

COM 402

Security and Privacy

Data Protection for Personalized Health

Prof. Jean-Pierre Hubaux

Head of LCA1 Lab

Academic Director of the Center for Digital Trust
School of Computer and Communication Sciences
EPFL

With gratitude to the biomedical and CS researchers I have the privilege to work with



“...and if anyone here suspects that the algorithm that put these two together might be flawed, speak now...”

Source: The New Yorker

GESUNDHEIT

Rückenwind für den E-Patienten: Parlament will obligatorische elektronische Patientendossiers

von Michel Burtscher - az Aargauer Zeitung • 18.2.2019 um 04:15 Uhr



Der gläserne Patient: Nun nimmt die Politik neuen Anlauf für elektronische Gesundheitsdossiers.

Zur Verfügung gestellt

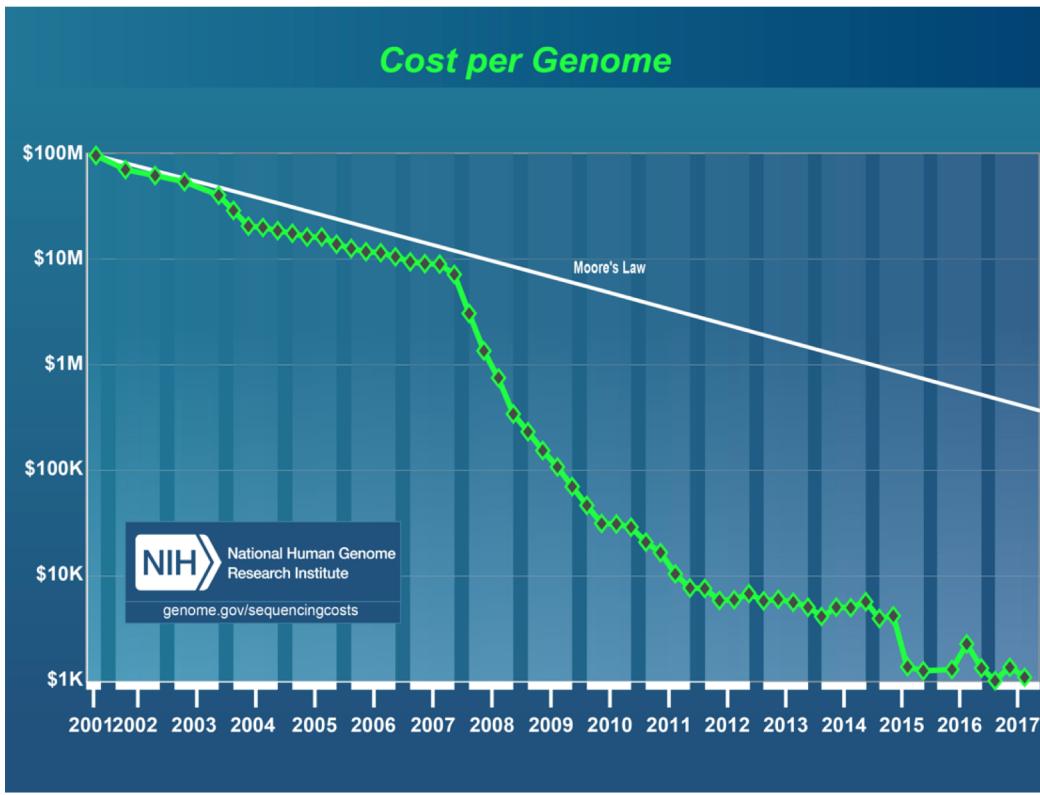
Damit sich das elektronische Patientendossier etablieren kann, müssen möglichst viele Ärzte, Spitäler und Kliniken mitmachen. Doch für Praxisärzte ist die Teilnahme bisher freiwillig. Das könnte sich bald ändern – zumindest für einen Teil von ihnen.

- Electronic Health Records are mandated by law in Switzerland
- Mandatory for hospitals as of 2020
- 2 systems will be deployed (by Swisscom and Swiss Post)

The Genomic Avalanche Is Coming...



Personalized Health



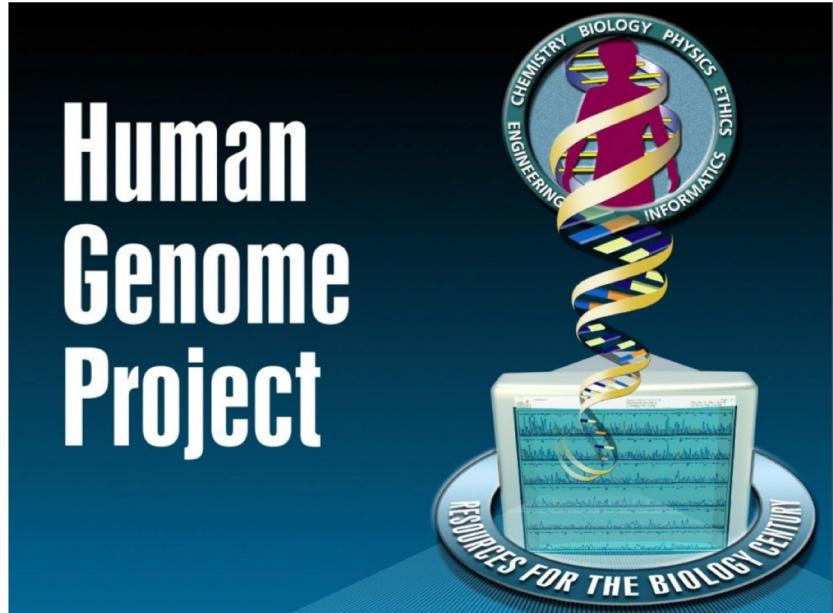
The massive digitalization of clinical and genomic information is providing unprecedented opportunities for improvements in diagnosis, preventive medicine and targeted therapies

Medical Use of Genetics

- Genetic disease risk tests help early diagnosis of serious diseases
- Pharmacogenomics → personalized medicine



