

# Security in a Distributed Network Environment

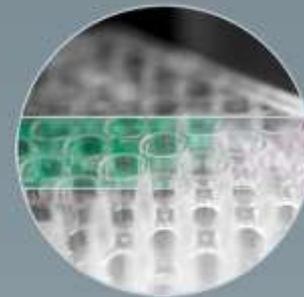
Kenneth R. Evans, PhD

CEO, OCBN;

Associate Professor, Dept of Pathology and Laboratory Medicine, Queen's University



Ontario Cancer  
Biomarker Network



# High Dimensional Medical Databases

- The opportunities are enormous...



# “Depressive Symptoms” Example

Autism

Depression

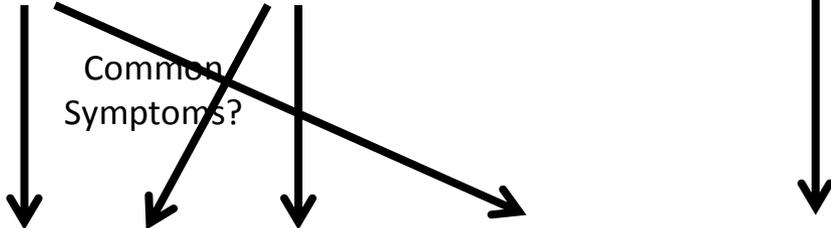
Neuro-degeneration

Cerebral Palsy

Addictions

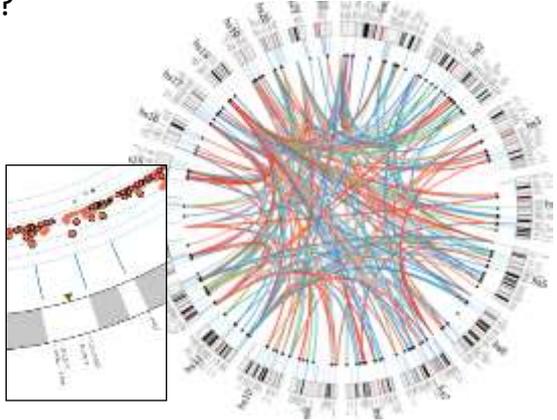
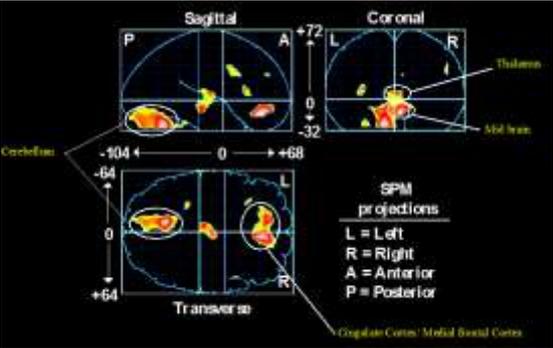
Epilepsy

Traumatic Brain Injury



“Worthlessness, Sadness, Helplessness”

Common mechanisms?



# Personalized medicine

- Changing what we mean by “disease”
- The opportunities are enormous...
- So, too are the challenges

# Security and Privacy

- Privacy is an enormous concern in medicine
  - But the full extent of the issue is only just beginning to be understood.
- Personal Health Information:
  - “Information in any format that identifies the individual, including demographic information collected from an individual that can reasonably be used to identify the individual.
  - Additionally, PHI is information created or received by a health care provider, health plan, employer, or health care clearinghouse; and relates to the past, present, or future physical or mental health or condition of an individual.”

HIPAA Guidance

# US Health Insurance Portability and Accountability Act (HIPAA): 18 “identifiers”

1. Name
2. Telephone numbers
3. Fax numbers
4. Electronic mail addresses
5. Social security numbers
6. Medical record numbers
7. Health plan beneficiary numbers
8. Account numbers
9. Certificate / license numbers
10. Vehicle identifiers and serial numbers, including license plate numbers
11. Device identifiers and serial numbers
12. Web Universal Resource Locators (URLs)
13. Internet Protocol (IP) address numbers
14. Location; all geographic subdivisions smaller than a state, including street address, city, county, precinct, zip code, and their equivalent geocodes.
15. Dates (all dates related to the subject of the information, e.g. birth dates, admission dates, discharge dates, encounter dates, surgery dates, etc.)
16. Biometric identifiers, including finger and voice prints
17. Full face photographic images and any comparable images
18. Any other unique identifying number, characteristic, or code

# Brain Image Data



- It is possible to reassemble an image of an individual from an MRI scan
- Studies have shown 30 – 70% success rates in identifying the individual from photographs
- Full face photographic images and ***any comparable images*** are PHI

...and evermore shall be so?

