# National Cybersecurity Center of Excellence

## NCCoE Virtual Workshop on Cybersecurity of Genomic Data

Wednesday, January 26, 2022, 11:00 AM – 4:30 PM (ET)

# AGENDA

| Segment | Time |
|---|---|
| Segment 1: Workshop Overview and Background | 11:00 AM – 11:40 AM |
| Segment 2: Keynotes | 11:40 AM – 12:20 PM |
| Segment 3: Challenges from the Field | 12:20 PM – 12:50 PM |
| Intermission | 12:50 PM – 1:30 PM |
| Segment 4: Challenges Sessions | 1:30 PM – 2:25 PM |
| Break | 2:25 PM – 2:35 PM |
| Segment 4 (Continued): Challenges Sessions | 2:35 PM – 3:45 PM |
| Break | 3:45 PM – 3:50 PM |
| Segment 5: Open Lightning Round | 3:50 PM – 4:20 PM |
| Segment 6: Next Steps | 4:20 PM – 4:30 PM |

# About the NCCoE

# WHO WE ARE

A **solution-driven**, **collaborative** hub addressing complex cybersecurity problems
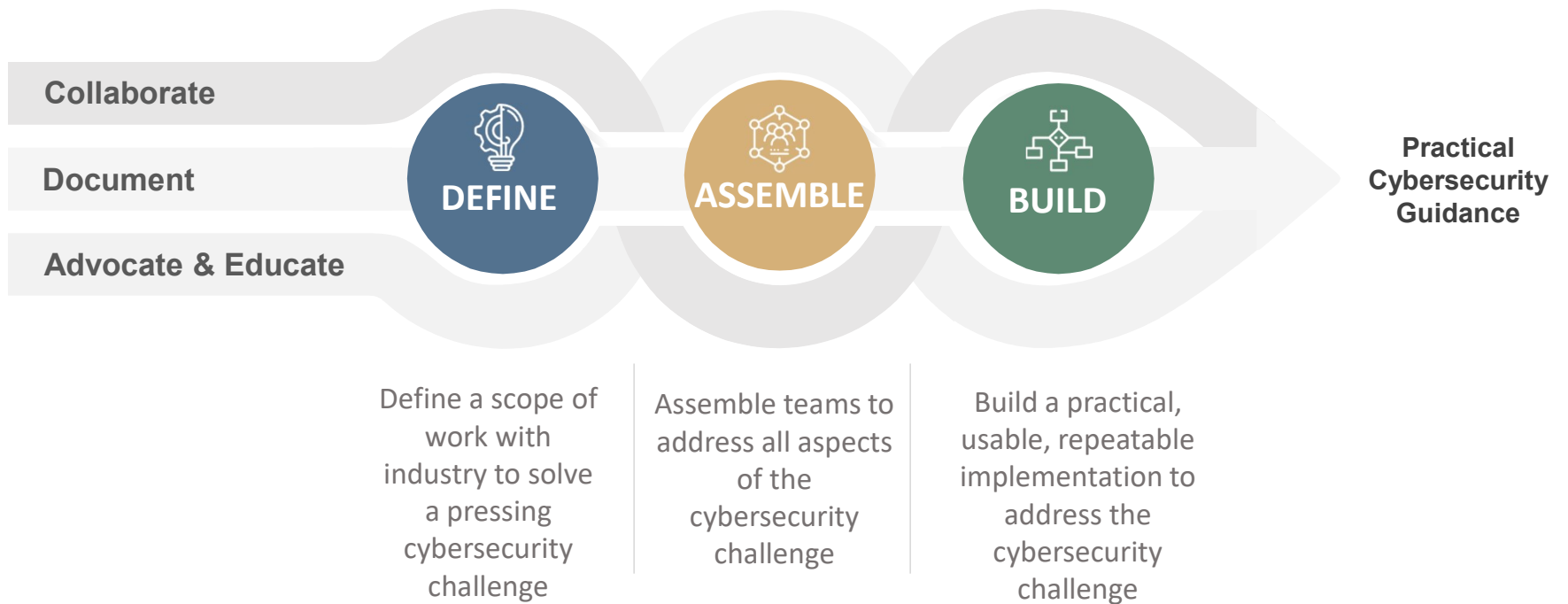
# OUR GOALS



**Improve cybersecurity** for businesses and commerce

**Lower the learning curve** for cybersecurity

**Spark innovation** in secure technology

# OUR APPROACH

| Collaborate | | DEFINE | ASSEMBLE | BUILD | Practical Cybersecurity Guidance |
| --- | --- | --- | --- | --- | --- |

**Collaborate**

**Document**

**Advocate & Educate**

**DEFINE**

**ASSEMBLE**

**BUILD**

**Practical Cybersecurity Guidance**

Define a scope of work with industry to solve a pressing cybersecurity challenge

Assemble teams to address all aspects of the cybersecurity challenge

Build a practical, usable, repeatable implementation to address the cybersecurity challenge

# Workshop Overview

Ron Pulivarti, NIST

# Housekeeping

- We support the health and well being for all.
  - We are supporting virtual collaboration.
  - We have three breaks planned for the day.

- We want audience engagement.
  - Please pose your questions for today's workshop using the Q&A window.
  - Please voice your insights in the Open Lightening Round from 3:50 – 4:30 PM.

- We intend to share our learnings today.
  - We are recording this session for future post on the NCCoE Website.
  - We will summarize key insights.

# Human Genomics at NIST

Samantha Maragh
Leader, Genome Editing Program

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

# The Human Genome

➢ The instruction code for humans

➢ ~6.4 billion letters long

➢ Present in each cell of a person

➢ ~ Half inherited from each biological parent

➢ Code is highly similar between people, but each person has a unique identifiable sequence

(Credit: Elymas/Shutterstock)

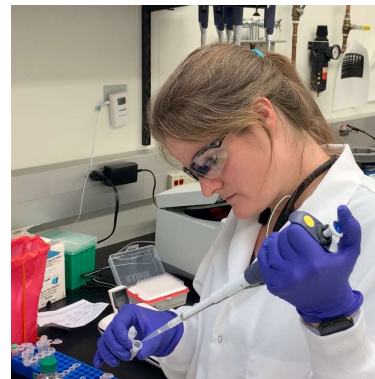# Uses of human genomic information
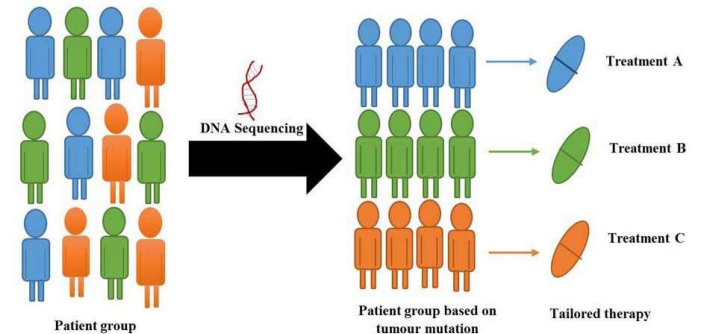


**Human Identification**

**Population Diversity and Ancestry**

**Human Health: Diagnostics**

**Scientific Research**

**Human Health: Treatment**

# NIST expertise with human genomics



Developing the standards for human identification

Human DNA Forensics

Leader, Pete Vallone

Leader, Justin Zook

Genome in a Bottle & Human Genomics

Whole human genome sequencing standard samples and datasets

Standards for cancer biomarker detection used for clinical diagnostics

Cancer Biomarkers

Leader, Hua-Jun He

Leader, Samantha Maragh

Genome Editing

Whole human genome samples from engineered human cells, standard datasets, and terminology standards

# NIST formed the Genome in a Bottle Consortium in 2012

GIAB has characterized variants in 7 human genomes
and released NIST whole genome DNA standards

# CRISPR technologies and uses

**CRISPR-Cas**
CRISPR - **C**lustered **R**egularly **I**nterspaced **S**hort **P**alindromic **R**epeats
Cas - **C**RISPR **as**sociated protein

**CRISPR-Cas** system were identified in nature as bacterial immune systems and have been pivoted to enable modification of the genetic code within cells at designed target positions (**genome editing**)



PAM

sgRNA

Cas9

+

Genome Editing System (e.g., CRISR-Cas) designed to target genome sequence

Cells of interest

DNA sequence change

Engineered/edited functional cells

**Human Gene Therapy applications**

Patient

Cells

Engineer cells

Engineered Cell therapy

**Autologous**

Donor

**Allogenetic** Patients

Patient

Direct delivery of genome engineering molecules

*in vivo*

# NIST Genome Editing Consortium
## (launched October 2018)

**NIST**

## MISSION

Convene experts across academia, industry, non-profit & government to addresses the measurements and standards needed to increase confidence of utilizing genome editing technologies in research and commercial products

### ORGANIZATION



WG1 Specificity Measurements — Physical benchmarks — Assay qualification

WG2 Data & Metadata — Community Norms — Standard datasets

WG3 Lexicon — Harmonized genome editing terminology

### MEMBER BENEFITS

- Access to a neutral forum for addressing pre-competitive needs

- Participation in the development of experimental benchmarks, guidelines and terminology

- Access to tools developed by the consortium ahead of public release

## MEMBERS

- Agilent
- Aldevron
- Applied StemCell
- AstraZeneca
- Bionano Genomics
- Bio-Rad
- Bluebird bio*
- Caribou Biosciences
- Catalytic Data Science
- Cergentis
- COBO Technologies
- **College of American Pathologists (CAP)**
- CRISPR Therapeutics
- DARPA
- DowDuPont Agroscience (Corteva)*
- Editas Medicine
- EMBL-EBI
- **FDA CBER**
- Genomic Vision
- Horizon Discovery
- Illumina
- Inscripta
- Integrated DNA Technologies
- Intellia Therapeutics

- KromaTiD
- Lonza
- Macrogen
- Mass General Hospital
- Mission Bio
- Novartis
- New England Biolabs
- NIH/NINDS
- NIH SCGE
- Precision Biosciences

- Sangamo Therapeutics
- SeQure Dx
- St. Jude Children's Research Hospital
- Synthego
- ThermoFisher Scientific
- Twinstrand Bioscineces
- UCSC
- WhiteLab Genomics

NIST coordinates with FDA and Center for Veterinary Medicine (CVM)

Cost sharing model. All members contribute $20,000 annually or *in-kind*

* Former members

# Importance of securing human genomics data

NIST

✓ Data integrity

✓ Data reliability

✓ Maintain limited access to individual genomic information

✓ Protect knowledge / intellectual property

✓ Data reusability / prevent loss

✓ Prevent against nefarious use or misuse

✓ Privacy

# Contact Us

## Thank You!

**Contact**

Samantha Maragh
Leader, Genome Editing Program

**Email address**

samantha@nist.gov

# Privacy at NIST

Naomi Lefkovitz

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

# Relationship Between Cybersecurity and Privacy Risk



**Cybersecurity Risks**

associated with cybersecurity incidents arising from loss of confidentiality, integrity, or availability

cyber security-related privacy events

**Privacy Risks**

associated with privacy events arising from data processing

**Data:** A representation of information, including digital and non-digital formats

**Privacy Event:** The occurrence or potential occurrence of problematic data actions

**Data Processing:** The collective set of data actions (i.e., the complete data life cycle, including, but not limited to collection, retention, logging, generation, transformation, use, disclosure, sharing, transmission, and disposal)

**Privacy Risk:** The likelihood that individuals will experience problems resulting from data processing, and the impact should they occur

# NIST Privacy Risk Assessment Methodology (PRAM)



Catalog of Problematic Data Actions and Problems

# NIST Privacy Engineering Objectives

**Predictability: enabling reliable assumptions by individuals, owners, and operators about data and their processing by a system, product, or service.**

**Manageability: providing the capability for granular administration of data, including alteration, deletion, and selective disclosure.**

**Disassociability: enabling the processing of data or events without association to individuals or devices beyond the operational requirements of the system.**

# Resources

**Websites**

https://www.nist.gov/privacyframework

**Mailing List**

List.nist.gov/privacyframework

**Contact Us**

PrivacyFramework@nist.gov

# National Institute of Standards and Technology - Cybersecurity

Ron Pulivarti, Senior Cybersecurity Engineer for the Healthcare Sector at the National Cybersecurity Center of Excellence (NCCoE), which is part of NIST

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

**NCCoE** — NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

Celebrating 50 years of Cybersecurity at NIST

# NINE PRIORITY AREAS:

Enhancing Risk Management
Trustworthy Networks
Strengthening Cryptographic Standards & Validation
Securing Emerging Technologies
Privacy
Trustworthy Platforms
Metrics & Measurement
Identity & Access Management
Awareness, Training, Education & Workforce Development

**Cybersecurity Framework released; Cybersecurity Enhancement Act assigns NIST workforce, other responsibilities**

**Law designates NIST to Federal Acquisition Security Council, produce supply chain guidance**

**NIST produces IoT guidance per IoT Cybersecurity Improvement Act; NIST issues Security and Privacy Controls, Rev 5; NIST updates NICE strategic plan**

## 2014    2016    2018    2019    2020    2021

**NIST launches public key Post-Quantum Cryptography Standardization initiative**

**NIST launches Small Business Cybersecurity Corner website following 2018 statute**

**NIST launches effort to enhance software supply chain security in response to Executive Order and technology supply chain partnership**

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

**NCCoE**
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# NIST Cybersecurity Framework 1.1



| Function | Category |
|---|---|
| Identify | Asset Management |
| | Business Environment |
| | Governance |
| | Risk Assessment |
| | Risk Management Strategy |
| | Supply Chain Risk Management |
| Protect | Identity Management and Access Control |
| | Awareness and Training |
| | Data Security |
| | Information Protection Processes & Procedures |
| | Maintenance |
| | Protective Technology |
| Detect | Anomalies and Events |
| | Security Continuous Monitoring |
| | Detection Processes |
| Respond | Response Planning |
| | Communications |
| | Analysis |
| | Mitigation |
| | Improvements |
| Recover | Recovery Planning |
| | Improvements |
| | Communications |

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

**NCCoE** NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

**NIST Special Publication 800-160, Volume 2**
Revision 1

# Developing Cyber-Resilient Systems:

*A Systems Security Engineering Approach*

| Adaptive Response | Analytic Monitoring | Coordinated Protection | Contextual Awareness | Deception | Diversity | Dynamic Positioning |
|---|---|---|---|---|---|---|
| Dynamic Reconfiguration | Monitoring and Damage Assessment | Calibrated Defense-in-Depth | Dynamic Resource Awareness | Obfuscation | Architectural Diversity | Functional Relocation of Sensors |
| Dynamic Resource Allocation | Sensor Fusion and Analysis | Consistency Analysis | Dynamic Threat Awareness | Disinformation | Design Diversity | Functional Relocation of Cyber Resources |
| Adaptive Management | Forensic and Behavioral Analysis | Orchestration | Mission Dependency and Status Visualization | Misdirection | Synthetic Diversity | Asset Mobility |
| | | Self-Challenge | | Tainting | Information Diversity | Fragmentation |
| | | | | | Path Diversity | Distributed Functionality |
| | | | | | Supply Chain Diversity | |

| Non-Persistence | Privilege Restriction | Realignment | Redundancy | Segmentation | Substantiated Integrity | Unpredictability |
|---|---|---|---|---|---|---|
| Non-Persistent Information | Trust-Based Privilege Management | Purposing | Protected Backup and Restore | Predefined Segmentation | Integrity Checks | Temporal Unpredictability |
| Non-Persistent Services | Attribute-Based Usage Restriction | Offloading | Surplus Capacity | Dynamic Segmentation and Isolation | Provenance Tracking | Contextual Unpredictability |
| Non-Persistent Connectivity | Dynamic Privileges | Restriction | Replication | | Behavioral Validation | |
| | | Replacement | | | | |
| | | Specialization | | | | |
| | | Evolvability | | | | |

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

**NCCoE**
NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

# Thank You.

# NIST Experiences in Genomics, Cybersecurity, and Privacy

Moderated Questions and Answers

In the toolbar at the bottom, click on the 3-dot button

On the menu, click Q&A

Q&A

Copy Event Link

Audio Connection

What color is the sky?