Human Genome Project



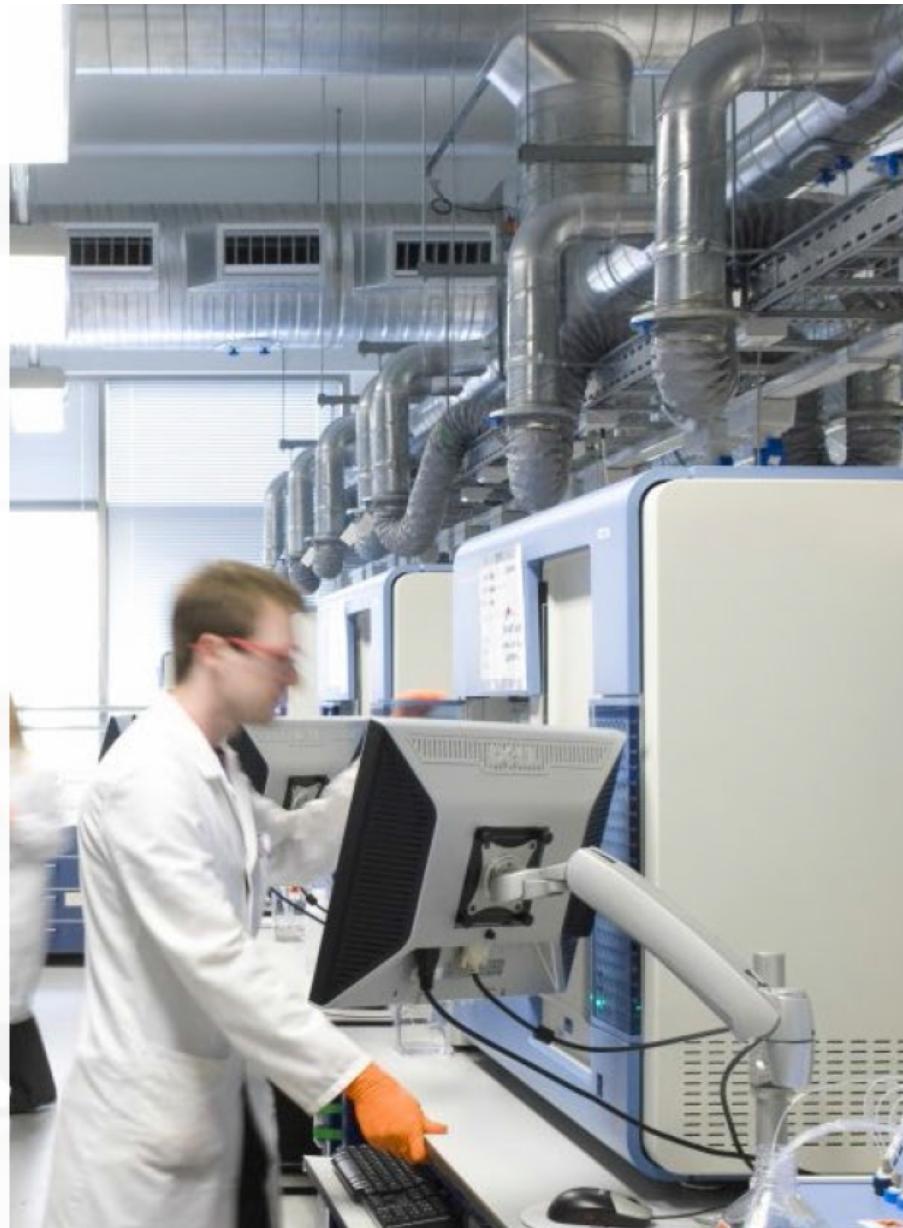
- 1990 – 2003
- Goals
 - sequencing and identifying all three billion chemical units in the human genetic instruction set,
 - Finding the genetic roots of disease
 - Developing treatments
- Mostly US/UK effort

NEXT GENERATION SEQUENCING

NEXT GENERATION SEQUENCING

**TECHNOLOGY ABOUT TO PROVIDE
AFFORDABLE (IN TIME AND COST)
AND PRECISE DETERMINATIONS OF
GENOME WIDE:**

- DNA sequence/variations (DNA-seq)
- Gene subregions' activity (RNA-seq)
- Protein-DNA interaction regions (ChIP-seq)