- The concept of PHI is evolving



# Erlich and Gymrek, 2013



## Identifying Personal Genomes by Surname Inference

Melissa Gymrek,<sup>1,2,3,4</sup> Amy L. McGuire,<sup>5</sup> David Golan,<sup>6</sup> Fran Halpern,<sup>7,8,9</sup> Yaniv Erlich<sup>1\*</sup>

Sharing sequencing data sets without identifiers has become a common practice in genomics. Here, we report that surnames can be recovered from personal genomes by profiling short tandem repeats on the Y chromosome (Y-STRs) and querying recreational genetic genealogy databases. We show that a combination of a surname with other types of metadata, such as age and state, can be used to triangulate the identity of the target. A key feature of this technique is that it entirely relies on free, publicly accessible Internet resources. We quantitatively analyze the probability of identification for U.S. males. We further demonstrate the feasibility of this technique by tracing back with high probability the identities of multiple participants in public sequencing projects.

Surnames are potentially inherited in most human societies, resulting in their co-segregation with Y-chromosomal haplotypes (Y-S). Based on this observation, multiple genetic genealogy companies offer services to remote distant paternal relatives by genotyping a few dozen

highly polymorphic short tandem repeats across the Y chromosome (Y-STRs). The association between surnames and haplotypes can be confounded by migration events, mutations, and adoptions of the same surname by multiple founders (5). The genetic genealogy community addresses these barriers with massive databases that list the test results of Y-STR haplotypes along with their corresponding surnames. Currently, there are at least eight databases and surname project Web sites that collectively contain hundreds of thousands of surname-haplotype records (table S1).

The ability of genetic genealogy databases to breach anonymity has been demonstrated in the past. In a number of public cases, male adoptees and descendants of anonymous sperm donors used recreational genetic genealogy services to genotype their Y-chromosomal haplotypes and to search the companies' databases (6-9). The genetic matches identified distant paternal relatives and pointed to the potential surnames of their biological fathers.

By combining other pieces of demographic information, such as date and place of birth, they fully exposed the identity of their biological fathers. Lander *et al.* (10) were the first to speculate that this technique could expose the full identity of participants in sequencing projects. Gschler (11) empirically approached this hypothesis by testing 31 Y-STR haplotypes of CEU participants in those databases and reported that potential surnames can be detected. [CEU participants are multigenerational families of northern and western European ancestry in Utah who had originally had their samples collected by CEPH (Centre d'Etude du Polymorphisme Humain) and were later recruited to participate in the HapMap project.] However, these surnames could reach thousands of individuals, and the study did not pursue full re-identification at a single-person resolution.

Our goal was to quantitatively approach the question of how readily surname inference might be possible in a more general population, apply this approach to personal genome data sets, and demonstrate end-to-end identification of individuals with only public information. We show that full identities of personal genomes can be imposed via surname inference from recreational genetic genealogy databases followed by Internet searches. In all cases in which individuals were studied who had donated DNA samples, the inferred consent statements they had signed stated privacy breach as a potential risk and the data usage terms did not prevent re-identification. Representatives of relevant organizations that funded the original studies were notified and confirmed the completion of this study with their guidelines (12).

As a primary resource for surname inference, we focused on Ysearch ([www.ysearch.org](http://www.ysearch.org)) and

- Accessed publicly available information from the 1000 Genomes project, as well as genealogy sites
- Successfully identified 50 individuals
- This is the reality: if we allow even “de-identified” data (ie without personal identifiers) into an open access environment, it will be possible to re-identify them.

Downloaded from www.sciencemag.org

<sup>1</sup>Harvard Medical School, Center for Genomic Medicine, Boston, MA 02115, USA. <sup>2</sup>Harvard-MIT Health Care and Biomedical Innovation Center, Cambridge, MA 02138, USA. <sup>3</sup>Harvard-MIT Center for Population and Family Studies, Cambridge, MA 02138, USA. <sup>4</sup>Department of Molecular Biology and Statistics Unit, Massachusetts General Hospital, Boston, MA 02114, USA. <sup>5</sup>Center for Medical Ethics and Health Policy, Baylor College of Medicine, Houston, TX 77030, USA. <sup>6</sup>Department of Statistics and Operations Research, Tel Aviv University, Tel Aviv 61020, Israel. <sup>7</sup>School of Computer Science, Tel Aviv University, Tel Aviv 61020, Israel. <sup>8</sup>Department of Molecular Biology and Biotechnology, Tel Aviv University, Tel Aviv 61020, Israel. <sup>9</sup>The International Computer Science Institute, Berkeley, CA 94704, USA.

