, and has access to, our information

- Legally, however, data collectors may be required to maintain your privacy



Data Collecting

Data Using

Data Sharing

---

# Privacy Everywhere

- Let's begin with data already in the system



Data Collecting

Data Using

Data Sharing

## Electronic Medical Records – Hooray!

- At V
  - [illegible bullets]

- And
  - [illegible]

- Incr
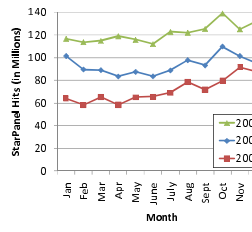  over [illegible]



## Clinical Information Systems Design
(Duncavage 2007; Mathe et al, 2008; Werner et al, 2007, 2008)
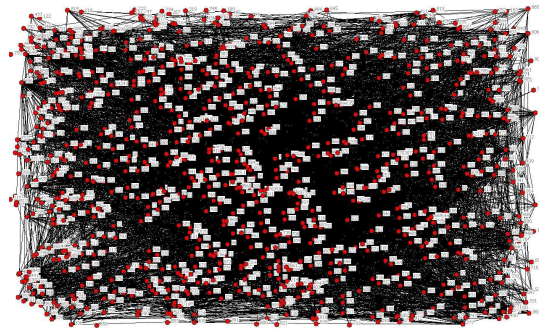


## Surveillance (Paulett, Malin)

Snooped

- Very little role-based access control in large academic medical centers! (why?)

- Most auditing is done manually! (why?)
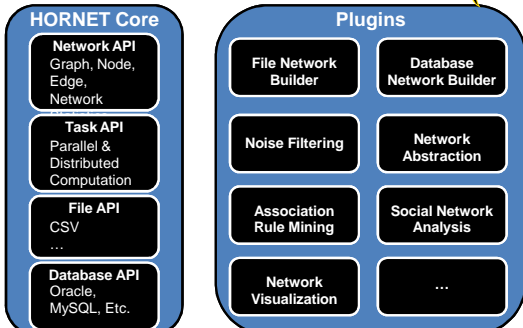
> 1.5 million patient records

> 20,000 authorized users



## Jan 1, 2006

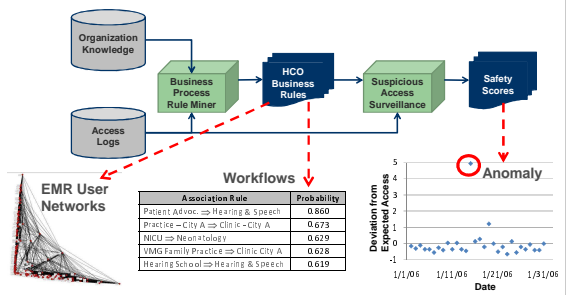- Users linked if accessed 1 common patient (~ 900 users, 2000 patients)



## HORNET: Healthcare Organizational Research Toolkit
(http://code.google.com/p/hornet/)

**HORNET Core**

- **Network API** Graph, Node, Edge, Network
- **Task API** Parallel & Distributed Computation
- **File API** CSV …
- **Database API** Oracle, MySQL, Etc.

**Plugins**

- **File Network Builder**
- **Database Network Builder**
- **Noise Filtering**
- **Network Abstraction**
- **Association Rule Mining**
- **Social Network Analysis**
- **Network Visualization**
- …

## Learning Policies and Anomalies from EMR Access Logs

- Extracts patterns from medical record access logs to model policies & detect "privacy" violations.