The Genomic Era

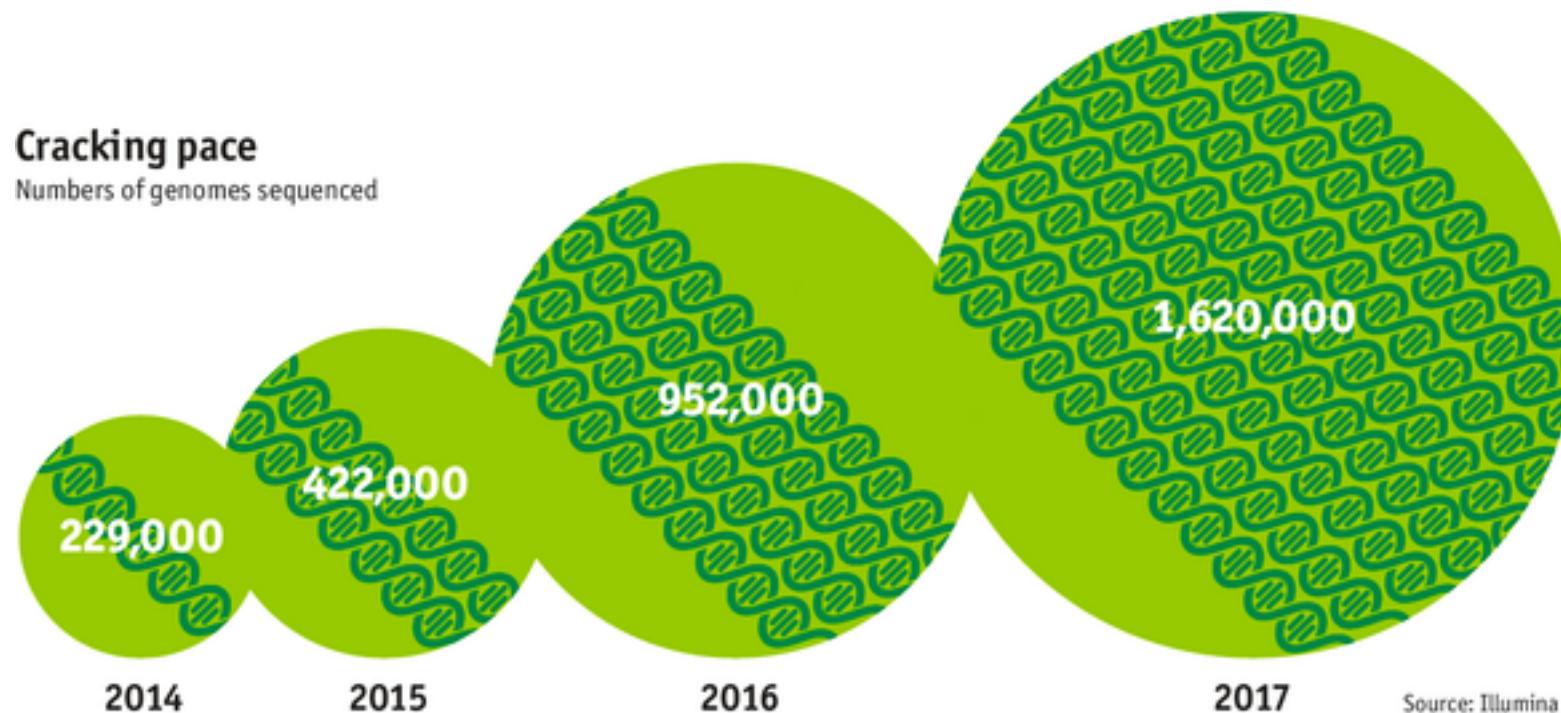
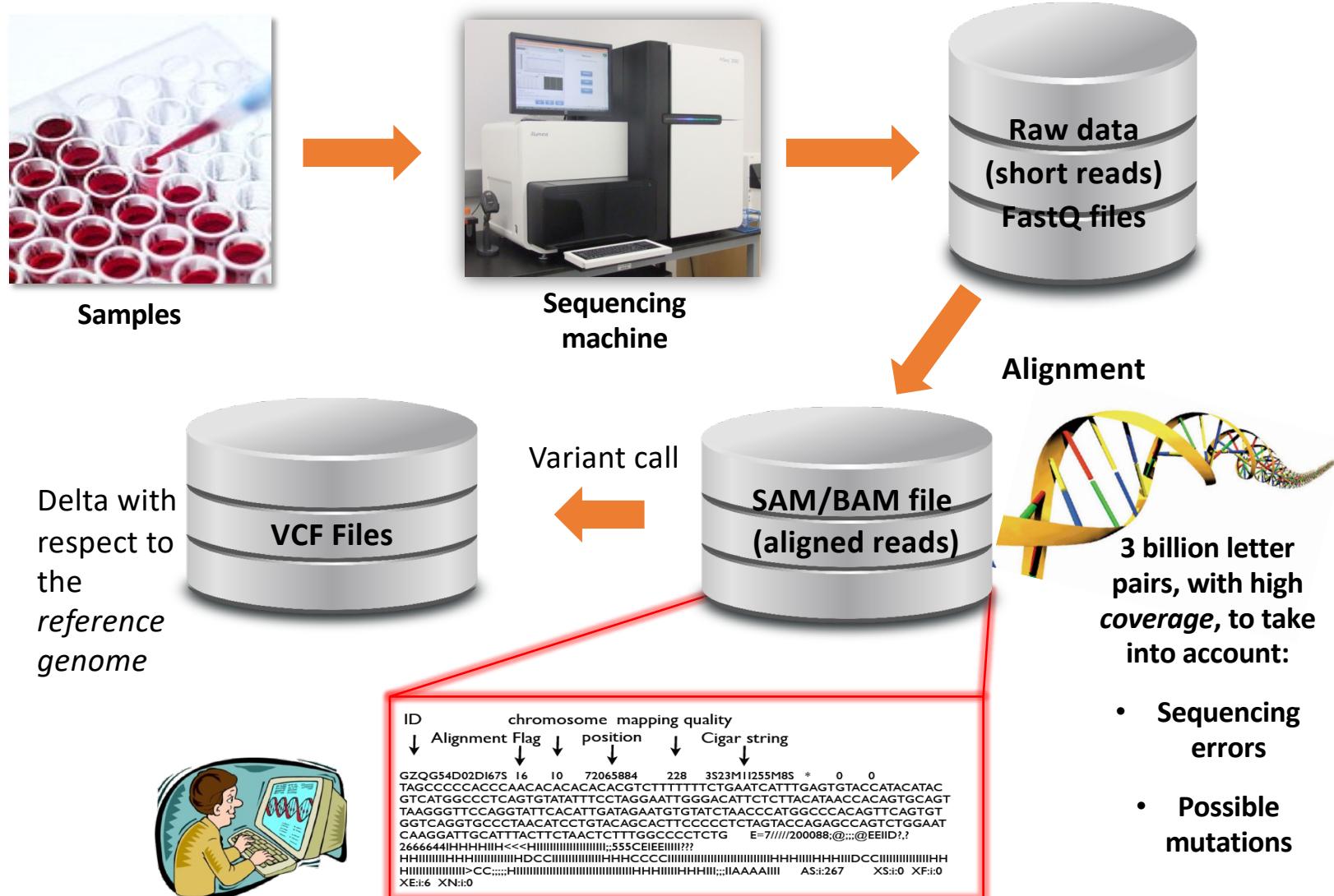


