# Genome Hacking
## Yaniv Erlich

@erlichya

# We need to share genetic information

**Hereditary Spastic Paraparesis**
(Erlich et al.)



**Joubert syndrome**
(Endevson et al.)


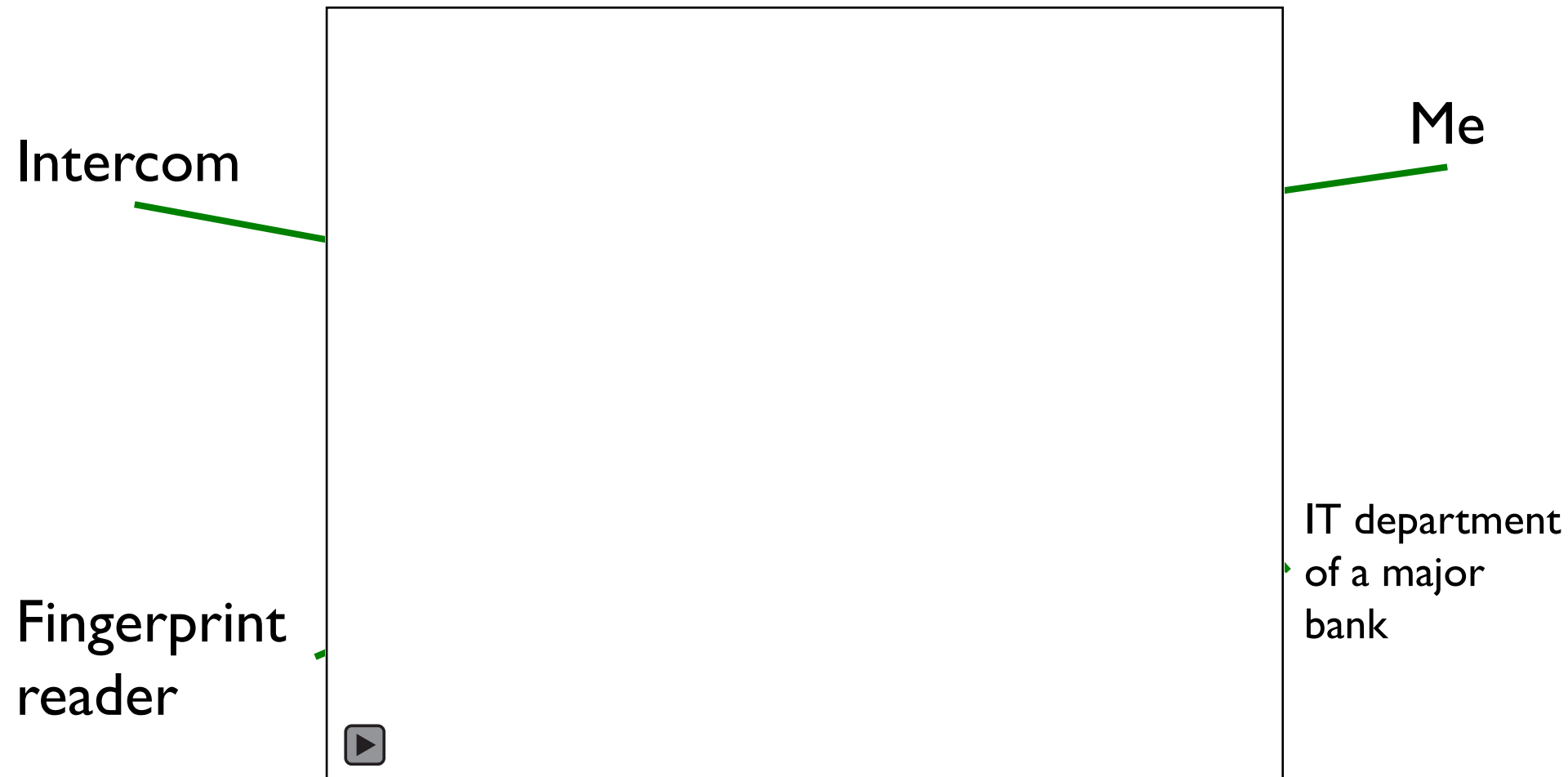
**Hemifacial Microsomia**
(Zielinski,.., & Erlich)
PLoS One

# Vulnerability research

Intercom

Me

Fingerprint
reader

IT department
of a major
bank

# Correlation between Y-chr and surnames



**www.ysearch.org:**

ACACACAC…

Erlich

# The main idea

A systematic study: can we recover the identity of <span style="color:red">anonymous</span> genomic datasets?