\*To whom correspondence should be addressed. E-mail: [yerlich@alum.mit.edu](mailto:yerlich@alum.mit.edu)

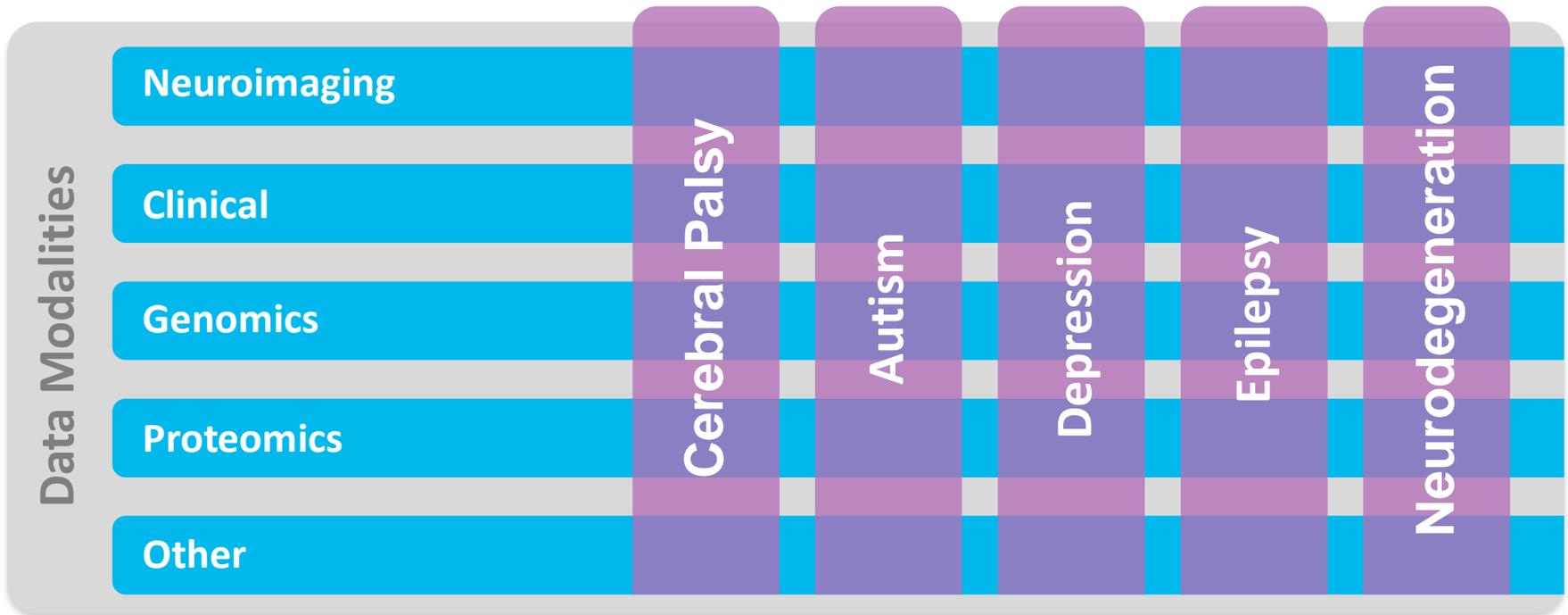
# Building a medical database

- Even with consent, is open access medical data a good idea?
  - ... or even ethical?
- Should we build a database under the assumption that the rules will remain as they are now?
- When OBI undertook the building of Brain-CODE these questions were foremost in all of our minds



# OBI's Big Data Needs

## *Integrated Discovery System*



- » Disease-themed, multidisciplinary collaborations collecting multiple modalities of data

Transformative, secure, virtual resource to advance understanding of CNS diseases

## What Brain-CODE means:

- To researchers
  - Collect, share and analyze 'big' data from anywhere, across multiple technologies, data sources, and diseases
  - Augment value of patient's research data through integration with accrued health information (e.g., ICES, OHS)
- To patients
  - Engagement in the research process
  - Education about their diseases, treatment options, and ongoing research
- To the economy
  - Accelerated development of novel treatments and diagnostics
  - Amplified government investment in CNS R&D

# Brain-CODE

- *Leverages existing infrastructure and expertise wherever possible*
- “InDOC” Consortium: □