Send    Send Privately...

Enter your question in the Q&A panel.

1. On the right side, click on Q&A header to open the Q&A panel.
2. Type in the box **your name, organization and question**.
3. Click send.

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

NCCoE NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

# National Counterintelligence and Security Center



**Mission**

- Lead and support the U.S. Government's counterintelligence and security activities critical to protecting our nation.

- Provide counterintelligence outreach to U.S. private sector entities at risk of foreign intelligence penetration.

- Issue public warnings regarding intelligence threats to the United States.

# Global Threat Picture

**Expanding Array of Adversaries**

**Improving Capabilities & Tradecraft**

**Expanded Range of Targets & Operations**

## Emerging Technologies: A Key Focus of Strategic Competitors

- U.S. leadership in emerging technology sectors -- such as biotech, AI, & quantum -- faces growing challenges from strategic competitors.

- China, Russia, and other nations recognize the economic & military benefits of these technologies & have enacted comprehensive national strategies to achieve leadership in these areas.

- To achieve their strategic goals, strategic competitors are using a wide variety of legal, quasi-legal, and illegal methods to acquire technology, talent, and know-how from the U.S. and other nations.



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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Foreign Exploitation of Genomic Data is Already Occurring

## Exploitation of DNA for Societal Control & Repression by the People's Republic of China (PRC)

- The PRC has conducted large-scale collection of DNA and other biometric data from residents of Xinjiang ages 12 to 65.

- DNA samples, fingerprints, iris scans, and blood types are linked to ID numbers and centralized in searchable database used by PRC authorities to carry out surveillance and detentions.

- Since 2017, between 1 million and 1.8 million Uyghurs & other minorities in Xinjiang have been placed in "re-education" centers.

- Multiple Chinese entities & companies, including two subsidiaries of BGI (the world's largest genomics company based in China), have been sanctioned by the U.S. Government for their roles in the PRC repression of Uyghurs in Xinjiang.

Bottom photo source: https://baijiahao.baidu.com/s?id=1564669932542581

# PRC Ambitions and Genomic Data

**NATIONAL POLICIES:** PRC has enacted national policies prioritizing the collection of healthcare data, including genetic data, both at home and abroad to achieve its goal of becoming a global biotech leader.

- **Precision Medicine Initiative:** In 2016, the PRC announced a $9 billion, 15-year project to collect, analyze, and sequence genomic data to become global leader in precision medicine under the "Healthy China 2030" initiative.

- **14th Five Year Plan:** In 2021, the PRC unveiled its 14th Five Year Plan, which listed genetics and biotechnology as among the cutting-edge science and technology research areas the PRC seeks to dominate in the years 2021-2025.

- **China Standards 2035:** The PRC's national strategy to set global rules and standards in emerging technologies, including those critical for future precision healthcare.

**ECONOMIC ADVANTAGE:** The PRC understands the collection and analysis of large genomic data sets from diverse populations helps foster new medical discoveries that can advance its AI, pharmaceutical, and precision medicine industries.

**MILITARY / SECURITY ADVANTAGE:** PRC has used genetic analysis for state surveillance, societal control, and has been conducting genetic research for military purposes and biodefense.

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# PRC Vectors to Access U.S. Genomic Data

**INVESTMENTS**
- China's largest genomics company, BGI, purchased U.S. genomic sequencing firm Complete Genomics in 2013.
- In 2015 WuXi Pharma Tech acquired U.S. genetic sequencing company NextCODE. WuXi NextCODE later received accreditation to perform molecular diagnostic and genetic testing in the U.S.

**PARTNERSHIPS**
- China's BGI has partnered with health institutions across America to provide low-cost genomic testing and sequencing services, while also gaining access to genetic data on persons in the U.S.
- According to a 2019 report prepared by Gryphon Scientific, 23 companies associated with China are certified to perform genetic testing in the U.S., giving them access to genetic data on patients in the US.

**COMPELLED ACCESS**
- All Chinese companies are subject to PRC laws requiring them to share data they acquire with the PRC government.

**CYBER INTRUSIONS**
- PRC has conducted cyber attacks on U.S. healthcare institutions and companies (such as Anthem and others) to acquire personal health information.

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Risks Associated with PRC Access to U.S. Genomic Data

**PRIVACY:** Your genetic data could end up in the hands of the PRC and used for purposes you never intended. The loss of your DNA is permanent and not only affects you, but your relatives, and potentially, generations to come.

**INTELLIGENCE:** By combining genetic data with other PII and healthcare / lifestyle data the PRC has acquired through cyberattacks and other means, the PRC could use this information to target U.S. personnel, dissidents, journalists, and others around the world for potential surveillance, manipulation, or extortion.

**ECONOMIC COMPETITION:** Large, diverse sets of genetic and health data from around the world can help the PRC enhance its AI, pharmaceutical, healthcare, and precision medicine industries at the expense of U.S. biotech industry.

- **No Reciprocity:** The PRC severely restricts U.S. and other foreign access to Chinese genetic data, putting America's biotech industry at a disadvantage.

**MILITARY & SECURITY CAPABILITIES:** There is growing concern over PRC research and exploitation of genetic data for bioweapons and biodefense, including to enhance the performance of soldiers in combat and more effectively support force readiness.

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

NCCoE
NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

# Direct-to-Consumer (DTC) Genetic Testing Companies



**DESIRABLE TARGETS:** DTC genetic testing companies, information exchanges, and data libraries are desirable targets for foreign adversaries, cyber criminals, and insider threats.

**HUGE GENETIC HOLDINGS:** DTC genetic testing companies hold large quantities of human genetic data and other personal information. Last year, the American Medical Association (AMA) projected that as many as 100 million individuals would undergo DTC genetic tests by the end of 2021.

**LESS REGULATED THAN U.S. HEALTHCARE PROVIDERS:** Data held by DTC genetic testing companies are not subject to HIPAA / privacy and security requirements that apply to health care providers, as consumers send samples directly to the companies without the involvement of a health care provider.

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Cyber Risks and DTC Genetic Testing Companies



CPO
MAGAZINE

HOME    NEWS    INSIGHTS    RESOURCES

CYBER SECURITY    NEWS    · 2 MIN READ

## DNA Testing Firm Data Breach Exposed Sensitive Information of More Than 2.1 Million People

ALICIA HOPE · DECEMBER 9, 2021

DNA Diagnostics Center (DDC) filed a data breach notification with the Maine Attorney General's office disclosing that hackers accessed sensitive details of more than 2.1 million people.

Image source: DNA Testing Firm Data Breach Exposed Sensitive Information of More Than 2.1 Million People - CPO Magazine

In November 2021, an Ohio-based DNA testing company reported to regulators that personal information on more than **2.1 million people** was acquired in a hacking incident. No genetic data reported stolen.

In July 2019, a California-based DNA testing company accidently exposed the personal data of 3,000 customers online, including some **300 files containing genetic data**.

In June 2018, an Israel-based DNA testing company, announced it had been breached and the email addresses of more than **92 million users** were compromised.  No genetic data reported stolen.

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Key Takeaways

**GREAT PROMISE, BUT KEY RISKS:** The collection and analysis of genomic data holds great promise for medical breakthroughs, but with it comes important risks to privacy as well as economic and national security.

- Large genetic databases that allow people's ancestry to be revealed and crimes to be solved also can be misused for surveillance and societal repression.
- Genomic technology used to design disease therapies tailored to an individual also can be used to identify genetic vulnerabilities in a population that potentially could be targeted.

**ADVERSARIES ALREADY EXPLOITING GENOMIC DATA:** Adversaries are already exploiting genomic data and have national plans to acquire and harness genomic data at home and abroad for their economic advantage and national security.

- Foreign companies and authoritarian regimes have already gained significant access to U.S. genomic data and related healthcare data through investments, research partnerships, contractual agreements, and other means.

**LEGAL / REGULATORY GAPS ON GENETIC DATA:** U.S. laws currently do not treat genetic data as a national security asset, but primarily focus on privacy and IP protection. Few restrictions prevent a U.S. company from selling genetic data to parties outside the U.S.

# Keynote:
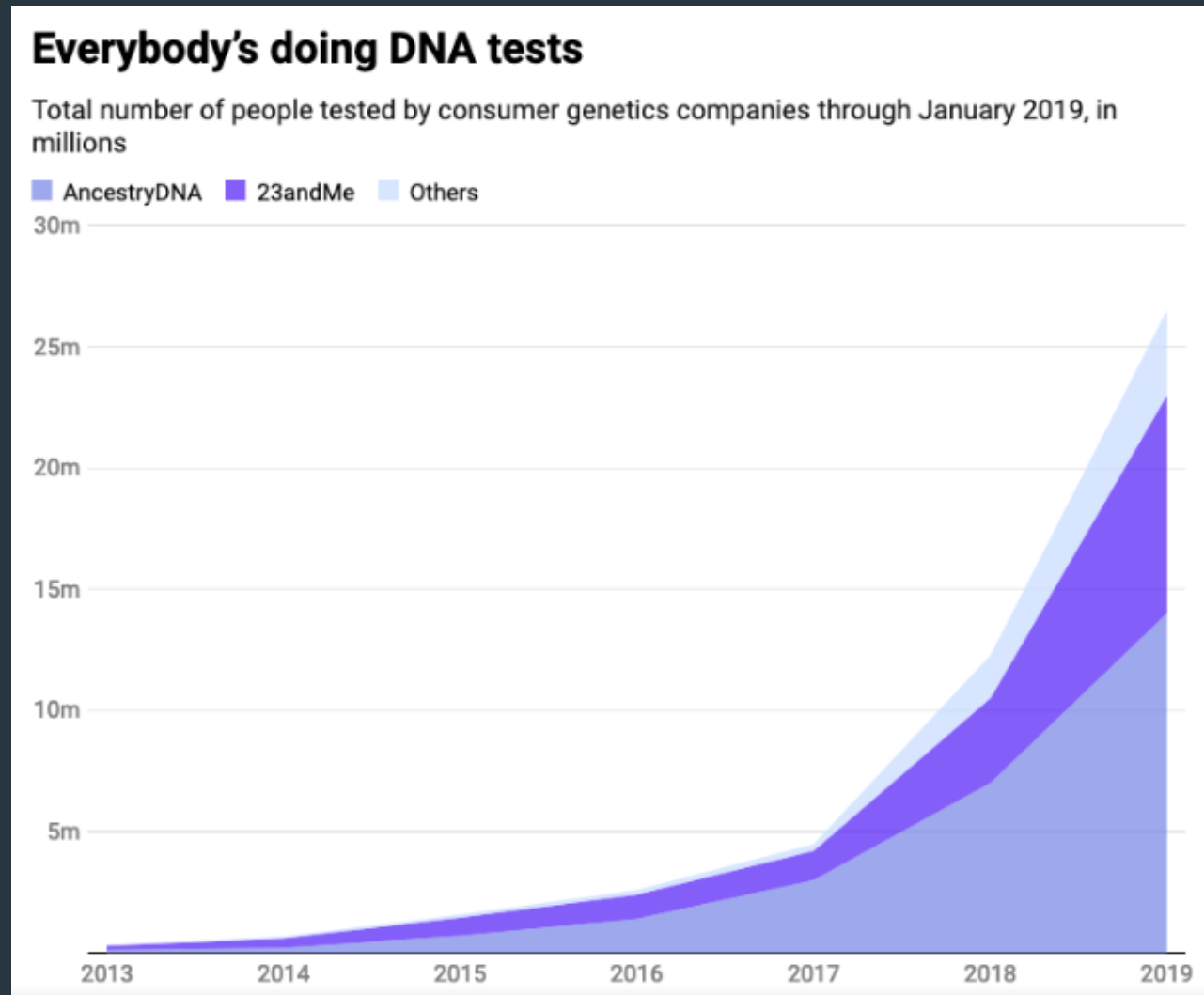# The Enabling Perspective

Yaniv Erlich (Eleven Therapeutics)

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Free-for-all genetic surveillance nation

- Dr. Yaniv Erlich
@erlichya

# The advent of consumer genomics



*MIT Technology review*

# Relative matching via shared IBD



Identity by –descent (IBD) segments

Modified from Huff et al., Genome Research, 2011

# Relative matching is the core of genetic genealogy

# 3ʳᵈ party support for relative matching

## Your raw genetic data

```
# MyHeritage DNA raw data.
# This file was generated on 2018-10-10 09:03:32
# For each SNP, we provide the identifier, chromosome
# number, base pair position and genotype. The genotype
# is reported on the forward (+) strand with respect to
# the human reference build 37.
# THIS INFORMATION IS FOR YOUR PERSONAL USE AND IS
# INTENDED FOR GENEALOGICAL RESEARCH
# ONLY. IT IS NOT INTENDED FOR MEDICAL OR HEALTH
# PURPOSES. PLEASE BE AWARE THAT THE
# DOWNLOADED DATA WILL NO LONGER BE PROTECTED BY OUR
SECURITY MEASURES.


#RSID,CHROMOSOME,POSITION,RESULT
"rs4477212","1","82154","AA"
"rs3094315","1","752566","--"
"rs3131972","1","752721","AG"
"rs12562034","1","768448","--"
"rs12124819","1","776546","--"
"rs11240777","1","798959","GG"
"rs6681049","1","800007","--"
"rs4970383","1","838555","AC"
"rs4475691","1","846808","TC"
"rs7537756","1","854250","AG"
"rs13302982","1","861808","GG"
"rs1110052","1","873558","TG"
"rs2272756","1","882033","GG"
```

Upload →

## MyHeritage (users: 3M)

## FTDNA (users: 1M)

## GEDmatch (users: 1.4M)

## DNA.Land (users: 150K)

# 3ʳᵈ party uploads are highly important

# Using genetic genealogy for forensic is not a new idea



**nature REVIEWS** GENETICS

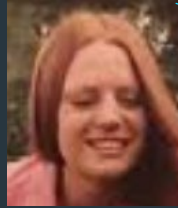## Routes for breaching and protecting genetic privacy

**REVIEWS**

*Yaniv Erlich[1] and Arvind Narayanan[2]*

Abstract | We are entering an era of ubiquitous genetic information for research, clinical care and personal curiosity. Sharing these data sets is vital for progress in biomedical research. However, a growing concern is the ability to protect the genetic privacy of the data originators. Here, we present an overview of genetic privacy breaching strategies. We outline the principles of each technique, indicate the underlying assumptions, and assess their technological complexity and maturation. We then review potential mitigation methods for privacy-preserving dissemination of sensitive data and highlight different cases that are relevant to genetic applications.