# Databases of interest

<span style="color:red">140,000 **publicly** accessible surname-Ychr records</span>

**www.smgf.org**                    **www.ysearch.org**

# How to find surnames?

Estimating the **time** to most recent common ancestor



$t_i$

$t_i$

surname

Target    i-th record in db

# Empirical test to determine the probability of recovering a US surname



Y-chr of a real person → Querying Ysearch and SMGF → Surname inference algorithm → **Inferring surname**

**Comparing the predicted surname to the true one**

**x900**

**For US Caucasian males:**

**12% Successful recoveries**

5% Wrong recoveries

83% Unknown

# Distribution of inferred surnames



Most of the inferred surnames are relatively rare

# Triangulate individuals with metadata

**Age**

**State**

**Surname**



**100,000 rounds**

**Adams**

■ Age+state+surname

■ Only age+state

The median of age+state+surname is 12 males.

# Putting it all together: the Venter case

lobSTR

We got a surname from whole genome sequencing data

**Method**

**lobSTR: A short tandem repeat profiler for personal genomes**

Melissa Gymrek,[1,2] David Golan,[3] Saharon Rosset,[3] and Yaniv Erlich[2,4]

[1]Harvard–MIT Division of Health Sciences and Technology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA; [2]Whitehead Institute for Biomedical Research, Cambridge, Massachusetts 02142, USA; [3]Department of Statistics and Operations Research, Tel Aviv University, Tel Aviv 69978, Israel

DYS 458
17

**Try it yourself: bit.ly/find_craig**

# Getting to Craig Venter

Searching for:

1.  Venter
2.  California
3.  Born in 1946
4.  Male

In USSearch.com

**Two matches, including:**

# Can we identify anonymous genomes?

# 1000 Genomes cases



**10 CEU (Utah) genomes**

Google    Winfield Utah

**Found an obituary that has the exact description of the pedigree**

surname predictions

**Probability of a random match < 5x10⁻⁹**

*Some of the details in this slide were modified to respect the identity of the family

# Beginner's luck?



p<5x10$^{-9}$

p<5x10$^{-6}$

p<10$^{-5}$

■ Successful surname recovery (targeted individual)

↗ Person tested by genetic genealogy service (source)

— Patrilineal line from source to target

Breaching the privacy of close to **50** CEU samples**.**

# Aftermath

## Our study



**Identifying Personal Genomes by Surname Inference**

Melissa Gymrek,[1,2,3,4] Amy L. McGuire,[5] David Golan,[6] Eran Halperin,[7,8,9] Yaniv Erlich[1]*

Sharing sequencing data sets without identifiers has become a common practice in genomics. Here, we report that surnames can be recovered from personal genomes by profiling short tandem



"..And Erlich respon... ...emplary way

findin... ...d other genetic...

resea...

an **important prece...** ...dealing with

newly disco... ...y... other researchers ...response

should ...

# The hitchhiker guide to genome hacking

**nature REVIEWS GENETICS**

REVIEWS

# Routes for breaching and protecting genetic privacy

*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.

# The hitchhiker guide to genome hacking

## LETTERS

### On Jim Watson's *APOE* status: genetic information is hard to hide

*European Journal of Human Genetics* (2009) **17**, 147–149; doi:10.1038/ejhg.2008.198; published online 22 October 2008

Genetic imputation

Alzheimer's disease

Barack Obama is the President

# The path forward

OPEN ACCESS Freely available online

**PLOS | BIOLOGY**

**Perspective**

## Redefining Genomic Privacy: Trust and Empowerment

Yaniv Erlich[1]*, James B. Williams[2], David Glazer[2], Kenneth Yocum[3], Nita Farahany[4], Maynard Olson[5], Arvind Narayanan[6], Lincoln D. Stein[7,8], Jan A. Witkowski[9], Robert C. Kain[3]

1 Whitehead Institute for Biomedical Research, Nine Cambridge Center, Cambridge, Massachusetts, United States of America, 2 Google Inc., Mountain View, California, United States of America, 3 Illumina Inc., San Diego, California, United States of America, 4 Duke University School of Law, Duke Science & Society, Durham, North Carolina, United States of America, 5 University of Washington, Port Orford, Oregon, United States of America, 6 Department of Computer Science, Princeton University, Princeton, New Jersey, United States of America, 7 Ontario Institute for Cancer Research, Toronto, Ontario, Canada, 8 Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada, 9 Banbury Center, Cold Spring Harbor Laboratory, Huntington, New York, United States of America

Key points: transparency, reputation system, compensation

# Acknowledgements

**Team Genetic Privacy**

Melissa Gymrek (HST – Harvard/MIT)

Amy McGuire (Baylor)

David Golan (Tel-Aviv University)

Eran Halperin (Tel-Aviv University)

## Funding:

WHITEHEAD INSTITUTE
NINE CAMBRIDGE CENTER