**Workflows**

| Association Rule | Probability |
| --- | --- |
| Patient Advoc. ⇒ Hearing & Speech | 0.860 |
| Practice – City A ⇒ Clinic – City A | 0.673 |
| NICU ⇒ Neonatology | 0.629 |
| VMG Family Practice ⇒ Clinic City A | 0.628 |
| Hearing School ⇒ Hearing & Speech | 0.619 |

## What's Going On?

- "Privacy" in various operational settings
  - Primary Care
    - Clinical Information Systems Design
    - "Intelligent" Auditing

  - Secondary Uses
    - De-identification / Re-identification / Anonymization
    - Secure Data Integration and Analysis

---



THE TENNESSEAN

TUESDAY, JUNE 20, 2006

VU to put patient DNA in vast research pool

---

## Information Integration



Electronic Medical Record System
- 80M entries on >1.5M patients

Clinical Notes · CPOE Orders (Drug) · ICD9, CPT · Clinical Messaging · Test Results

Discarded blood
- 50K per year

Updated Weekly

Extract DNA

Clinical Resource

---

## Research Support & Data Collection



Genotyping, genotype-phenotype relations

cases

controls

Investigator query

cases

controls

scrubbed

Data analysis

---

## Holy Hand Grenades! How Did You…

- Initially an institutionally funded project

- Office for Human Research Protections designation as Non-Human Subjects Research under 45 CFR 46 ("HIPAA Common Rule")*
  - *Samples & data not linked to identity*
  - Conducted with IRB & ethics oversight

*Roden D et al. Development of a large-scale de-identified DNA biobank to enable personalized medicine. Clin Pharmacol Ther. 2008; 84(3): 362-369.

---

## HIPAA
### (the ELEPHANT in the room)

- Health Insurance Portability & Accountability Act

- Resolve state laws hampering standardization, transfer, & sharing of health information

- "covered entity" cannot use or disclose protected health information (PHI)
  - data "explicitly" linked to a particular individual, or
  - could reasonably be expected to allow individual identification

## HIPAA - Secondary Data Sharing

- Safe Harbor

- Limited Release

- Statistical or Scientific Standard

## HIPAA Safe Harbor

- Data that can be given away without oversight
- Requires removal of eighteen attributes
  - Names / Initials
  - Street address, city, county, precinct code and equivalent geocodes
  - All elements of dates, except year, and all ages over 89
  - #'s: Phone, Fax, Social Security, Medical Record, Health Plan ID, Account, License, Serial, Device
  - Web: Email, URL, IP addresses
  - Biometric identifiers: finger, voice prints
  - Full face photo images and comparable images
  - Any other unique identifying number, characteristic, or code
    - A code is an identifier if the person holding the coded data can re-identify the individual

## HIPAA Limited Dataset

- Includes more specific information than Safe Harbor
- Can include
  - Dates of birth, death, service
  - Geographic Info: Town, Zip code, County

- **Requires Contract:** Research entity provides assurances that it will not use or disclose the information for purposes other than research and will not identify or contact the individuals who are the subjects

## HIPAA Statistical / Scientific Standard

- Certify via "generally accepted statistical and scientific principles and methods, that the risk is very small that the information could be used, alone or in combination with other reasonably available information, by the anticipated recipient to identify the subject of the information."

- "Must document the methods and results of the analysis that justify such a determination"

- "Must not disclose the key or other mechanism that would have enabled the information to be re-identified"
  - includes pseudo-random number algorithms and seed values

## "Scrubbing" Medical Records

MR# is removed

Replaced SSN and phone #

**Rules***
**Regular Expressions**
**Dictionaries**
**Exclusions**

↓↓↓↓↓

**Machine Learning – Conditional Random Fields****

*Gupta D, et al. Evaluation of a deidentification (De-Id) software engine to share pathology reports and clinical documents for research. Am J Clin Pathol. 2004 Feb; 121(2): 176-186.*
*Wellner et al. Rapidly retargetable approaches to de-identification in medical records. Journal of the American Medical Informatics Association. 2007.*

## "Scrubbed" Medical Record

MR# is removed

Replaced SSN and phone #

Unknown residual re-identification potential (e.g. "the mayor's wife")

Substitute names

Shifted Dates

## Technology + Policy

- Databank access restricted to Vanderbilt employees

- Must sign use agreement that prohibits "re-identification"

- Operations Advisory Board and Institutional Review Board approval needed for each project