## Erlich and Narayanan, 2014

An open research question is the use of non-Y-chromosome markers for genealogical triangulation. The Mitosearch and GEDmatch websites run open, searchable databases for matching mitochondrial and autosomal genotypes, respectively. Our expectation is that mitochondrial data will not be very informative for tracing identities. The resolution of mitochondrial searches is low owing to the small size of the mitochondrial genome, which means that a large number of individuals share the same mitochondrial haplotypes. In addition, matrilineal identifiers (such as surname or clan) are fairly rare in most human societies, which complicates the use of mitochondrial haplotype for identity tracing. By contrast, autosomal searches can be powerful. Genetic genealogy companies have started to market services for dense genome-wide arrays that enable the identification of distant relatives (on the order of third to fourth cousins) with fairly sufficient accuracy[43]. These hits would reduce the search space to no more than a few thousand individuals[44]. The main challenge of this approach would be to derive a list of potential people from a genealogical match. As stated above, family trees of most individuals are not publicly available; such searches are therefore demanding and would require indexing a large number of genealogical websites. With the growing interest in genealogy, this technique might be easier in the future and should be taken into consideration.

# Long range familial searches



| Case | Announcement | By |
|---|---|---|
| Buckskin Girl | April 9, 2018 | DNA Doe Project |
| Golden State Killer | April 24, 2018 | Barbara Rae-Venter |
| Lyle Stevik | May 8, 2018 | DNA Doe Project |
| William Earl Talbott II | May 21, 2018 | Parabon |
| Joseph Newton Chandler III | June 21 2018 | DNA Doe Project |
| Gary Hartman | June 22, 2018 | Parabon |
| Raymond "DJ Freez" Rowe | June 25, 2018 | Parabon |
| James Otto Earhart | June 26, 2018 | Parabon |
| John D. Miller | July 15, 2018 | Parabon |
| Matthew Dusseault | July 28, 2018 | Parabon |
| Spencer Glen Monnett | July 29, 2018 | Parabon |
| Darold Wayne Bowden | August 23rd, 2018 | Parabon |
| Michael F. Henslick | August 29th, 2018 | Parabon |

# The probability of finding a relative?

*Repeat 1,280,000 times*

Find genetic relatives for a MyHeritage participant
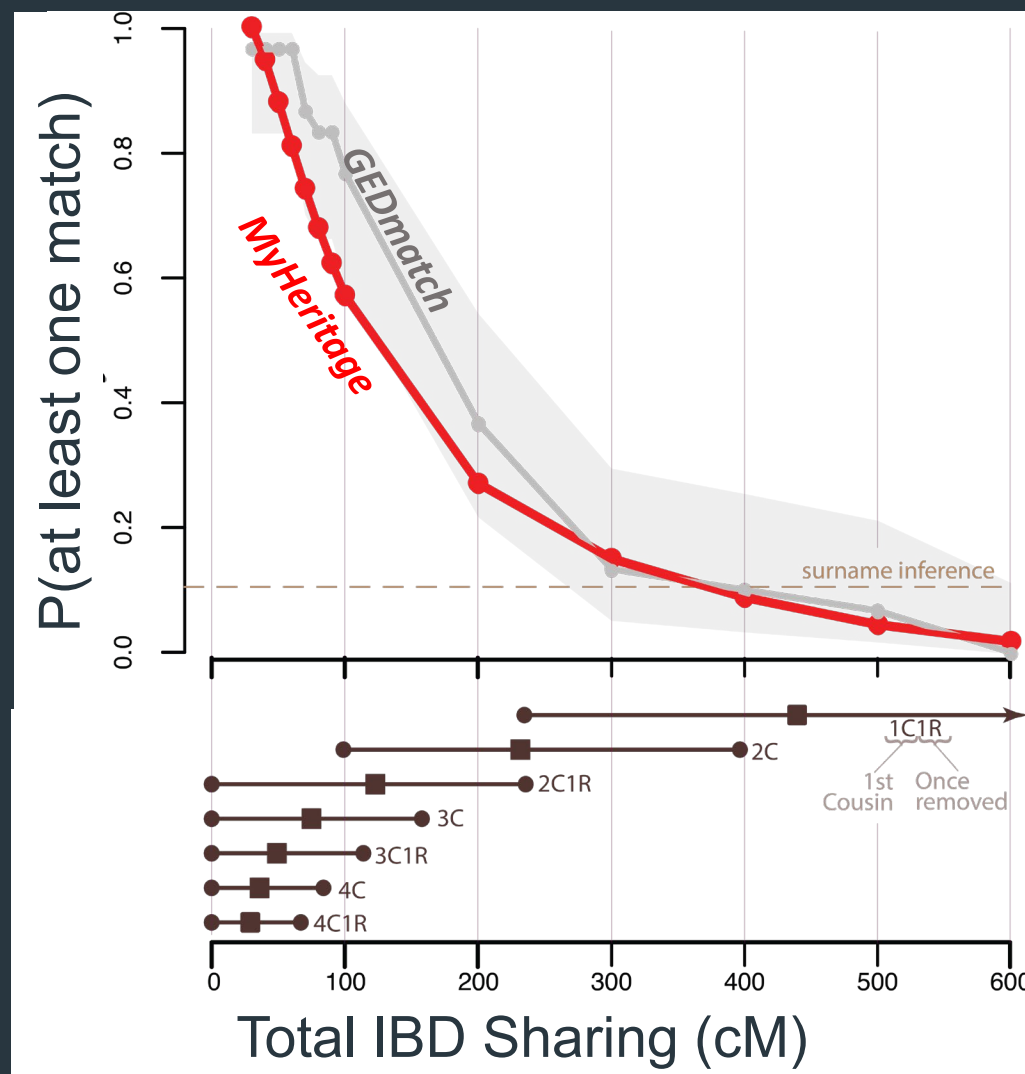
↓

Exclude >700cM matches

↓

What is the top match of the person?

*Repeat 30 times*

Find genetic relatives for a GEDMatch participant

↓

Exclude >700cM matches

↓

What is the top match of the person?



**Estimate: ~60% of US individuals of European heritage have a 3rd cousin match**

# Small scale study confirms our projection



SCIENCE

## We Tried To Find 10 BuzzFeed Employees Just Like Cops Did For The Golden State Killer

The Golden State Killer case has triggered a boom in "genetic genealogy" for solving crimes. But how hard is it to find people by sleuthing in their family trees?
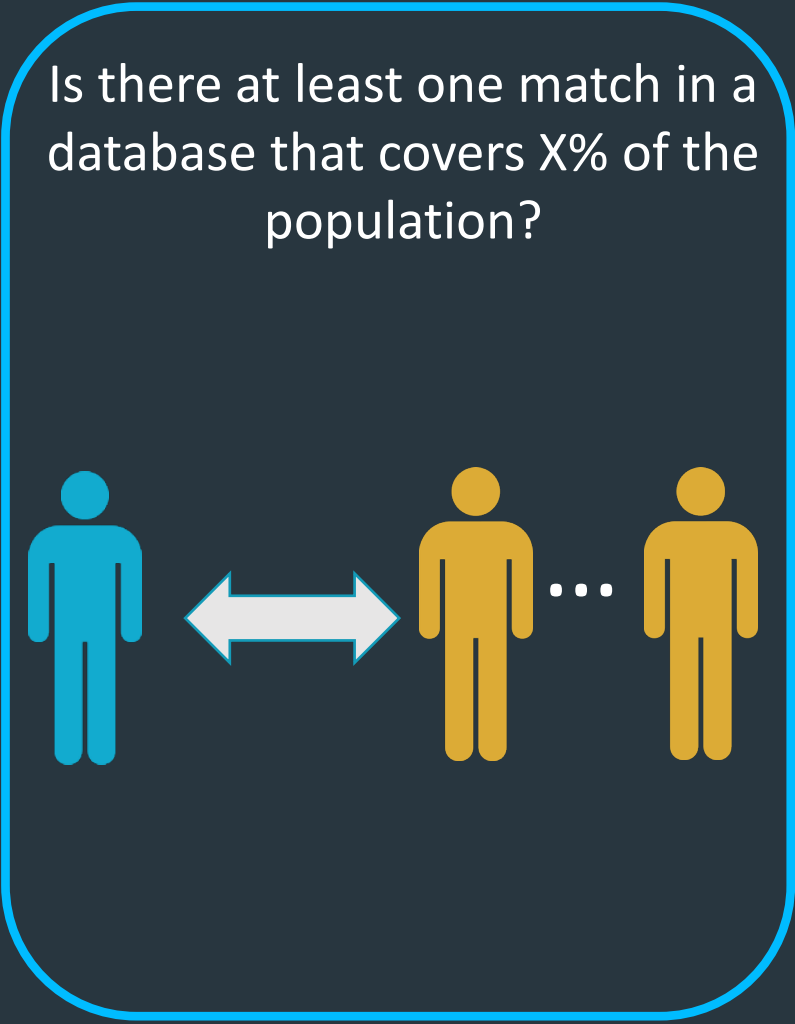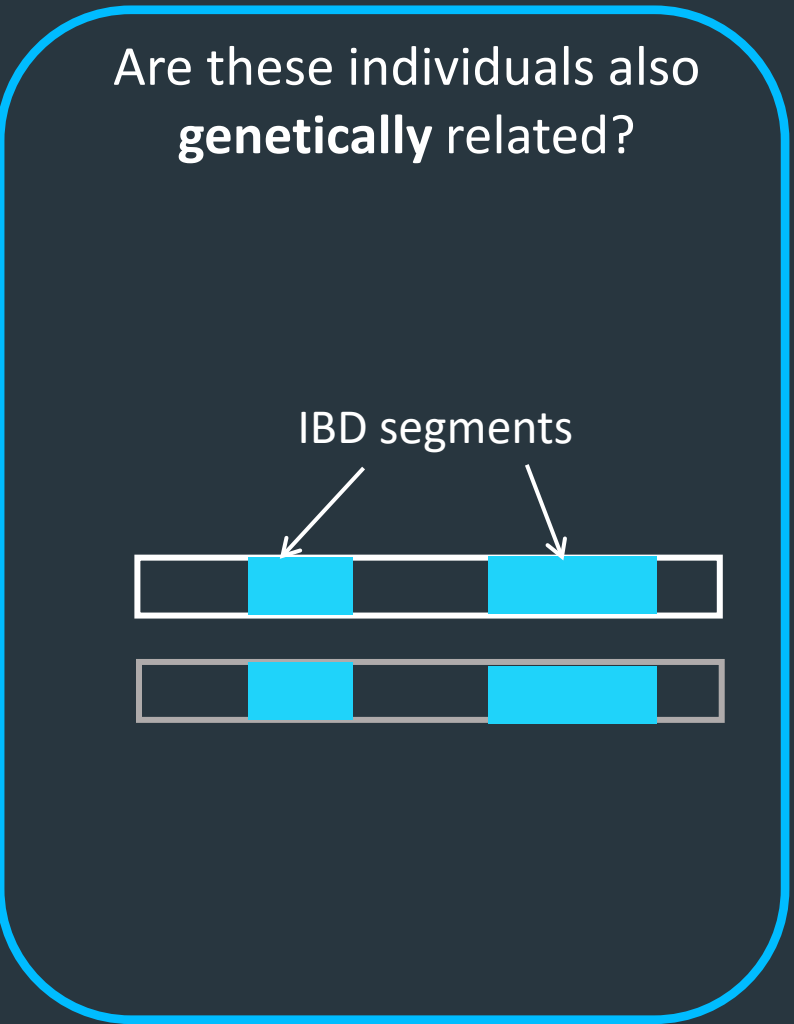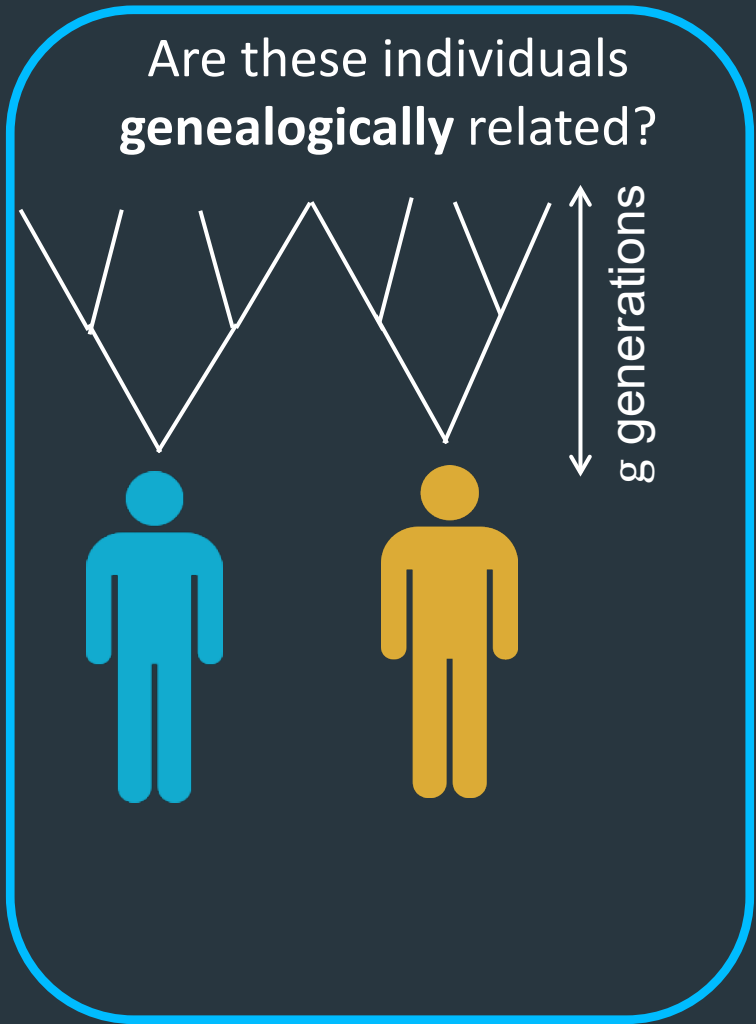
Peter Aldhous
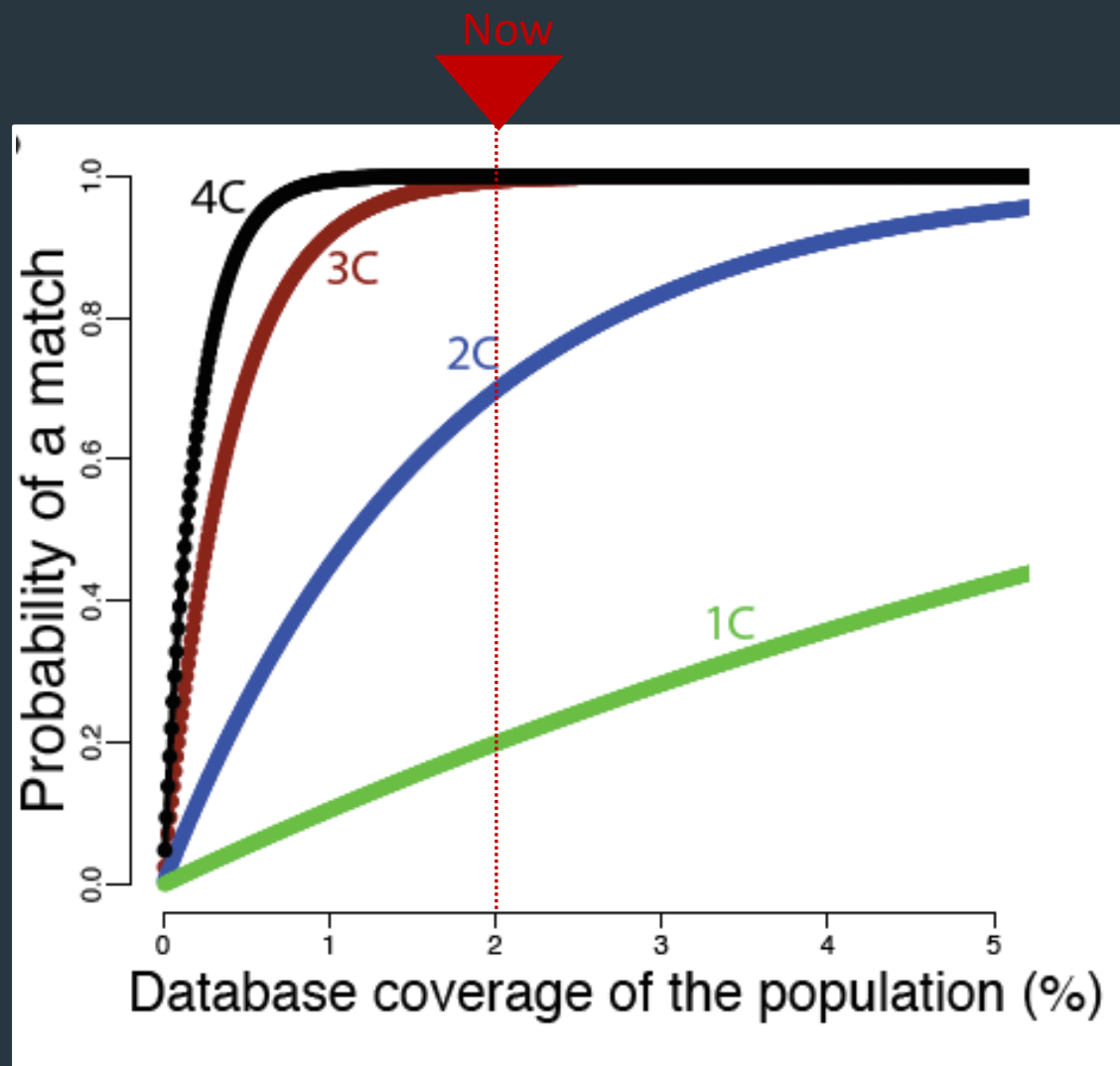BuzzFeed News Reporter