Group	Lead	Expertise
<b>Ontario Cancer Biomarker Network (OCBN)</b>	Dr. Kenneth Evans (Exec Dir)	Molecular biomarker and clinical specimen data management; Overall leadership of BrainCODE
<b>Rotman Research Institute</b>	Dr. Stephen Strother	Brain imaging data management
<b>Applied Health Research Centre (AHRC) at St. Michael’s Hospital</b>	Dr. Muhammad Mamdani	Clinical trials infrastructure
<b>High Performance Computing Virtual Laboratory (HPCVL) in Kingston</b>	Dr. Ken Edgecombe	Secure, regulatory-compliant computing and large scale data storage
<b>E-Health Information Laboratory at the Children’s Hospital of Eastern Ontario (CHEO)</b>	Dr. Khaled El-Emam	Tools for de-identification, encryption, secure linkage, secure computation

# Privacy and Security

- 
- » The OBI and Brain-CODE *adhere to the highest levels of data privacy and security:*
    - **Security:** HPCVL, controlled access, security of data transfer and storage higher than bank-level
    - **Encryption:** Sophisticated algorithms for record linkage
    - **Privacy systems:** Full PIA completed in collaboration with Ontario Privacy Commissioner: designated “*Privacy by Design Ambassador*” by same
    - **Governance:** Extensive policy framework and governance apparatus ensures ongoing compliance with regulatory and REB requirements

# Brain-CODE Portal



## Brain-CODE

### Brain-CODE Tools



The **Brain-CODE Subject Registry** provides a secure way of entering personal health information from study participants. Use this tool to enter encrypted health card numbers and subject identifiers, which will allow your research data to be integrated across modalities and linked with other health databases.

[Login to Subject Registry](#)



**Medidata RAVE** is a commercial, regulatory-compliant clinical data management system for web-based electronic data capture of clinical assessment information collected from clinical trials, clinical research studies and registries.

[Enter Clinical Data](#)



**OpenClinica Enterprise** is a commercial, regulatory-compliant clinical data management system for web-based electronic data capture of clinical assessment information collected from clinical trials and clinical research studies.

[Enter Clinical Data](#)



**REDCap (Research Electronic Data Capture)** is a secure web application for building and managing online surveys and databases.

[Login to REDCap](#)



**SPReD** is a comprehensive database for the storage and management of neuroimaging data including MRI, PET, EEG, and CT.

[Upload Neuroimaging Data](#)



**BASE** is a database application for tracking genomics workflows, and manage and share resulting datasets. Use this tool to upload data generated from your microarray and next-generation sequencing studies.

[Upload Genomics Data](#)

#### New forum posts - test

- [Test Topic in CP-NET Forum](#)
- [Resetting my password](#)
- [Uploading neuroimaging data](#)

# *Designed to seamlessly link with autonomous databases*

Patient enrolls in study



Patient contributes data to other study



Mathematical properties of encrypted numbers allow them to be identified as coming from same person, without ever de-encrypting the health card number itself

# Encryption issues

- Some efforts to encrypt even highly complex data (e.g. Imaging, sequencing)
- Perhaps a solution to security issues, but what happens when you need to analyze the data?
- For many large data sets you have another issue: they're too large for widespread downloading to be practical
- One solution:
  - Allow access to data for mining, but within a secure environment
  - Only allow results to be extracted, not data
  - Brain-CODE approach, but also being proposed by others (e.g. NCI)

# Secure, Comprehensive, High Performance Data Management Platform

## Brain-CODE Platform

Access restricted to  
approved project  
personnel

Locked down  
database: no outside  
access

Data accessible only  
within workspaces;  
analytics tool support

ZONE 1  
**IDP Staging Areas**

ZONE 2  
**DataBank**

ZONE 3  
**Data Mining Workspaces**



# Brain-CODE Team

## Executive Management

Ken Evans  
Muhammad Mamdani  
Stephen Strother  
Ken Edgecombe  
Khaled El-Emam

## Project Management

Anthony Vaccarino  
Moyez Dharsee  
Barbara McQuaid

## High Performance Computing

Chris MacPhee  
Costa Dafnas  
Michael Hanlan

## Data Management

Rino La Grassa  
Natascha Kozlowski  
Anthony Vaccarino  
Stephen Strother  
Tom Gee  
Paulo Nuin

## Software Development

Tom Gee  
Nima Nourhaghighi  
Fan Dong  
Fletcher Johnson  
Paulo Nuin  
Rashad Badrawi  
Rino La Grassa  
Chris Ducharme  
Ben Eze  
Aleks Essex

## Common Data Elements

Anthony Vaccarino  
Stephen Strother  
Tom Gee  
Natascha Kozlowski  
Moyez Dharsee