- All data access logged and audited per project

## To Vanderbilt …. and Beyond

### The eMERGE Network
#### electronic Medical Records & Genomics
*A consortium of biorepositories linked to electronic medical records data for conducting genomic studies*

- **Consortium members (http://www.gwas.net)**
  - □ Group Health of Puget Sound (UW)   □ Mayo Clinic
  - □ Marshfield Clinic   □ Northwestern University
  - □ Vanderbilt University

- **Funding condition: contribute <u>de-identified</u> genomic and EMR-derived phenotype data to database of genotype and phenotype (dbGAP) at NCBI, NIH**

## Data Sharing Policies

- **Feb '03: National Institutes of Health Data Sharing Policy**
  - *"data should be made as widely & freely available as possible"*
  - *researchers who receive >= $500,000 must develop a data sharing plan or describe why data sharing is not possible*
  - Derived data must be shared in a manner that is devoid of "identifiable information"

- **Aug '06: NIH Supported Genome-Wide Association Studies Policy**
  - □ Researchers who received >= $0 for GWAS

# The Face that Launched a Thousand Ships

29

## Healthcare Reform At Work

- In 1997, approx. 44 of 50 states collected and disseminated hospital discharge data
- In 2005, approx. 47 of 50 states " "
- Attributes recommended by *National Association of Health Data Organizations* for disclosure

| | |
|---|---|
| – Patient Zip Code | □ Principle Diagnosis Codes (ICD-9) |
| – Patient Birth Date | □ Procedure Codes |
| – Patient Gender | □ Physician ID Number |
| – Patient Racial Background | □ Physician Zip Code |
| – Patient Number | □ Total Charges |
| – Visit Date | |

J. Schoenman et al. The value of hospital discharge databases. NORC & NAHDO. 2005.
http://www.hcup-us.ahrq.gov/reports/final_report.pdf

### Case Study – "Quasi-identifier"

*Back in the '90s*

Ethnicity

Visit date    Zip Code

Diagnosis    Birthdate

Procedure    Gender

Medication

Total charge

Hospital Discharge Data

L. Sweeney. Journal of Law, Medicine, and Ethics. 1997.    31

### Case Study – "Quasi-identifier"

*Back in the '90s*

Name

Address

Date registered    Zip Code

Party affiliation    Birthdate

Date last voted    Gender

City of Cambridge, MA Voter Registration Records

L. Sweeney. Journal of Law, Medicine, and Ethics. 1997.    32

### Case Study – "Quasi-identifier"

Re-identification of William Weld

Ethnicity          Name
Visit date         Address
Diagnosis   Zip Code   Date registered
Procedure   Birthdate  Party affiliation
Medication  Gender     Date last voted
Total charge

Hospital          Voter List
Discharge Data

L. Sweeney. Journal of Law, Medicine, and Ethics. 1997.    33

### 5-Digit Zip Code
### + Birthdate
### + Gender

**63-87% of US estimated to be unique**

L. Sweeney. Uniqueness of Simple Demographics in the U.S Population. 2000.

P. Golle. Revisiting the Uniqueness of US Population. ACM WPES. 2006.    34

# And Now, It's A Phenomenon!!!

### The AOL Search Log Case (2006)

- Goal: Support web search research
- 650k customers, 20 million queries, 3 month period
- Names replaced with persistent pseudonyms

| Name | Query | Date | Time |
|---|---|---|---|
| John Doe | Books | 1/2/05 | 16:52 |
| Bob Smith | Payscale | 1/4/05 | 23:41 |
| John Doe | Porn | 1/8/05 | 03:15 |