Posted on April 9, 2019, at 9:16 a.m. ET

In the end, I identified 6 out of our 10 volunteers. Four of those cases I solved by tracking them down through their relatives' family trees, much as the cops did with DeAngelo. In a twist I didn't anticipate, I found two more not through their relatives, but simply because their ancestry indicated that their family came from a specific country — raising uncomfortable questions about genetic racial profiling.

# Modeling the probability of finding relatives

Are these individuals **genealogically** related?

g generations

Are these individuals also **genetically** related?

IBD segments

Is there at least one match in a database that covers X% of the population?

...

Caveats: no-population structure or consanguinity; assumes random samples.

# The probability of a match in the future



**Virtually any US person of European heritage will have a 3ʳᵈ cousin in these databases.**

# Can we get to a single person?

325M → | **3rd cousin match** | ... ex → 16.5



| | | | | G3-grandparent | | | |
|---|---|---|---|---|---|---|---|
| | | | G2-grandparent | G3 A/U | | | |
| | | Great-Grandparent | G2 A/U | 1C3R | | | |
| | Grandparent 1360-2136 | Great A/U 487-1770 | 30x 1C2R 96-385 | 150x 2C2R 0-158 | | | |
| Parent 3380-37.. | Parent 3380-3718 | Aunt/Uncle 1452-2034 | 1C1R 236-657 | 75x 2C1R 96-385 | 3C1R 0-112 | | |
| Sibling 2342-2917 | Match | | 1C 638-1119 | 38x 2C 102-406 | 188x 3C 0-155 | 4C 0-82 | |
| Nie/Nep 1452-2034 | Child 3380-3718 | | 1C1R 236-657 | 94x 2C1R 36-235 | 3C1R 0-112 | 4C1R 0-65 | |
| G. Nie/Nep 487-1770 | Grandchild 1360-2136 | | 17x 1C2R 96-385 | 234x 2C2R 0-157 | 3C2R 0-90 | 4C2R 0-55 | |

**Even three simple pieces of**

# Crime Scene and Distance Correlates of Serial Rape

Janet Warren,[1,7] Roland Reboussin,[2] Robert R. Hazelwood,[3] Andrea Cummings,[4] Natalie Gibbs,[5] and Susan Trumbetta[6]

This study, derived from a sample of 108 serial rapists (rapes = 565), examines the relationship between demographic, crime scene, and criminal history variables and the distance traveled by serial rapists in order to offend. The pattern of offenses perpetrated by each of the 108 serial offenders as it relates to his place of residence is also analyzed in terms of known characteristics of the offender and his offenses. The theoretical focus of the study integrates premises derived from criminal investigative analysis, environmental criminology, ethnographic geography, journey to crime research, and criminal geographic targeting to explore the cognitive symmetry between the "how" and the "where" of serial sexual offenses. These components or dimensions of serial crime are explored in an attempt to aid law enforcement in their investigation of hard-to-solve serial crimes.

**KEY WORDS:** serial rape; journey to crime; crime scene analysis; criminal investigative analysis; spatial analysis of crime; environmental criminology; criminal geographic targeting; geographic profiling.

# Summary so far

We expect 2$^{nd}$ - 3$^{rd}$ cousin for virtually every person in the US with European descent (if access is allowed)

Basic demographic information can substantially narrow the search space to handful of individuals

The method is extremely powerful

# Paper

# Identity inference of genomic data using long-range familial searches

**Yaniv Erlich[1,2,3,4]\*, Tal Shor[1], Itsik Pe'er[2,3], Shai Carmi[5]**

[1]MyHeritage, Or Yehuda 6037606, Israel. [2]Department of Computer Science, Fu Foundation School of Engineering, Columbia University, New York, NY, USA. [3]Center for Computational Biology and Bioinformatics (C2B2), Department of Systems Biology, Columbia University, New York, NY, USA. [4]New York Genome Center, New York, NY, USA. [5]Braun School of Public Health and Community Medicine, The Hebrew University of Jerusalem, Jerusalem, Israel.

\*Corresponding author. Email: erlichya@gmail.com

Consumer genomics databases have reached the scale of millions of individuals. Recently, law enforcement authorities have exploited some of these databases to identify suspects via distant familial relatives. Using genomic data of 1.28 million individuals tested with consumer genomics, we investigated the power of this technique. We project that over 60% of the searches for individuals of European-descent will result in a third cousin or closer match, which can allow their identification using demographic identifiers. Moreover, the technique could implicate nearly any US-individual of European-descent in the near future. We demonstrate that the technique can also identify research participants of a public sequencing project. Based on these results, we propose a potential mitigation strategy and policy implications to human subject research.

# So why am I worried?

# Genetic genealogy can be weaponized by counter-intelligence

1. Everyone can uploads data to GEDmatch/FTDNA/etc...

2. Also adversaries of the US (they don't give a damn to Toc)

3. Counter-intelligence and other players can exploit genetic genealogy to cast population-scale genetic surveillance over the US

4. The risk is asymmetric (US is substantially affected but not other countries)

# Intelligence services are interested in DNA

## WikiLeaks: are Chinese spies stealing Iceland's genetic database?

by Jared Yee | 18 Dec 2010 | Link

Another bioethics angle has emerged in Wikileaks. Chinese spies are investigating genetic research companies in Iceland, according to cables written in authorities said that intelligence gathering included bugging phone lin hacking into databases.

## Novichok bottle could hold attackers' DNA



Dawn Sturgess died last week from novichok poisoning; Charlie Rowley remains in hospital

AFP PHOTO/FACEBOOK

# Differentiate good vs. bad actors

Can we differentiate legitimate searches from illegitimate searches?

- Legitimate datasets are produced with a regular DTC lab or authorized crime labs

- Illegitimate datasets are produced by research labs, unauthorized crime labs, etc.

- Idea: ask authorized labs to sign datasets before letting users downloading the data.

# How it works?

**Before**                                                                 **After**

```
# MyHeritage DNA raw data.
# This file was generated on 2018-10-10 09:03:32
# For each SNP, we provide the identifier, chromosome
# number, base pair position and genotype. The genotype
# is reported on the forward (+) strand with respect to
# the human reference build 37.
# THIS INFORMATION IS FOR YOUR PERSONAL USE AND IS
# INTENDED FOR GENEALOGICAL RESEARCH
# ONLY. IT IS NOT INTENDED FOR MEDICAL OR HEALTH
# PURPOSES. PLEASE BE AWARE THAT THE
# DOWNLOADED DATA WILL NO LONGER BE PROTECTED BY OUR
SECURITY MEASURES.


#RSID,CHROMOSOME,POSITION,RESULT
"rs4477212","1","82154","AA"
"rs3094315","1","752566","--"
"rs3131972","1","752721","AG"
"rs12562034","1","768448","--"
"rs12124819","1","776546","--"
"rs11240777","1","798959","GG"
"rs6681049","1","800007","--"
"rs4970383","1","838555","AC"
"rs4475691","1","846808","TC"
"rs7537756","1","854250","AG"
"rs13302982","1","861808","GG"
"rs1110052","1","873558","TG"
"rs2272756","1","882033","GG"
```

```
# MyHeritage DNA raw data.
# This file was generated on 2018-10-10 09:03:32
# For each SNP, we provide the identifier, chromosome
# number, base pair position and genotype. The genotype
# is reported on the forward (+) strand with respect to
# the human reference build 37.
# THIS INFORMATION IS FOR YOUR PERSONAL USE AND IS
# INTENDED FOR GENEALOGICAL RESEARCH
# ONLY. IT IS NOT INTENDED FOR MEDICAL OR HEALTH
# PURPOSES. PLEASE BE AWARE THAT THE
# DOWNLOADED DATA WILL NO LONGER BE PROTECTED BY OUR SECURITY
MEASURES.
#SIGNATURE=RZTcitAZ1bneCfURL5gsC5yRghb9=
#RSID,CHROMOSOME,POSITION,RESULT
"rs4477212","1","82154","AA"
"rs3094315","1","752566","--"
"rs3131972","1","752721","AG"
"rs12562034","1","768448","--"
"rs12124819","1","776546","--"
"rs11240777","1","798959","GG"
"rs6681049","1","800007","--"
"rs4970383","1","838555","AC"
"rs4475691","1","846808","TC"
"rs7537756","1","854250","AG"
"rs13302982","1","861808","GG"
"rs1110052","1","873558","TG"
"rs2272756","1","882033","GG"
```

Seamless for the user!