Figure from The Economist

From Blood Sample to Genome Analysis





IEEE
SPECTRUM

FOR THE TECHNOLOGY INSIDER. 01.15

A COOKE'S GIFT
TO HIS SON
A SWIMMER TACKLES
A Diabetes GBC
P. 21

WEARABLES WORN
BY PRO ATHLETES
The tech that gives
top stars an edge
P. 44

THE DOG IN
YOUR POCKET
Inside the race to
make competition
P. 46

WHY THE EBOLA
MODELS FAILED
Big gaps in info
couldn't be fixed
P. 62



IEEE

HACKING THE HUMAN **OS**

How Big Data
Will Transform
Medicine
and Health

SPECIAL REPORT

MIT Technology Review

VOL. 118 NO. 3 MAY/JUNE 2015 \$6.99

Feature p. 48
HP Tries to Reinvent
the Computer

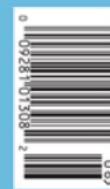
Business Report p. 63
Persuasion

Review p. 72
The Problem with
Fake Meat



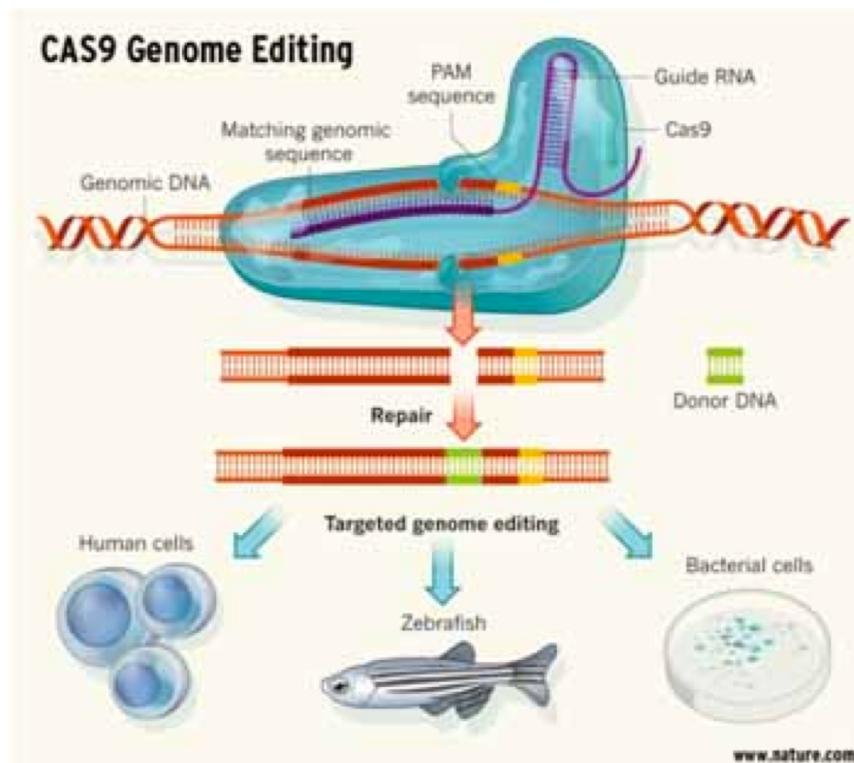
WE CAN
NOW
ENGINEER
THE
HUMAN
RACE

p26



Genome Editing (CRISPR-CAS9)

- Potential to alter the **human** genome
- Strong potential for treatment of (human) genetic diseases
- Moratorium pronounced in December 2015 for edition of inheritable parts of the human genome
- Moratorium not fully respected...



CRISPR: Clustered regularly interspaced short palindromic repeats
CAS9 is a protein

China's CRISPR twins might have had their brains inadvertently enhanced



CRISPR: Gene-editing technique

- Deletion of gene CCR5 was intended to make babies HIV-resistant
- This deletion could increase cognition capabilities

Biotechnology May 6, 2019

...

Google backs a bid to use CRISPR to prevent heart disease

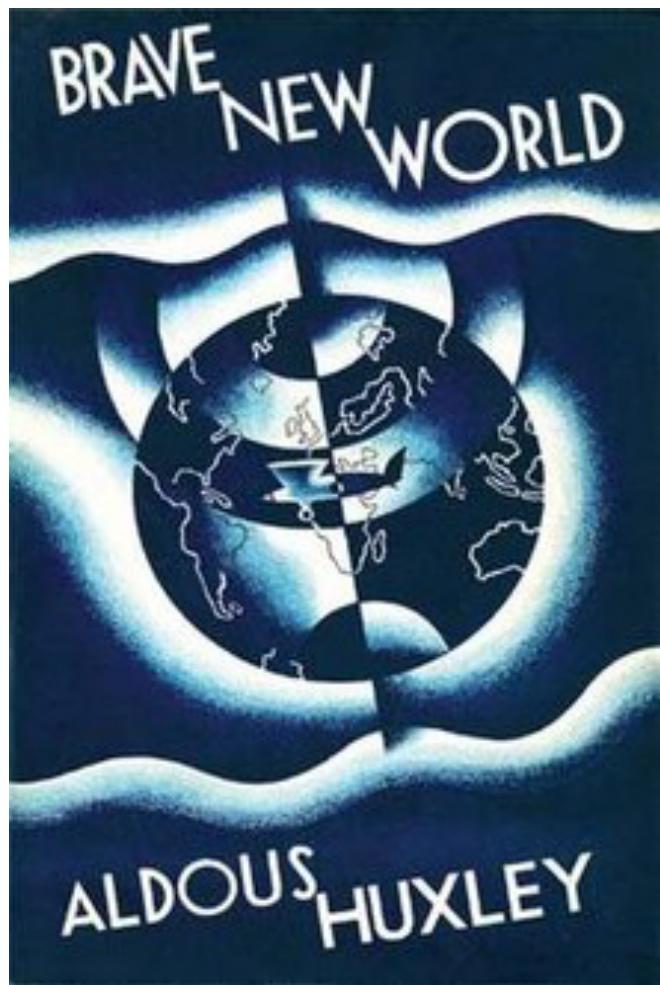


Ever wonder why some fortunate people eat chips, don't exercise, and still don't get clogged arteries? It could be because they've got lucky genes.