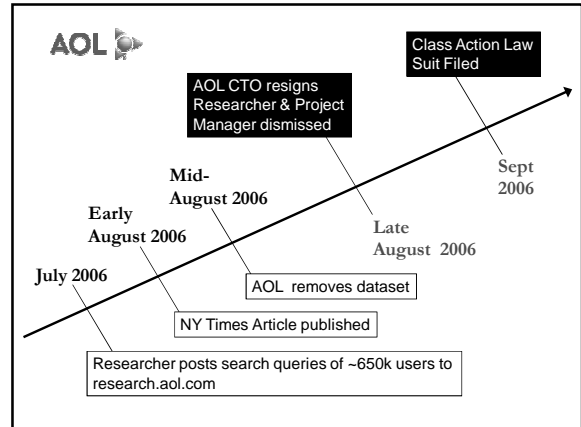| User | Query | Date | Time |
|---|---|---|---|
| 8123 | Books | 1/2/05 | 16:52 |
| 9010 | Payscale | 1/4/05 | 23:41 |
| 8123 | Porn | 1/8/05 | 03:15 |

Barbaro & Zeller. A face exposed for AOL searcher no. 4417749. **New York Times**.  Aug 9, 2006.
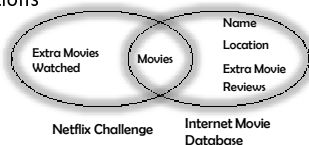
Thelma Arnold & Dudley

---



Class Action Law Suit Filed

AOL CTO resigns
Researcher & Project Manager dismissed

Mid-August 2006

Sept 2006

Early August 2006

Late August  2006

July 2006

AOL  removes dataset

NY Times Article published

Researcher posts search queries of ~650k users to research.aol.com

---

## The Netflix Challenge (2008-2009)

- Netflix published movie selections of ~450,000 pseudonymized subscribers
- Re-identification via uniqueness of movie combinations



Name
Location
Extra Movie Reviews

Extra Movies Watched

Movies

Netflix Challenge

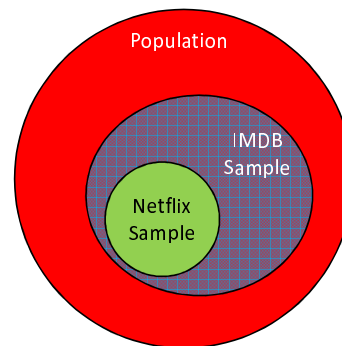Internet Movie Database

- Class action filed December 2009

A. Narayanan & V. Shmatikov *IEEE Security and Privacy Conference*. 2008.    39

---
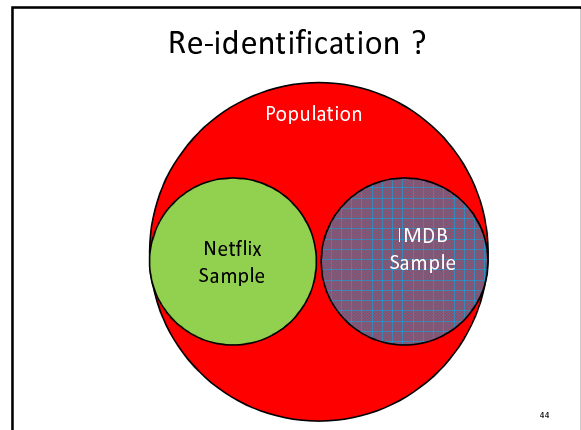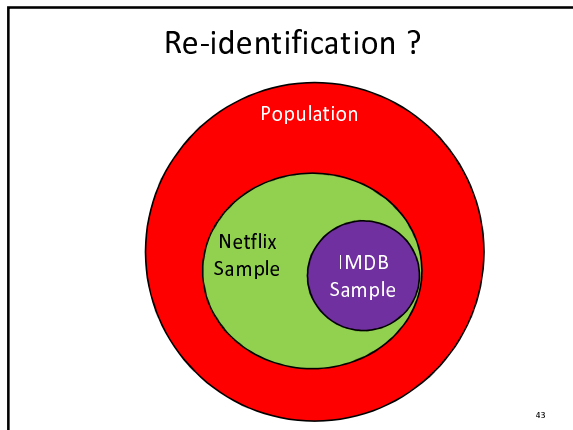