# Acknowledgments

**Tal Shor**
MyHeritage
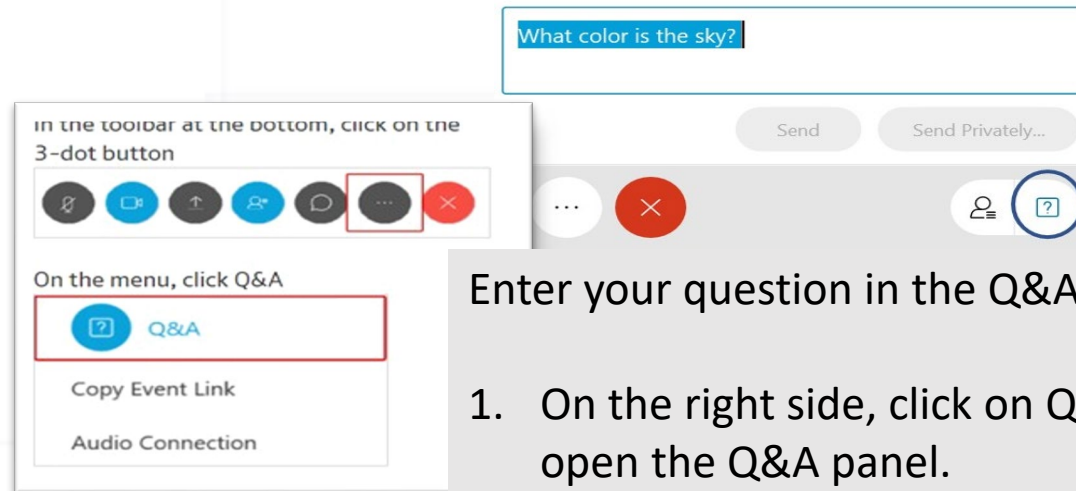
**Itsik Pe'er**
Columbia University

**Shai Carmi**
Hebrew University

# Keynotes: The Protection and Enabling Perspectives

Moderated Questions and Answers

In the toolbar at the bottom, click on the 3-dot button

On the menu, click Q&A

? Q&A

Copy Event Link

Audio Connection

What color is the sky?

Send    Send Privately...

Enter your question in the Q&A panel.

1. On the right side, click on Q&A header to open the Q&A panel.
2. Type in the box **your name, organization and question**.
3. Click send.

# Protecting and Sharing Genomic Data: a Swiss/European Perspective

Prof. Jean-Pierre Hubaux, EPFL

Co-Founder of Tune Insight SA

Work done in close collaboration with Lausanne University Hospital (CHUV)

With gratitude to all the colleagues I have had the privilege to work with

# About Switzerland

8.5 inhabitants

26 cantons (states), each with its own laws

Most of the health system managed by the cantons, not the federal government; the latter defines the overall policy and strategy

Data protection laws very similar to EU GDPR

Very strong political decentralization

One of highest GDP/capita in the world

Strong pharma: Roche, Novartis,…

5 university hospitals

2 federal institutes of technology: EPFL (Lausanne), ETH (Zurich)

# Use case for Swiss Personalized Oncology Project: federated analytics platform for research and molecular tumor board



**Q1**: How many adult cancer patients consenting on reuse of routine data for research with diagnosis of a malignancy on or after 1st January 2015, mutations in BRAF gene and under anti-PD-1 are there?

**Explore**

**Q2**: Among these patients, what is the overall survival for patients with and without a mutation on position 600 of the BRAF gene?

**Analysis**

# The Main Challenges we Faced

- Multi-disciplinary nature of the problem: bio-informaticians, clinicians, geneticists, hospital IT specialists, hospital lawyers, data protection authorities, ethicists, computer scientists
- Mess of the health data
- Financial sustainability of the solution

# Distributed Learning - Current Approaches

## a) Fully centralized

Data Leakage

*Raw data*

Examples:
- All of Us
- EGA
- Genomics England

## Meta-analysis

Data Leakage

*Aggregated data*

Trusted Party

https://covidclinical.net/

## Decentralized

Data Leakage

*Aggregated data*

http://www.datashield.ac.uk
Personalized Health Train (PHT)

## Differential Privacy
### Decentralized

Introduce Bias

*= Partial Results Obfuscation*

Examples:
- M. Kim et al. "Secure and Differentially Private Logistic Regression for Horizontally Distributed Data," TIFS 2019
- M. Abadi et al. Deep learning with differential privacy. In ACM CCS, 2016.
- Chaudhuri and C. Monteleoni. Privacy-preserving logistic regression. In NIPS, 2009.

## (e) Cryptographic (SMC, HE)
### Decentralized

Limited #parties

*Secret shared/encrypted*

Examples:
- A. Gascón et al.. Privacy-preserving distributed linear regression on high-dimensional data. PETS, 2017.
- P. Mohassel and Y. Zhang. SecureML: A system for scalable privacy-preserving machine learning. In IEEE S&P, 2017.

## Our approach

→ **Data + Model Confidentiality as long as 1 entity is honest**
→ No data outsourcing
→ Scales with #parties
→ Exact results

# Privacy-Preserving Federated Neural Network Learning

**Solution:** The data providers (DPs) collaborate to enable a joint gradient descent while protecting their security/privacy and **obtain a global and accurate model**
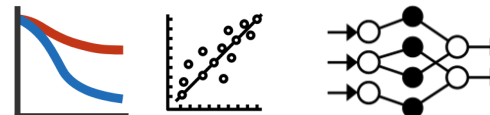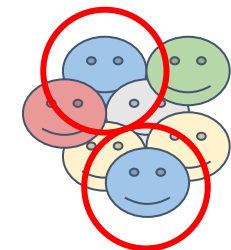


**Compute gradients**

**2**. Each DP **performs several training iterations** on its data**.**

iterations **global model**

**1**. Define the **task**, e.g., training of a neural network model.

**Compute gradients**

**3**. The DPs collectively and iteratively **combine their gradients (update) in a global model.**

**Compute gradients**

**4**. After the training, **the model** is **kept secret for oblivious predictions**

S.Sav, A. Pyrgelis, J.R. Troncoso-Pastoriza, D. Froelicher, J.P. Bossuat , J.S. Sousa and J.P. Hubaux,

**POSEIDON: Privacy-Preserving Federated Neural Network Learning.** NDSS, 2021

- **Distributed software platform** for federated cohort exploration and analytics of clinical and genomic data

- Co-developed by EPFL and CHUV

- Built on top of the i2b2 cohort explorer (i2b2 is used by 250+ hospitals worldwide)

- Relies on **advanced cryptographic techniques**
  → Multi-party homomorphic encryption (MHE)

- Code-reviewed and pen-tested by third-party industrial companies, compliant with hospitals' information security policies

- Main functionalities
  - **MedCo-Explore: cohort exploration**
    - Obtaining cohort sizes for clinical research studies based on inclusion/exclusion criteria
  - **MedCo-Analysis: federated analytics**
    - Survival analysis
    - ML training and testing

# April 2020: MedCo deployed at 3 hospitals

**EPFL software to enable secure data-sharing for hospitals**

The MedCo system aims to facilitate medical research on pathologies — such as cancer and infectious diseases — by enabling secure computations on decentralized data. The unique software has recently been deployed at three Swiss hospitals.

02.04.20

LINKS

- MedCo
- LDS
- Video

- First application:

  Swiss Personalized Oncology project:

  → melanoma data and beyond

- Planned deployment at Zurich University Hospital

- Ongoing international deployments: USA, NL, Italy, France

# Data Protection Impact Assessment (DPIA) for multisite medical data analysis (June 2021)

## Centralized approach with standard pseudonymization

| Threat | Threat likelihood | Threat impact | Risk | Risk level |
|---|---|---|---|---|
| Unlawful access to the system | Unlikely | High | Loss of data confidentiality | Moderate |
| Malicious use of the system | Possible | High | Loss of data confidentiality | High |
| Loss of data | Unlikely | Minor | Loss of data integrity, data unavailability | Minor |
| Data leak of host/cloud | Possible | High | Loss of data confidentiality | High |
| Collusion of host/cloud | Possible | High | Loss of data confidentiality | High |
| Corrupted or malicious host/cloud | Possible | High | Data unavailability, loss of data integrity, loss of data confidentiality, loss of data correctness | High |
| Unavailability of host/cloud | Possible | Minor | Data unavailability, loss of data correctness | Moderate |
| Re-identification/attribute inference | Possible | High | Loss of data confidentiality | High |

## Federated approach enhanced **with MedCo**

| Threat | Measure introduced with MedCo | Threat likelihood | Threat Impact | Risk | Risk level |
|---|---|---|---|---|---|
| Unlawful access to the system | 1 | Unlikely | Minor | Loss of data confidentiality | Low |
| Malicious use of the system | 1, 2, 4, 10 | Possible | Minor | Loss of data confidentiality | Low |
| Loss of data | 3, 5 | Unlikely | Minor | Loss of data integrity, data unavailability | Low |
| Data leak | 4, 5, 8, 9, 10 | Unlikely | Minor | Loss of data confidentiality | Low |
| Collusion between nodes | 4, 9 | Unlikely | Moderate | Loss of data confidentiality | Moderate |
| Corrupted or malicious nodes | 2, 5, 6, 7, 8, 9 | Unlikely | Moderate | Data unavailability, loss of data integrity, loss of data confidentiality, loss of data correctness | Moderate |
| Unavailability of of nodes | 6, 7 | Possible | Minor | Data unavailability, loss of data correctness | Moderate |
| Re-identification or attribute inference | 1, 2, 4, 9, 10 | Unlikely | Minor | Loss of data confidentiality | Low |

# Feedback from Swiss authorities on MedCo DPIA

Schweizerische Eidgenossenschaft
Confédération suisse
Confederazione Svizzera
Confederaziun svizra

**Federal Data Protection and Information Commissioner**

"... the threat impact of most risks with the MedCo system shows to be clearly lower than with traditional systems. Since data processed within the Medco framework remain encrypted during computation, an attacker would cause little damage. **As no entity has the full decryption key, it seems indeed unlikely that he could decrypt and abuse the stolen data**. ..."

13 September 2021

# GDPR legal compliance: partial aggregates are not personal data anymore, they are anonymous

**JMIR** Publications
Advancing Digital Health & Open Science

| Articles ▾ | Search articles | 🔍 |

🏠 Journal of Medical Internet Research    ↓    Journal Information ▾    Browse Journal ▾    Subn

Published on 25.2.2021 in Vol 23, No 2 (2021): February

📌 Preprints (earlier versions) of this paper are available at https://preprints.jmir.org/preprint/25120, first published October 19, 2020.

## Revolutionizing Medical Data Sharing Using Advanced Privacy-Enhancing Technologies: Technical, Legal, and Ethical Synthesis

James Scheibner [1, 2] 🆔; Jean Louis Raisaro [3, 4] 🆔; Juan Ramón Troncoso-Pastoriza [5] 🆔; Marcello Ienca [1] 🆔; Jacques Fellay [3, 6, 7] 🆔; Effy Vayena [1] 🆔; Jean-Pierre Hubaux [5] 🆔

13

# Truly privacy-preserving federated analytics for precision medicine with multiparty homomorphic encryption

David Froelicher, Juan R. Troncoso-Pastoriza, Jean Louis Raisaro, Michel A. Cuendet, Joao Sa Sousa, Hyunghoon Cho, Bonnie Berger, Jacques Fellay & Jean-Pierre Hubaux ✉

Metrics

## Abstract

Using real-world evidence in biomedical research, an indispensable complement to clinical trials, requires access to large quantities of patient data that are typically held separately by multiple healthcare institutions. We propose FAMHE, a novel federated analytics system that, based on multiparty homomorphic encryption (MHE), enables privacy-preserving analyses of distributed datasets by yielding highly accurate results without revealing any intermediate data. We demonstrate the applicability of FAMHE to essential biomedical analysis tasks, including Kaplan-Meier survival analysis in oncology and genome-wide

14

# FAMHE: Privacy-Preserving Federated Analytics for Precision Medicine with MHE - GWAS



(a) *Original* Approach

(b) FAMHE-GWAS

(c) *Meta-analysis* Approach

(d) FAMHE-FastGWAS

(e) Independent Approach

[Original approach] McLaren, P. J. et al. Polymorphisms of Large Effect Explain the Majority of the Host Genetic Contribution to Variation of HIV-1 Virus Load. Proc. Natl. Acad. Sci. 112, 14658–14663 (2015).

[FAMHE] Froelicher et al. Truly Privacy-Preserving Federated Analytics for Precision Medicine with Multiparty Homomorphic Encryption.

# FAMHE: Genome-wide association study

**Default**: 1857 patients spread among 12 data providers.

**→ scale in all dimensions**

a. *With the number of data providers*

b. *With the number of patients*

c. *With the number of variants*

[Centralized] McLaren, P. J. et al. Polymorphisms of Large Effect Explain the Majority of the Host Genetic Contribution to Variation of HIV-1 Virus Load. Proc. Natl. Acad. Sci. 112, 14658–14663 (2015).
[FAMHE] Froelicher et al. Truly Privacy-Preserving Federated Analytics for Precision Medicine with Multiparty Homomorphic Encryption. Submitted to Nat. comms. https://doi.org/10.1101/2021.02.24.432489

# FAMHE: Privacy-Preserving Federated Analytics for Precision Medicine with MHE - Survival curves (Kaplan-Meier)

Data split among 3
data providers:



TMB — Bottom 80% — Top 10%-20% — Top 10%

| Number at risk | | | | | |
|---|---|---|---|---|---|
| Time [months] | 0 | 12 | 24 | 36 | 48 |
| Bottom 80% | 1302 | 584 | 230 | 85 | 33 |
| Top 10%-20% | 186 | 103 | 41 | 17 | 4 |
| Top 10% | 172 | 99 | 41 | 13 | 4 |

[Centralized] Samstein, R. M. et al. Tumor Mutational Load Predicts Survival after Immunotherapy across Multiple Cancer Types. Nat. genetics 51, 202–206 (2019).

[FAMHE] Froelicher et al. Truly Privacy-Preserving Federated Analytics for Precision Medicine with Multiparty Homomorphic Encryption.

17

# Share without Sharing: Available Options

Hardware ✅

Software

"Vanilla" Federated Learning ❗

Protection of Partial Aggregates

Homomorphic Encryption (HE) ❗

By Crypto

By Adding Noise (differential privacy) ❗

Secure Multiparty Computation (SMC) ❗

**Multiparty Homomorphic Encryption** HE + SMC ✅

**TUNE INSIGHT**

# TUNE INSIGHT

## Enterprise Data & Analytics

$100B Market
(GlobalData, Enterprise Data and Analytics Market, 2020))

**56%**
say own data is not enough, expanding ability to source external data
(Forrester, The Insights Professional's Guide To External Data Sourcing, 2019)

However, organizations are **prevented** to enter valuable data collaborations due to fear of **data leaks** and **data protection regulations**

Cross-vertical enterprise SaaS enabling organizations to make better decisions, together, by orchestrating secure collaborations around their sensitive data.