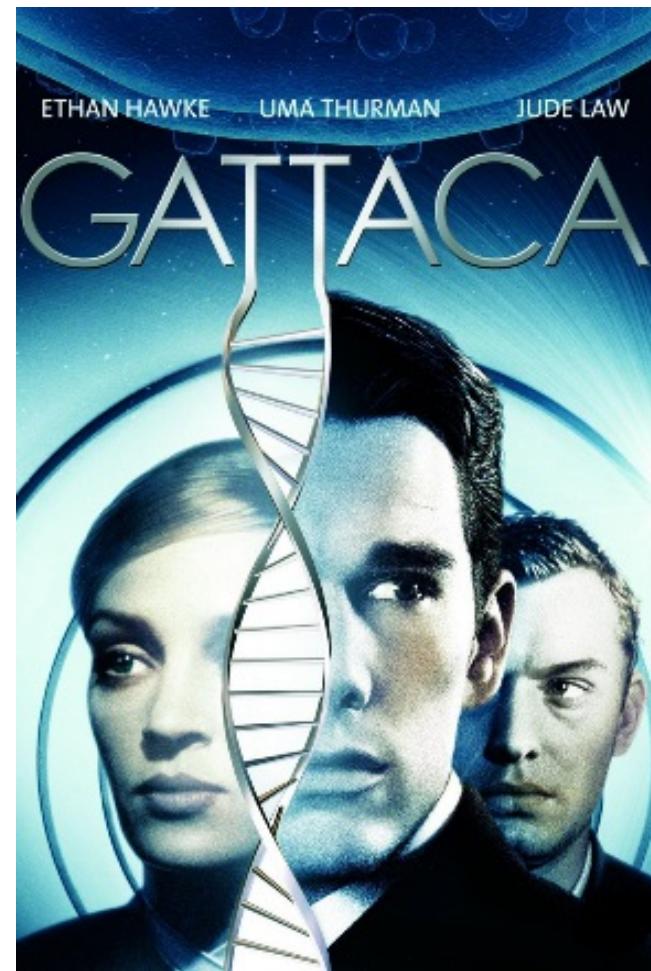
Now Alphabet (Google's parent company) is bankrolling a startup company that plans to use gene editing to spread fortunate DNA variations with "one-time" injections of the gene-editing tool CRISPR.

Heart doctors involved say the DNA-tweaking injections could "confer lifelong protection" against heart disease.

Big idea: The startup, Verve Therapeutics, said it had raised \$58.5 million from investors including Alphabet's venture fund, GV. What makes Verve different? Most gene-therapy companies have gone after rare diseases like hemophilia. But Verve thinks editing people's DNA could instead help solve the most common cause of death.



Novel, 1931



Movie, 1997

Medical Data Breach in Switzerland

The screenshot shows a news article from the Swiss newspaper 'Le Matin'. The header includes the logo 'Le Matin' and a navigation bar with links to MONDE, SUISSE, SPORTS, FAITS DIVERS, PEOPLE, LOISIRS, SOCIÉTÉ, HIGH-TECH, and Images. The main headline reads 'Bilans de santé en balade sur le net'. Below it is a sub-headline: 'GAFFE — Des données médicales ultraconfidentielles de patients romands ont été librement accessibles durant des jours sur Internet. Le groupe Synlab déplore une erreur humaine.' A byline indicates the article was written by Raphaël Pomey and last updated on April 8, 2015. The image below the text shows a laboratory technician in a white coat working at a computer terminal with multiple monitors displaying medical data. A caption at the bottom left of the image states: 'Tests du papillomavirus, dépistages du sida ou d'une hépatite ont circulé pendant des jours sur le Web.' and credits 'Image: Jason Butcher / Corbis / Montage'.

- 8300 files of medical analyses freely available online, for days
 - With full patient identifications
 - Including HIV and other tests
 - Sheer carelessness (no attack)
- (April 2015)

Ransomware Attack against German Hospitals

20. März 2016, 10:05 Uhr Klinikum Neuss

Wenn Cyberkriminelle ein Krankenhaus lahmlegen



Das Lukaskrankenhaus der Städtischen Kliniken in Neuss wurde Opfer von Cyberkriminellen. (Foto: dpa)

“WannaCry” Ransomware Virus (May 2017)

Do state institutions have the resources to fight hackers?

Public sector has lessons to learn as hospital trusts and GPs struggle to recover from ransomware attack



① A ransomware attack bought computers to a standstill across the world on Friday. Photograph: Ritchie B. Tongo/EPA

The Guardian,
14 May 2017

US Healthcare Official “Wall of Shame”

https://ocrportal.hhs.gov/ocr/breach/breach_report.jsf

Around 5 declared breaches per week, each affecting 500+ people



U.S. Department of Health and Human Services
Office for Civil Rights
Breach Portal: Notice to the Secretary of HHS Breach of Unsecured Protected Health Information

Welcome | File a Breach | HHS | Office for Civil Rights | Contact Us

Under Investigation Archive Help for Consumers

As required by section 13402(e)(4) of the HITECH Act, the Secretary must post a list of breaches of unsecured protected health information affecting 500 or more individuals. The following breaches have been reported to the Secretary:

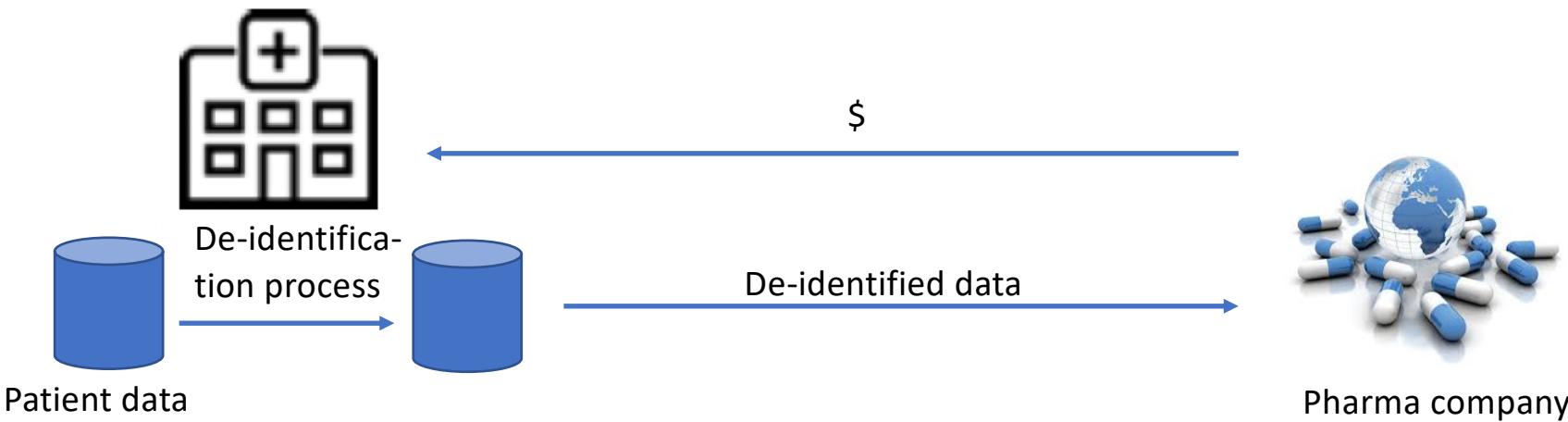
Cases Currently Under Investigation

This page lists all breaches reported within the last 24 months that are currently under investigation by the Office for Civil Rights.

[Show Advanced Options](#)

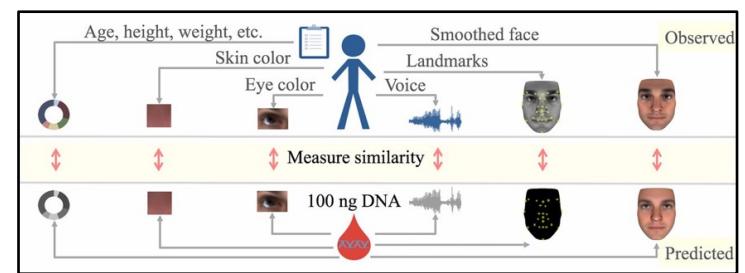
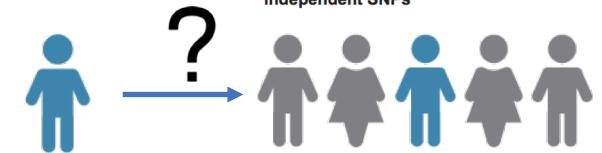
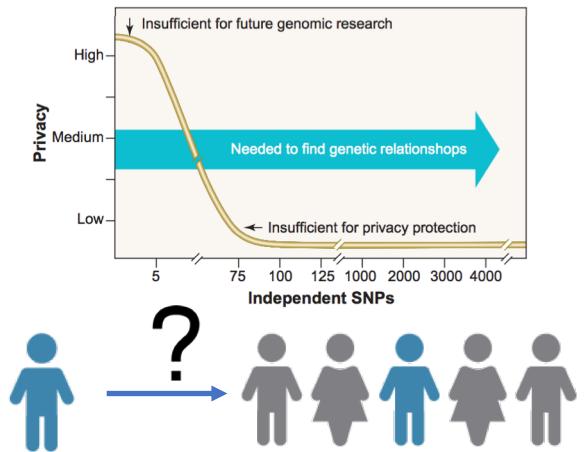
Breach Report Results							
Expand All	Name of Covered Entity	State	Covered Entity Type	Individuals Affected	Breach Submission Date	Type of Breach	Location of Breached Information
1	Ohio Living	OH	Healthcare Provider	6510	09/07/2018	Hacking/IT Incident	Email
1	Rockdale Blackhawk, LLC d/b/a Little River Healthcare	TX	Healthcare Provider	1494	09/07/2018	Unauthorized Access/Disclosure	Electronic Medical Record, Other
1	J.A. Stokes Ltd.	NV	Healthcare Provider	3200	09/05/2018	Hacking/IT Incident	Desktop Computer, Electronic Medical Record, Network Server
1	Reliable Respiratory	MA	Healthcare Provider	21311	09/01/2018	Hacking/IT Incident	Email
1	Port City Operating Company doing business as St. Joseph's Medical Center	CA	Healthcare Provider	4984	08/31/2018	Loss	Other Portable Electronic Device
1	Carpenters Benefit Funds of Philadelphia	PA	Health Plan	20015	08/31/2018	Hacking/IT Incident	Email

Hospital / pharma interface



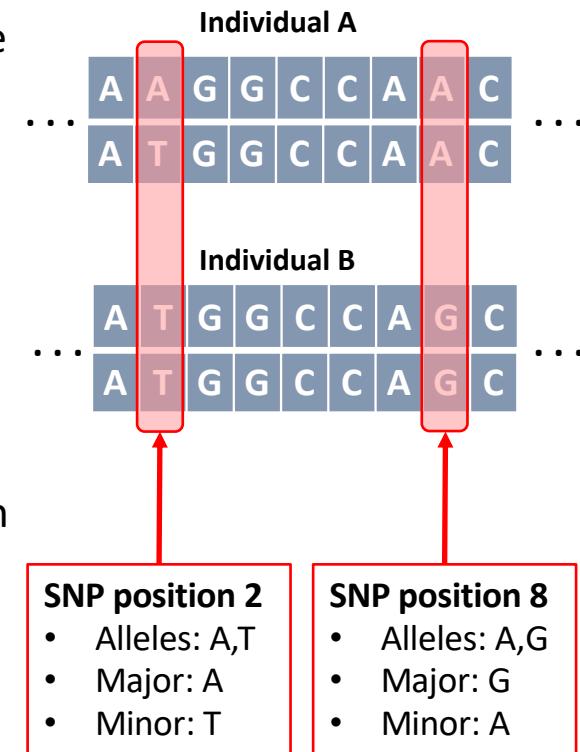
De-identification of genomic data is impossible