---

Samples
Samples
Samples
Samples
Samples
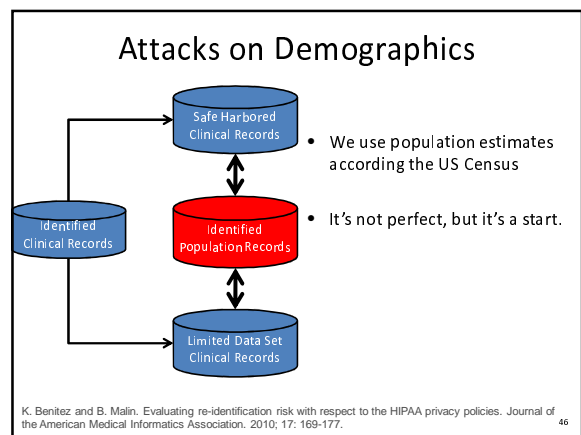and
Populations

41

---

## Re-identification ?

Population

IMDB Sample

Netflix Sample

42

## Re-identification ?



43

## Re-identification ?



44

# Now Back to Your Regularly Scheduled Programming

45

## Attacks on Demographics



- We use population estimates according the US Census

- It's not perfect, but it's a start.

K. Benitez and B. Malin. Evaluating re-identification risk with respect to the HIPAA privacy policies. Journal of the American Medical Informatics Association. 2010; 17: 169-177.

46

## The Census?    http://factfinder.census.gov/



# Beyond "unique"

48

## Case Study: Tennessee



Group size = 33

Limited Dataset
{Race, Gender, **Date** (of Birth), **County**}

Fraction of population identifiable under Limited Dataset

Fraction of population identifiable under Safe Harbor

Safe Harbor
{Race, Gender, **Year** (of Birth), **State**}

49

## All U.S. States

Safe Harbor        Limited Data set



50

## Policy Analysis via "Trust Differential"...

$$\frac{\text{Risk(Limited Dataset)}}{\text{Risk(Safe Harbor)}}$$

- Uniques
  - *Delaware's risk increases by a factor    ~1,000*
  - *Tennessee's    "    "    "    "    ~2,300*
  - *Illinois's "    "    "    "    "    ~65,000*

- ≤20,000
  - *Delaware's risk does not increase*
  - *Tennessee's risk increases by a factor of ~8*
  - *Illinois's risk increases by a factor of ~37*

51

## State Policy

| | IL | MN | TN | WA | WI |
|---|---|---|---|---|---|
| WHO??? | Registered Political Committees (ANYONE – In Person) | MN Voters | Anyone | Anyone | Anyone |
| Format | Disk | Disk | Disk | Disk | Disk |
| Cost | $500 | $46; "use ONLY for elections, political activities, or law enforcement" | $2500 | $30 | $12,500 |
| Voter ID | ● | ● | ● | ● | ● |
| Name | ● | ● | ● | ● | ● |
| Address | ● | ● | ● | ● | ● |
| Voter Status | ● | ● | ● | ● | ● |
| District Information | ● | ● | ● | ● | ● |
| Election History | ● | ● | ● | ● | ● |
| Date of Birth | ● | ○ | ● | ● | |
| Date of Registration | ● | | ● | ● | |
| Sex | ● | | ● | ● | |
| Race | | | ● | | |
| Phone Number | ● | ● | | | |

52

## Cost?

| | Limited Dataset | | Safe Harbor | |
|---|---|---|---|---|
| State | Marker Risk | Cost per Re-id | Total Risk | Cost per Re-id |
| VA | 3159764 | $0 | 221 | $0 |
| NY | 2905697 | $0 | 221 | $0 |
| SC | 2231973 | $0 | 1386 | $0 |
| WI | 72 | $174 | 2 | $6,250 |
| WV | 55 | $309 | 1 | $17,000 |
| NH | 10 | $827 | 1 | $8,267 |

53

## A Couple of Parting Thoughts

- The application of technology must be considered within the systems and operational processes they will be applied

- One person's vulnerability is another person's armor (variation in risks)

- It is possible to inject privacy into health information systems – but it must be done early (see "privacy by design")!

- Sometimes theory needs to be balanced with practicality

54