- CHF400k in customer-paid projects including with Swiss Re, Armasuisse
- Pilot deployed at Swiss hospitals
- CHF100k EPFL Innogrant
- State-of-the-art post-quantum encryption technology
- Raised pre-seed with Wingman Ventures

**1.** Query connected organisations

**2.** Receive collective insight

**3.** Apply new knowledge

Encrypted computation

Better patient treatment
More personalized services
Better network management
Better intrusion & attack detection
More accurate premiums
Better risk estimation
Better market predictions
Reduced fraud losses

**Access to insights**    **Immediacy**    **Compliance**

**Personalization**    **Scalability**    **Control**

# MHE: mathematical proofs instead of vendor lock-in and side-channel attacks

| | Software-based solutions (MHE) | Hardware-based solutions (e.g., Intel SGX) |
|---|---|---|
| **System and trust model** | **Decentralized** (federated computing, edge computing) or **centralized** (outsourced) systems | **Only centralized** systems (data has to be transferred to the TEE) |
| **Assumptions** | Protection against passive adversaries with quantum computing power: **processing infrastructure (including side-channels) and other data providers** | Protection against passive adversaries (other tenants); **limited protection against the processing infrastructure**; protection against side-channels is implementation-dependent |
| **Implementation cost** | **Tailored solution**; application-specific design; composition of cryptographic building blocks; limited range of efficient functionalities | **Available SDKs**; relatively easy conversion to secure enclave; general-purpose solutions; limited libraries and memory inside the enclave |
| **Performance and overhead** | **Less than 10x** overhead when full packing capacity is utilized (federated training of GLMs and NNs). Up to 4-5 orders of magnitude overhead for non-optimized or non-packed solutions | **Negligible** overhead for **regular instructions**; **4x overhead for memory** copy operations; **35x overhead for syscalls** to/from enclave |
| **Response to newly discovered vulnerabilities** | **Software patch** with protocol update; usually, no re-encryption of the data is needed | **Firmware** patch with variable **performance impact** (1x to 20x slow-down); **architecture change and hardware replacement**; **enclave code update** (update signatures, keys, and require new attestation) |

| | | | |
|---|---|---|---|
| GLM | : Generalized Linear Model | SDK | : Software Development Kit |
| MHE | : Multi-party homomorphic encryption | SGX | : Software Guard eXtensions |
| NN | : Neural Network | TEE | : Trusted Execution Environment |

# International collaborations

- Prof. Xiaoqian Jiang, UT Health

- GA4GH Data Security Work Stream

- MedCo now part of the i2b2 official community projects

- Prof. Shawn Murphy, HMS, and the ACT Network

- Broad Inst. + MIT

- Cancer Institute of the Netherlands

- …

# Events devoted to the topic

- ## GenoPri.org: International workshop on genome privacy and security

  - Yearly workshop, typically co-located with GA4GH main annual event

- ## iDash - http://www.humangenomeprivacy.org/2021

  - Annual event with technical challenges on genome data protection and sharing

# Conclusion

- We have solved the problem of GDPR-compliant federated learning for medical data, including genomic data

- Solution: Multi-party homomorphic encryption (MHE)

  - Perform computations without "seeing" the data
  - Rely on decentralized trust and mathematical proofs
  - No need to transfer the data

- Scalability with the number of data providers and the size of the datasets

- Green light from the Swiss federal data protection authority

- Support and development of new features: provided by Tune Insight

Contact me at jean-pierre.hubaux@epfl.ch

More information at https://medco.epfl.ch

# Privacy Risks and Challenges from the Perspective of Individual Rights

**John Verdi**, Senior Vice President of Policy at the Future of Privacy Forum.

# Privacy Risks and Challenges from the Perspective of Individual Rights

*FPF Work:*

- In July 2018, the Future of Privacy Forum released [Privacy Best Practices for Consumer Genetic Testing Services](#)
- FPF developed the Best Practices following consultation with technical experts, regulators, leading consumer genetic and personal genomic testing companies, and civil society
- On January 1, 2022, California's Genetic Information Privacy Act (GIPA) became effective, codifying many of FPF's best practices

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Privacy Risks and Challenges from the Perspective of Individual Rights

*FPF Requirements:*

- Express consent for collection, use, and retention of genetic data;
- Separate express consent for transfer of to third parties and for incompatible uses;
- Informed consent for research;
- Educational resources about the basics, risks, benefits, and limitations of genetic and personal genomic testing;
- Access, correction, and deletion rights;
- Valid legal process for disclosure to the government and transparency reports;
- Ban on sharing genetic data with third parties (such as employers, insurance companies, educational institutions, and government agencies) without consent or as required by law;
- Restrictions on marketing based on genetic data; and
- Strong data security protections and privacy by design

# Privacy Risks and Challenges from the Perspective of Individual Rights

*Privacy Risks and Challenges of genomic data:*

- Unique, immutable biometric
- Potentially reveals information about identity
- Potentially reveals information about heritage
- Potentially reveals information about health
- Potentially reveals information about relatives' identities, heritage, and health
- Difficult or impossible to de-identify without undermining utility

# Privacy Risks and Challenges from the Perspective of Individual Rights

*Privacy Risks and Challenges of genomic data:*

- False identifications in criminal matters (evidence mishandling)
- Unexpected family connections and non-connections
- Dept. of Defense warning re: health tests and readiness reporting
- False identifications in criminal matters (remote relatives)
- Data breaches – e.g. 2020 GED Match law enforcement breach
- Re-identification attacks, e.g. cross-referencing clinical, research, and publicly available data sets

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Digital Biosecurity in a Modern Context

Charles Fracchia <charlesfracchia@gmail.com>

PGP: 4A06 3D3A B157 C3DE C31C 91B0 6E76 3F6A DA35 06C4

CEO BioBright / VP Data Dotmatics

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# PUBLIC DOMAIN ATTACKS ON BIOINFRASTRUCTURE

Indian Biomanufacturing

Tissue Regenix

Düsseldorf University

Miltenyi Biotec

Moderna EUA

Oxford University

Health Service Ireland

Pfizer/BioNTech Full Approval

WHO

**Jan 2020**

**Sep 2020**

**Dec 2020**

**Feb 2021**

**Aug 2021**

**Oct 2021**

First Cases

US: January 20th 2020
UK: 29th January 2020

Pfizer/BioNTech EUA

Janssen EUA

6.87 Billion Doses Administered

Malware on Sputnik News

Siegfried Holding

Vaccine Manufacturing

Americold

European Medicines Agency

Life Sciences Institute

# WE CANNOT FIX THE ISSUE WITHOUT TACKLING DEVICE INSECURITY



https://www.youtube.com/watch?v=7du1TltZOJg

**Devices are embedded at every step of the biological process.**

# TARDIGRADE APT ON THE BIOECONOMY



BIO-ISAC Tardigrade Amplify Alert

# Who is HudsonAlpha?

HudsonAlpha Institute for Biotechnology is a nonprofit institute founded in 2008 specializing in genetics and genomics research and biotech education.

Tens of thousands of genomes (human and non-human) sequenced per year on campus

Sequencing use cases include:
- Genetic testing
- Clinical genome sequencing
- Genomic screening programs
- Original plant genome sequencing

We also host more than 50 associate companies on our campus - all of them involved in bioscience and many performing genomic sequencing in their own labs

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

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

NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# Cybersecurity Challenges

**Campus and Entrepreneurial Mission**

- **Associate companies provide their own sequencers**

- **Combines all the challenges of IoT and BYOD**

# Cybersecurity Challenges

**Sequencers are essentially IoT devices**

- **Internet connection required**

- **Dedicated PC**

- **Unknown software**

- **Software firewalls**

# Cybersecurity Challenges

**Availability vs. Security**

- Security takes a back seat to availability and accuracy

- Software updates can be problematic

# Cybersecurity Challenges

**Lack of Security Standards**

- **No guidelines or standards (e.g., STIG)**

**Anticipated proliferation**

- **Decreasing cost and more widespread use will lead to more attacks**



*Cost per Human Genome*

# CONSIDER A HOLISTIC APPROACH

# SECURITY REQUIREMENTS

- Architecture
  - Consider zero-trust architecture
- Robust and auditable authentication and authorization
- Continuous Monitoring
  - (24x7x365) with a high degree of automation
  - Automate the Incident Response Procedure as much as possible
- Accountability to industry standards:
  - FedRAMP "Moderate," "High" or similar standards (NIST 800-family)
  - ISO27000 family of standards
  - Discovery & sampling audit at frequent intervals (trust & verify)

# PRIVACY REQUIREMENTS

- Compliant with the applicable privacy regulations
  - GDPR for citizens of the EU
  - Federal, state, and local regulations regarding PII and PHI
- Consent management
  - Ongoing management of contributor's consent with audit trails
  - Managing familial scans (law enforcement)
- Automated data privacy scans
  - Support for DSARs, etc.



**Are long read genomic sequences inherently identifiable?**

# DATA REQUIREMENTS

- **C**onfidentiality, Integrity and Availability (CIA)
  - Encrypted data in transit and at rest
  - Evidence proving the accuracy and consistency throughout the lifespan of the data
  - Code/data are available to only those who are authorized within time limits

- **G**overnance
  - Establishing the appropriate authorizations on code/data
  - What must be retained? Where?  For how long?

- **P**rovenance
  - Privacy concerns, alignment with consent, ability to track and prove compliance
  - Metadata – source of sequence, sequencer, processing steps

- **Q**uality
  - How do  you measure quality per use case?

# SOFTWARE REQUIREMENTS

- Unconstrained by the size of the sequences and "-omics" type

- Horizontal and vertical scaling – application, metadata processing, etc.

- Walled Garden vs Open Federation
  - Users add their own software, which could contain malware, crypto mining, unsupported and vulnerable supporting libraries (log4j)

- License Management in leveraged software (Asset Mgmt)
  - Open Source Software – permissive, viral light, highly viral, Affero
  - Commercial Software – terms need to balance with usage

- Auditability
  - Immutable logs showing *every action performed on every object*

# NETWORK REQUIREMENTS

- Volume and velocity
  - Increasing at 40+% per year

- Centralized repository or Federated repository
  - Network usage varies on the overall design

- Enforcement of data localization regulations and agreements
  - Privacy regulations require data to stay within a country/region's jurisdictions and how the data can be used.
  - Commercial agreements control where the data resides, how it can be accessed and used.

# CLOSING THOUGHTS

- Zero-Trust - verify and corroborate

- Track-and-Trace – anything you say, be prepared to prove it!

- Governance

Balancing Security, Privacy, Quality, and Science

# Privacy in the Genomic Era

XiaoFeng Wang, IEEE Fellow, Rudy Professor at IUB

http://www.informatics.indiana.edu/xw7

# Genomic Revolution

- Fast drop in the cost of genome-sequencing
  - ➢ 2000: $3 billion
  - ➢ Mar. 2021: $800 - $1,000
  - ➢ Genotyping 1M variations: below $200

- Unleashing the potential of the technology
  - ➢ Healthcare: e.g., disease risk detection, personalized medicine
  - ➢ Biomedical research:  e.g., geno-phono association
  - ➢ Legal and forensic
  - ➢ DTC:  e.g., ancestry test, paternity test
  - ……



Legend:
- Cost per Genome
- Moore's Law

# Genome Privacy

- Privacy risks
  - Genetic disease disclosure
  - Collateral damage
  - Genetic discrimination …
- Protection
  - Clear access policies
  - Accountability
  - Data anonymization
  - Best practice for data privacy
  - Privacy awareness ……



WORLD ECONOMIC FORUM
COMMITTED TO IMPROVING THE STATE OF THE WORLD

White Paper

**Genomic Data Policy Framework and Ethical Tensions**

June 2020



**PRIVACY** and **PROGRESS** in Whole Genome Sequencing

Presidential Commission *for the* Study of Bioethical Issues

October 2012

**For more information: Privacy and Security in the Genomic Era** by Naveed,  E. Ayday, E. Clayton, J. Fellay, C. Gunter,  JP  Hubaux, B. Malin and X. Wang

Available at http://arxiv.org/pdf/1405.1891v1.pdf

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

NCCoE
NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

# Disclosed Genomic Data can be Abused

- Hawasupai case (1989): use of Indian tribe genome data without proper informed consent, with impacts on NIH's All of Us project (2020)

- Genomic data for solving crimes (with privacy implications)
  - E.g., Capture of the Golden State Killer (through GEDmatch)
  - But privacy concern is raised: how one's individual choice affects others?



The New York Times Magazine

**Your DNA Test Could Send a Relative to Jail**

Thanks to "genetic genealogy," solving crimes with genomic databases is becoming mainstream — with some uncomfortable implications for the future of privacy.

# Unauthorized Disclosure of DNA/Meta Data Continue to Happen

- DNA Diagnostics Center (DDC), breach more than 2.1 million people (2021)

- GEDmatch hack causes email addresses from its users to be used in a phishing attack on another leading genealogy site (2020)

- Veritas Genetics claim a data breach resulted in unauthorized access of some customer information (2019)

……

# Genomic Privacy: Technical Challenges

- Dissemination: privacy protection is difficult !

  ➢ Anonymization is hard: genotype to phenotype mapping

  ➢ Impact of genetic genealogy

  ➢ Extremely high dimensions: hard to balance between privacy and utility


- Computing: big data analysis

  ➢ Beyond the capability of existing secure computing technologies
    ➢ NIH originally disallows reads with human DNA to be given to the public Cloud
    ➢ Now, use the cloud at your own risk

# Challenge in Privacy-preserving Genomic Data Sharing

- Old problems:
  - Statistical inference control, access control, query auditing…
- However, genome data are special:
  - Special structures, e.g. linkage disequilibrium
  - Existence of reference genomic data that are publicly available (e.g. large population studies as HapMap, WTCCC, 1000 Genome)

- Examples:
  - Homer's attack and NIH's responses (2008)
  - Our analysis on test statistics released by GWAS papers (2009)
  - Shringarpure and Bustamante's attack on beacons (2015)

# iDASH Genomic Data Privacy and Security Protection Competition

**Since 2014, http://www.humangenomeprivacy.org**

An interdisciplinary challenge on genomic privacy research

Motivated by real world biomedical applications and with participation of privacy technology experts, Biomedical and ELSI researchers (academia and industry)

Develop practical solutions for privacy preserving genomic data sharing and analysis

Demonstrate the feasibility of secure genome analysis and dissemination using DP, MPC, HE, TEE

Reported in the media (e.g., Nature News)

# Topics and Trend from 2014 to 2021

| | |
|---|---|
| Privacy-preserving Data Sharing | Encryption Testing |
| Secure Release | De-duplication |
| Secure Outsourcing | Software Guard Extensions |
| Homomorphic Encryption | Secure Search |
| Secure Collaboration | Blockchain and Smart Contract |
| Secure Multiparty Computation | Secure Machine Learning |
| Beacon Service | Privacy-preserving Machine Learning |
| Privacy-preserving Search | |



Participants and countries in iDash

# Participation Around the World



- Academia:  Cornell,  MIT, UTHealth, UCSD, Yale, Purdue, Vanderbilt, EPFL, SNU , CUHK, Manitoba …
- Industry:  IBM, MSR, Samsung, Alibaba, Tencent, Baidu …
- Government: Sandia National Lab, French Alternative Energies and Atomic Energy Commission …

# Contributions to the Progress in Genome Privacy

For the task secure multi-label tumor classification using Homomorphic Encryption in 2020, most teams are utilizing linear/logistic   regression models to implement cancer classification. These models have been improved significantly over the past few years in the HE competition, which is quite scalable and efficient now. The top solutions achieved a Micro-AUC of 0.97 to classify 11 cancer types from encrypted genetic variants of 909 samples within 5 minutes.