- **Lin et al. 2004 *Science***: 75 or more SNPs (Single Nucleotide Polymorphisms) are sufficient to identify a single person
- **Homer et al. 2008 *PLOS Genetics***: aggregated genomic data (i.e., allele frequencies) can be used for re-identifying an individual in a case group with a certain disease
- **Gymrek et al. 2013 *Science***: surnames can be recovered from personal genomes, linking “anonymous” genomes and public genetic genealogy databases
- **Lipper et al. 2017 *PNAS***: Anonymous genomes can also be identified by inferring physical traits and demographic information
- **Many more to come...**



Most common genetic variation: Single Nucleotide Polymorphism (SNP)

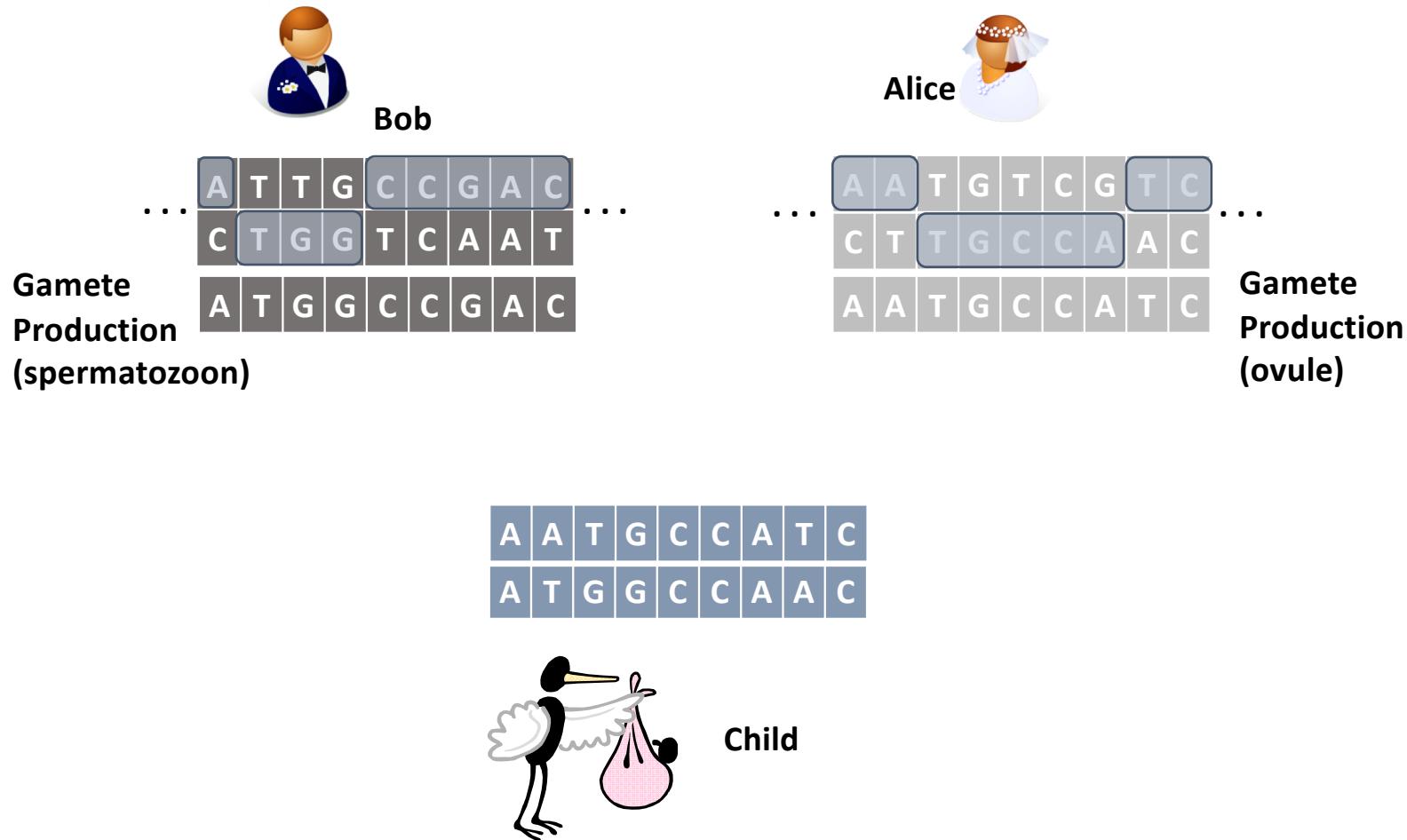
- Occurs when, at a specific position, at least a single nucleotide (A,C,G, or T) differs between members of the same species in more than 1% of the population
- Potential nucleotides for a SNP are called **alleles**
- 2 different alleles can be observed for each SNP:
 - **Major** allele (M)
 - **Minor** allele (m)
- Every genome carries 2 alleles at each SNP position
A SNP can be either:
 - **Homozygous minor** [m,m]
 - **Heterozygous** [m,M] or [M,m]
 - **Homozygous major** [M,M]



Alice and Bob: The Long Awaited Happy End

After having extensively authenticated each other,
after having exchanged thousands of highly private
messages,
after having established numerous secure channels
between each other,
after years of intense but platonic relationship, finally,
finally... ❤️

... Alice and Bob got closer to each other



Direct to Consumer Genomics (1/2)

- Ancestry.com (millions of customers)

AncestryDNA—The World's Largest Consumer DNA Database.
Get started in a few simple steps.

Order your complete kit with easy-to-follow instructions.

Return a small saliva sample in the prepaid envelope.

Your DNA will be analyzed at more than 700,000 genetic markers.

Within 6-8 weeks, expect an email with a link to your online results.

Uncover your ethnic mix.

When your results arrive, you'll see a breakdown of your ethnicity—and it may contain a few surprises. Then, you can start learning more about the places where your family story began.

See all 26 ethnic regions covered by the AncestryDNA test.

Find relatives you never knew you had.

Once you've taken your test, we'll search our network of AncestryDNA members and identify your cousins—the people who share your DNA. And if you're lucky, you might even make a **New Ancestor Discovery™***

*Some features may require an Ancestry subscription.