For the task differentially private federated learning for the cancer prediction model in 2020, the submitted solutions achieved almost perfect model accuracies while enforcing a high differential privacy standard (privacy budget of 3.0 or lower). The training process of the best-performing solution is very fast, comparable with the efficiency of training a machine learning model with all data by a single party.

For the task of data sharing consent for health-related data using contracts on blockchain in 2021, it is feasible to store patient consent sharing preference records for seven categories for a given clinical/genomic study on blockchain up to ~6,800 records per hour (or ~1.889 records per second).

# Acknowledgement

# More Information:

1. Learning Your Identity and Disease from Research Papers: Information Leaks in Genome Wide Association Study, 2009, ACM CCS

2. Addressing Beacon re-identification attacks: Quantification and mitigation of privacy risks, 2017, JAMIA

3. Real-time Protection of Genomic Data Sharing in Beacon Services, 2018, AMIA

4. A Secure Alignment Algorithm for Mapping Short Reads to Human Genome, 2018, RECOMB

5. MBeacon: Privacy-Preserving Beacons for DNA Methylation Data, 2019, NDSS

6. Haplotype-based membership inference from summary genomic data, 2021 Bioinformatics

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

NCCoE
NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

# Cybersecurity Challenges

Moderated Questions and Answers

In the toolbar at the bottom, click on the 3-dot button

On the menu, click Q&A

Q&A

Copy Event Link

Audio Connection

What color is the sky?

Send    Send Privately...

Enter your question in the Q&A panel.

1. On the right side, click on Q&A header to open the Q&A panel.
2. Type in the box **your name, organization and question**.
3. Click send.

# Million Veteran Program (MVP) Overview

- MVP is a **national VA research program**, launched in 2011, designed to advance precision health care by learning how genes, lifestyle, and military experiences and exposures affect health and illness
  - Establish a comprehensive, diverse cohort of at least one million Veterans
  - Provide broad access to the data for scientific discovery
  - Establish pipelines to translate discoveries to the clinic to improve the health of Veterans
- MVP is one of the world's largest healthcare system-based research programs of its kind with **over 864,000 Veterans enrolled (as of Dec. 2021)**

# MVP Data Universe

# MVP Biospecimen Data Overview

*GOAL: Generate the maximum amount of data from biospecimens
to enhance scientific discovery*

- Baseline genetic data profile (genotype) generated for all participants
  - Data from 650,000 samples are currently provided to approved researchers
- Genetic data for specific ethnic groups (Blacks, Hispanics and Asians) using a customized analytic tool currently underway for ~ 200K participants

- Whole genome sequences have been generated on ~ 140,000 samples

  - Processing underway

- Other data such as proteomics and metabolomics are being piloted

# Balancing Data Privacy/Security and Access

- Bring researchers to the data in a central secure scientific computing platform
  - Computing infrastructure within the VA meets VA IT Privacy and Security requirements; DOE and the University of Chicago VA Data Commons have an approved VA authority to operate (ATO)
- Biospecimen (blood sample) and data (surveys) collected are labeled using a code instead of identifiable information
- New ship ID created for sample send-outs to vendors
- Crosswalk to identity of participant is held by few authorized core staff
- Researchers access only coded data (no direct identifiers such as SSN, name, date of birth, street address)
- Researchers sign rules of behavior and analyze data in a central, secure computing system
- No data leaves the system; only summary results can be taken out

# MVP Data Access Model

# External access to MVP Data

- VA Data Commons : allow data access broadly to investigators within and outside the VA

- Contract with the University of Chicago (UoC)

- Data deidentified at the VA and moved to UoC

  - Safe harbor method plus

  - Formal expert statistical determination

- Data will be migrated to a cloud compute infrastructure for many simultaneous approved users

- Beta-testing in FY 2022 ; piloting in FY23

# Reidentification Risk

- **Re-Identification:** the ability to determine whether an individual is included in a pooled sample, based on the allele frequencies in the pool -- without the need to access individual-level genotype data of that pooled data set

- All the published references discussing re-identification are theoretical, not actual case reports of participant re-identification

- In order for re-identification to occur, the user must already have access to that person's genetic information from another source

- Accuracy of re-identification is determined by:
  - the size of the population (small sample size = better accuracy)
  - the diversity of the population (homogenous population = better accuracy)
  - the frequency of the genetic variants (rare genetic variants = better accuracy)

# MVP Risk Mitigation Strategies

- MVP is sufficiently large and diverse, therefore theoretical re-identification risk is extremely low
  Only aggregate results will be shared, no individual-level data

- Additional steps taken to further reduce risk
  - Results filtered to only include genetic variants with a minor allele count > 30 or minor allele frequency > 0.005, whichever is less (metrics should be based on the subset of the study population actually used for the analysis, not the general population)
  - Total study population used for the analysis must be >3000 participants
  - If a case-control study, there must be >500 cases in the analysis

# Thank you!

# Privacy Challenges for Genomic Data

Moderated Questions and Answers

What color is the sky?

Send    Send Privately...

In the toolbar at the bottom, click on the 3-dot button

On the menu, click Q&A

Q&A

Copy Event Link

Audio Connection

Enter your question in the Q&A panel.

1. On the right side, click on Q&A header to open the Q&A panel.
2. Type in the box **your name, organization and question**.
3. Click send.

# American Society of Human Genetics

- **Mission**: *Advance human genetics and genomics in science, health, and society through excellence in research, education and advocacy*

- **Vision**: *People everywhere realize the benefits of human genetics and genomics research*

- **Annual Meeting**: *Attracts up to 9,000 attendees*

- **Year-Round Scientific Programs**

- **Two Scientific Journals**:
  - *American Journal of Human Genetics*
  - *Human Genetics and Genomics Advances*

# Data-sharing Fuels Progress in Human Genetics & Genomics Research

- Broad data-sharing a hallmark of the human genetics community
  - Essential for completion of Human Genome Project

- Data-sharing fundamental for continued advances in research & medicine

# Policies and Systems Need to Maintain Privacy of Research Participants

- Acquisition, analysis, sharing of human genetic data, use of genetic tools, need to be conducted responsibly

- ASHG supports policies that strengthen research participant privacy
  - Genetic Information Nondiscrimination Act
  - 21st Century Cures Act
  - Common Rule
  - NIH Genomic Data Sharing Policy

# Privacy/Security Essential for Public Participation in Genetics Research

What are the biggest barriers or concerns to your participation in human genetics research? (Choose all that apply)

| | |
|---|---|
| Concerns about the security of the database where my data would be stored | **48%** |
| Concerns about the privacy of my genetic information | **47%** |
| Concerns that my genetic information could influence my access to health care | **40%** |
| There are no studies currently available to me | **34%** |
| Concerns that my genetic information could influence my access to insurance or a job | **34%** |
| I cannot participate due to financial reasons | **28%** |
| Concerns that my genetic information could reveal unwelcome information about others or myself | **25%** |
| Participating in a study will take too much time | **24%** |

*Source: A Research!America poll of U.S. adults conducted in partnership with Zogby Analytics in December 2019*

# Policies and Systems Must Enable Science, Maintain Privacy

**Data-sharing**

**Privacy**

Advance science

Protect participants

# DNA researchers question Senate bill's security provisions

Measure aimed at stopping China from misusing human genome data could harm research efforts, groups argue

By **Jocelyn Kaiser**

A provision buried in a 2400-page bill approved last week by the U.S. Senate to help the United States compete with China is drawing fire from human genome researchers. It would require the National Institutes of Health (NIH) to develop new security protocols aimed at preventing the misuse of U.S.-funded genomic data by China and other nations.

The provision is not based on substantiated security risks, and "could slow biomedical advances and impose unintended burdens," the American Society of Human Genetics (ASHG) warned last week in a letter to lawmakers. The Association of American Medical Colleges cautioned in a statement that "any additional protections or restrictions ... should be commensurate with the actual risk."

Research advocates are applauding many provisions of the huge Senate bill, the United States Innovation and Competition Act (S. 1260), which calls for increasing federal research spending and creating a technology directorate at the National Science Foundation (*Science*, 21 May, p. 777). But they're less enthusiastic about a provision reflecting concerns that China is amassing DNA data on

> "Any additional protections or restrictions ... should be commensurate with the actual risk."
>
> Association of American Medical Colleges

national security risks." NIH must work with intelligence agencies to issue, within 1 year, "a comprehensive framework" for managing risks, such as requiring more training for NIH-funded investigators and peer reviewers and including security experts on data access panels.

In the past, NIH has argued that existing security measures are adequate. Researchers already strip identifying information from genome data, and NIH reviews, and sometimes rejects, scientists' requests for access. But in 2019, the Office of Inspector General (OIG) of the Department of Health and Human Services, NIH's parent agency, suggested NIH do more, for example by adding controls on foreign scientists who use U.S. genome data.

In a response to OIG, NIH questioned the severity of the threat. It noted security worries were largely based on "a single Congressional testimony," from FBI agent Edward You, who has long warned of the risks of sharing genomic research data. Fears of economic harm were "theoretical," NIH said, noting that many experts argue that sharing data promotes innovation. And it scoffed at the "improbability" of weaponizing human genetics data. Research would "come to a halt," NIH said, if it had to write craft

# Thank you!

*Contact:  gjarvik@medicine.washington.edu*

# AWS for Genomics

## Solving Challenges in Genomic Data Sharing

*Ankit Malhotra, Ph.D.*
*Genomics Lead,*
*AWS Worldwide Public Sector Health*

National Institute of Standards and Technology
U.S. Department of Commerce
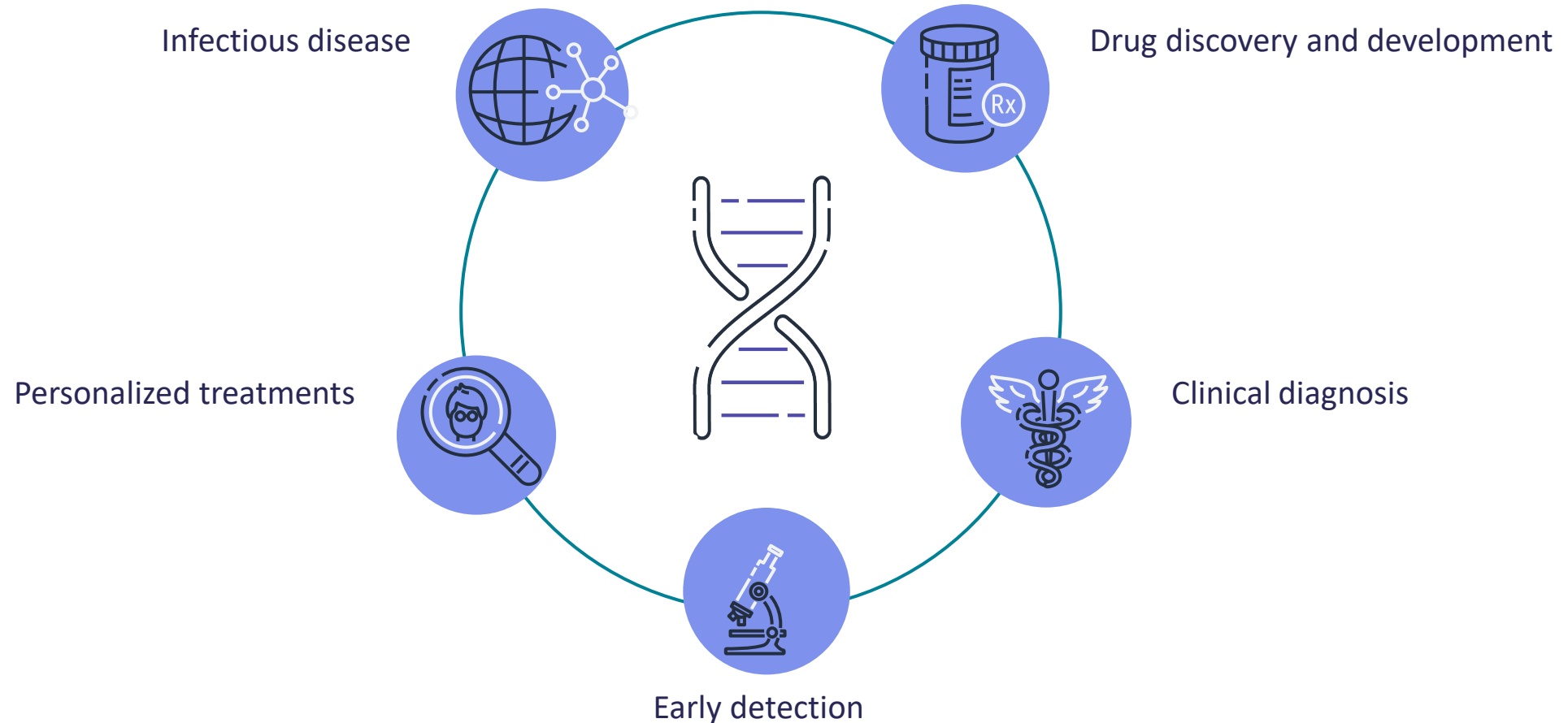
NCCoE
NATIONAL CYBERSECURITY
CENTER OF EXCELLENCE

# The precision medicine revolution

Transformative technologies in sequencing and computing is driving innovation across healthcare and enabling precision medicine.

Clinical Genome Sequencing

Genetic Risk Scores

Targeted Therapeutics

Induced Pluripotent Stem Cells

CRISPR Genome Editing

# Genomics—a catalyst for personalized health



Infectious disease

Drug discovery and development

Personalized treatments

Clinical diagnosis

Early detection

# Challenges in leveraging genomics data

Large volumes of data needs to be transfer, stored, analyzed

Sequencing and analysis requires immense processing power, time

Frequently requires integration of multi-modal datasets

Protected health information must be secured



NovaSeq6000™, v1.5 reagents have delivered the $600 genome

# Genomics on AWS

**Data Transfer & Storage**

Trusted partner for secure data transfer, life cycle management, storage cost optimization and digital preservation

**Secondary Analysis & Workflow Automation**

Manage multiple workflows, accelerate, simplify and scale data analysis with both flexibility and reproducibility

**Data Aggregation & Governance**

Harmonize multi-omic datasets and govern robust data access controls and permissions across a global infrastructure

**Interpretation & Deep Learning**

Turn big genomic data into actionable insights with a rich layer of sophisticated solutions and services

## BUILD

AWS Modern Genomics Data Platform

AWS Genomics CLI

SageMaker, HealthLake, Sequera Labs, Cromwell

## BUY

databricks | ancestry | genuity science | illumina | paradigm4

Oxford NANOPORE Technologies | Helix | IQVIA | REALM

Seven Bridges | GRAIL | DNAnexus | LIFEOMIC | lifebit

## AWS Open Data Program



Registry of Open Data on AWS — The Human Microbiome Project

Registry of Open Data on AWS — 1000 Genomes

Registry of Open Data on AWS — Encyclopedia of DNA Elements (ENCODE)

Registry of Open Data on AWS — TCGA on AWS

At last count there were 301 datasets (87 life science) hosted on AWS S3 as part of Registry of open data on AWS

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

**NCCoE** NATIONAL CYBERSECURITY CENTER OF EXCELLENCE

# Open Access to Top Genomics Datasets

**AWS hosts a variety of public datasets that anyone can access for free. Below are just a few examples**

- 1000 Genomes Project
- The Cancer Genome Atlas
- International Cancer Genome Consortium
- 3000 Rice Genome
- Genome in a Bottle (GIAB)
- The Genome Modeling System
- Medicare Drug Spending
- The Human Connectome Project
- The Human Microbiome Project
- OpenNeuro
- Physionet
- Tabula muris
- gnoMAD
- and more….

# Security

AWS supports 98 security standards and compliance certifications, including HITRUST, GDPR compliance, FedRAMP, ISO 27001, and HIPAA.

**Whitepaper -** **https://docs.aws.amazon.com/whitepapers/latest/aws-overview/security-and-compliance.html**

## AWS shared responsibility model



### AWS shared responsibility model

**Customer**

Responsibility for security **"IN"** the cloud

AWS is responsible for protecting the infrastructure that runs all of the services offered in the AWS Cloud.

**AWS**

Responsible for security **"OF"** the cloud

Genomics organizations control access to and management of their data, includes data access permissioning.

Applications

Building Blocks

Regulated Landing Zone

AWS Services

# AWS security, identity, and compliance solutions

| Identity and access management | Detective controls | Infrastructure protection | Data protection | Incident response | Compliance |
|---|---|---|---|---|---|
| AWS Identity and Access Management (IAM) | AWS Security Hub | AWS Firewall Manager | Amazon Macie | Amazon Detective | AWS Artifact |
| AWS Single Sign-On | Amazon GuardDuty | AWS Network Firewall | AWS Key Management Service (KMS) | Amazon EventBridge | AWS Audit Manager |
| AWS Organizations | Amazon Inspector | AWS Shield | AWS CloudHSM | AWS Backup | |
| AWS Directory Service | Amazon CloudWatch | AWS WAF – Web application firewall | AWS Certificate Manager | AWS Security Hub | |
| Amazon Cognito | AWS Config | Amazon Virtual Private Cloud | AWS Secrets Manager | CloudEndure Disaster Recovery | |
| AWS Resource Access Manager | AWS CloudTrail | AWS PrivateLink | AWS VPN | | |
| | VPC Flow Logs | AWS Systems Manager | Server-Side Encryption | | |
| | AWS IoT Device Defender | | | | |

## Genomics England Develops Genomic and Health Information Platform on AWS to Turn Science into Healthcare

### Challenge

Through the 100,000 Genomes Project alone, GEL amassed 50 petabytes of data. Seeking to make the data accessible to the research community, GEL is in the process of migrating its data to AWS to enable democratized access.

### Solution

GEL is working with AWS to use compression technologies and other advanced tools to optimize cloud storage and analysis of genomic data based on the field's specific needs

### Benefits

- To make genomic healthcare a reality, GEL is transitioning from project to platform, using Amazon Web Services (AWS) tools to give researchers reliable, comprehensive, and privacy-compliant access to these massive datasets. Through secure collaboration and analysis, this initiative will inform diagnoses, drive drug development, and unlock the future of precision medicine.

## AstraZeneca is Raising the Bar with Running its Genome Sequencing Pipeline on AWS

### Challenge

AstraZeneca's Centre for Genomic Research (CGR) has a bold target to analyze 2m genomes by 2026. However

- On-premise compute resources, which limit the performance capacity
- Dependency upon 3rd party informatics providers
- Orchestration of bioinformatics pipeline was time-intensive therefore costly and hard to scale

### Solution

With support from AWS Professional Services, AstraZeneca built a highly scalable and high performance sequence data processing pipeline on AWS. The solution leveraged FPGA instances for compute and extensive use of Step Functions, Lambda, SQS and AWS Batch. The output is stored in scalable and highly secure AWS managed databases and S3 storage.

### Benefits

- The bespoke pipeline was able to increase processing time by 2400%
- The results has been used to provided scientists advanced access to the clinical effects of natural mutations in humans that mimic drug inhibition/suppression

# Impact of the pandemic

The SARS-Cov-2 pandemic caused widespread impact for healthcare systems around the world, and brought genomic sequencing and testing into the public eye.

This has also brought forth challenges in global data sharing:

- Privacy concerns - data misuse may lead to infringement of privacy for individuals and their relatives
  - Need for novel approaches to data anonymization for research use

- Compatibility and aggregation
  - accessing and reconciling duplications/differences from distributed data sources hindered meta-analyses

- Real / Near real time data ingestion

# Use cases: genomic information in the cloud

**NCI Genomic Data Commons**
*https://aws.amazon.com/solutions/case-studies/university-of-chicago-case-study/*

**University of Chicago Biomedical Research Hub (Gen3)**
*https://academic.oup.com/jamia/advance-article/doi/10.1093/jamia/ocab247/6432980*

**Hong Kong Genome Project (LifeBit)**
*https://lifebit.ai/blog/lifebit-awarded-a-four-year-contract-for-hong-kongs-genome-project/*

**CanCOGeN, Genome Canada, Illumina**
*https://www.genomecanada.ca/en/cancogen*

**GISAID**
*https://www.gisaid.org/*

**Undiagnosed Disease Network, Harvard School of Medicine (Service Workbench)**
*https://aws.amazon.com/blogs/publicsector/solving-medical-mysteries-aws-cloud-medical-data-sharing-innovation-undiagnosed-diseases-network/*

**UK Biobank (DNANexus)**
*https://www.ukbiobank.ac.uk/enable-your-research/research-analysis-platform*

# NIH PRIVACY & DATA SHARING RESEARCH

**NHGRI - Mission of responsible Genomics Data Sharing**
- ELSI – "Ethical, Legal, and Social Implications" program
- Technical privacy portfolio - research grants and small business
  - Homomorphic encryption, Secure Multiparty Computation, Differential Privacy, Secure Enclaves, Machine learning with privacy, etc.

**NIH - Public Trust is the currency of the realm**
- Data sharing through cloud commons (one copy instead of many copies)
  - AnVil, BioData Catalyst, CRDC, Kids First etc.
- Federated data sharing
- RAS Researcher Auth System
- DUOS - pilot of data access automation
- Privacy Preserving Record Linkage (PPRL)

**ODSS - Catalyze modern computing at NIH**

- **Office of Data Science Strategy – Susan Gregurick**

- https://datascience.nih.gov/

# GENOMICS PRIVACY PORTFOLIO AT NIH

| | | |
|---|---|---|
| NCI | NEI | NHLBI |
| NHGRI | NIA | NIAAA |
| NIAID | NIAMS | NIBIB |
| NICHD | NIDCD | NIDCR |
| NIDDK | NIDA | NIEHS |
| NIGMS | NIMH | NIMHD |
| NINDS | NINR | NLM |
| CC | CIT | CSR |
| FIC | NCATS | NCCIH |
| OD | | |

## 27 NIH Institutes and Centers

NHGRI National Human Genome Research Institute, NCI National Cancer Institute, NHLBI National Heart, Lung, & Blood Institute

NICHD National Institute of Child Health & Development
- Genomics & health data

NEI National Eye Institute
- Retinal scans

NIDCR National Institute of Dental & Craniofacial Research
- Facial and dental images

NIBIB National Institute of Biomedical Imaging & Bioingineering
- Imaging & signal data

NIEHS National Institute of Environmental Health Sciences
- Geolocation

NCATS National Center for Advancing Translational Sciences
- N3C COVID patient data

# HUMAN PANGENOME FAQS

**International collaboration at foundation of genomics**

- 1'st Human genome $3B
- Now millions at $1000 each
- Simple codes build complexity
  - Genome: 4-letter code (A,C,T,G)
  - Computers: (0,1)
- Need huge numbers to decipher signals and interpret genome data

# NHGRI HUMAN PANGENOME REFERENCE

- Pangenome graph replaces linear "single genome" reference
- Represents global human diversity
- Enables population scale analysis
- Graph compresses large number of genomes into compact form
- "Subway map" of human journey



HUMAN PANGENOME



Reference structural variant pangenome graph

https://humanpangenome.org/

# FACIAL RECOGNITION OF GENETIC SYNDROMES



- Many (1000's) genetic syndromes have distinct facial dysmorphisms.
- Sometimes gene defect detected in facial scan of "normal" relative

Hong et al. Genetic syndromes screening by facial recognition technology, *Orphanet J. of Rare Diseases* (2021)

Face2Gene: "...facial recognition software to aid clinical diagnoses of thousands of genetic conditions, such as Sotos syndrome (cerebral gigantism), Kabuki syndrome, intellectual disability, Down syndrome, etc."

# POLYGENIC RISK SCORES



Ali Torkamani et al. "The personal and clinical utility of polygenic risk scores" *Nat. Rev. Genet*. (2018)

## NHGRI Consortia

GREGoR (previously Mendelian Centers) to find single gene cause of disease (cystic fibrosis, progeria, etc)

PRIMED to study polygenic ("complex") diseases using advances in Polygenic Risk Scores (diabetes, heart disease, autism, etc.)

# PRECISION HEALTH POWERED BY GENOMICS

**PATH TO PERSONALIZATION**
To tailor health care to individuals, information from various sources must be brought together. These data, both genetic and environmental, should be drawn from diverse populations.



Mark McCarthy & Ewan Birney, "Personalized profiles for disease risk must capture all facets of health" *Nature* 597, 175-177 (2021)

- Polygenic & Pangenome with diverse genomes to represent the human family

- Environmental, lifestyle, income, access to health care, exposures, culture

- "This will inevitably bring the realms of research and clinical care together, and will require us to address fundamental questions about data ownership, privacy, equality of access, fairness and social responsibility. Global efforts to create such standards are in place, for example the Global Alliance for Genomics and Health."

# GLOBAL DATA STANDARDS

Global Alliance
for Genomics & Health

Collaborate. Innovate. Accelerate.

➢ GA4GH is international collaboration on standards for genomics and health data with human rights foundation

➢ Many social and technical "onramps" for inclusion and adoption

➢ Example: GA4GH Passports & Visas

➢ *Cell Genomics* special issue Nov 2021

GA4GH: https://www.ga4gh.org/

# Alzheimer's Disease Sequencing Project (ADSP) and NIAGADS

- 17k complete genomes released in 2021

- 36k (2022), 75k (2023) planned

- Data releases are managed by NIAGADS
  - National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site at University of Pennsylvania
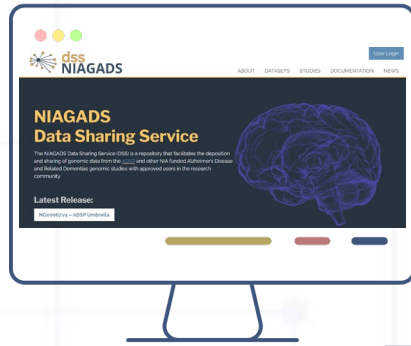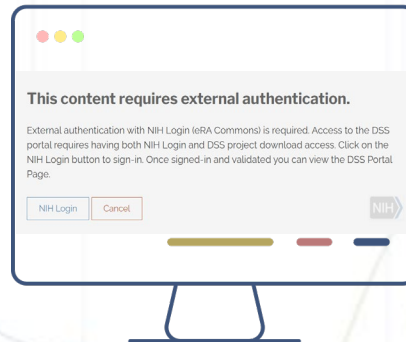


R3 ADSP WGS- 16,906 genomes

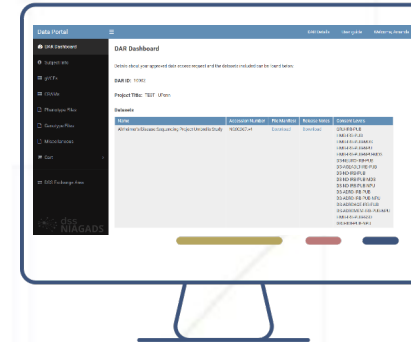# Authentication of data requesters



Step 1: https://dss.niagads.org

Step 2: NIH eRA Authentication

Step 3: Data manifest browsing on NIAGADS DSS

Step4: Data accessing

# Learning curve for FISMA is steep

- Scope

- Amount of work
  - 10 months preparation
  - 1 year for 3 times external assessment

- Regulations to be met

- Cost benefit analysis

# Respecting Informed Consent

- NIH Genome Sharing Policy
- Institutional certification to capture informed consent conditions
- Data access committee
- How to split data based on informed consent

# De-identification of genomic data

- Genome sequencing is identifiable

- What about functional genomics data? Theoretically RNA-Seq is identifiable because it carries variants. What other types of sequencing data?

- What does it take to de-identify data?

# What will be helpful if we do it again

- Guidelines for FISMA requirements that are specific to human genome data

- Tutorials and FAQs on how to set up a FISMA compliant framework: timeline, amount of work, cost (budget and staff)