Direct to Consumer Genomics (2/2)

- 23andMe.com
(millions customers)



Name	Confidence	Your Risk	Avg. Risk
Atrial Fibrillation	★★★★★	33.9%	27.2%
Prostate Cancer ♂	★★★★★	29.3%	17.8%
Alzheimer's Disease	★★★★★	14.2%	7.2%
Age-related Macular Degeneration	★★★★★	11.1%	6.5%
Colorectal Cancer	★★★★★	7.8%	5.6%
Chronic Kidney Disease	★★★★★	4.2%	3.4%
Restless Legs Syndrome	★★★★★	2.5%	2.0%
Parkinson's Disease	★★★★★	2.2%	1.6%



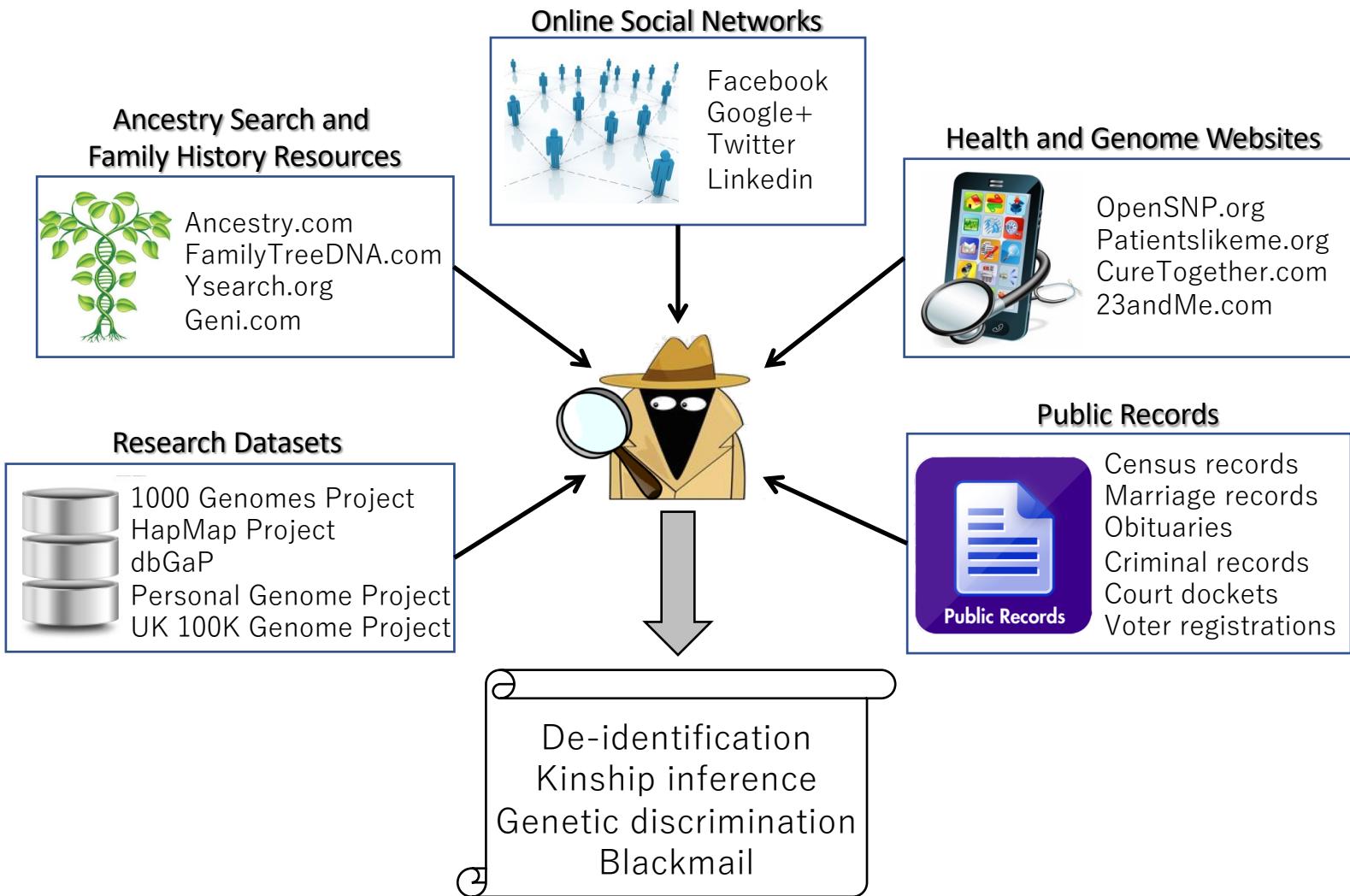
With genetic testing, I gave my parents the gift of divorce

Updated by George Doe on September 9, 2014, 7:50 a.m. ET

 TWEET (2,073)

 SHARE (15K)

+



Re-identification attacks on genomic data have a long history

The image shows a rectangular box containing a scientific abstract. At the top left is the text "OPEN ACCESS Freely available online". At the top right is the "PLOS GENETICS" logo. The main title of the article is "Resolving Individuals Contributing Trace Amounts of DNA to Highly Complex Mixtures Using High-Density SNP Genotyping". Below the title is a short abstract: "Individuals contributing trace amounts of DNA to highly complex mixtures can be resolved using high-density SNP genotyping. This approach was used to resolve individuals contributing to a sample from a forensic mixture of DNA from 100 individuals. The results show that 10,000–50,000 SNPs are sufficient to determine if an individual was part of a cohort, even when he contributed < 0.1% of the data. The findings also suggest that composite statistics across cohorts, such as allele frequency or genotype counts, do not mask identity within genome-wide association studies. The implications of these findings are discussed." A large callout box is overlaid on the bottom half of the abstract, containing the text: "10,000 – 50,000 SNPs are sufficient to determine if an individual was part of a cohort, even when he contributed < 0.1% of the data".

10,000 – 50,000 SNPs are sufficient to determine if an individual was part of a cohort, even when he contributed < 0.1% of the data

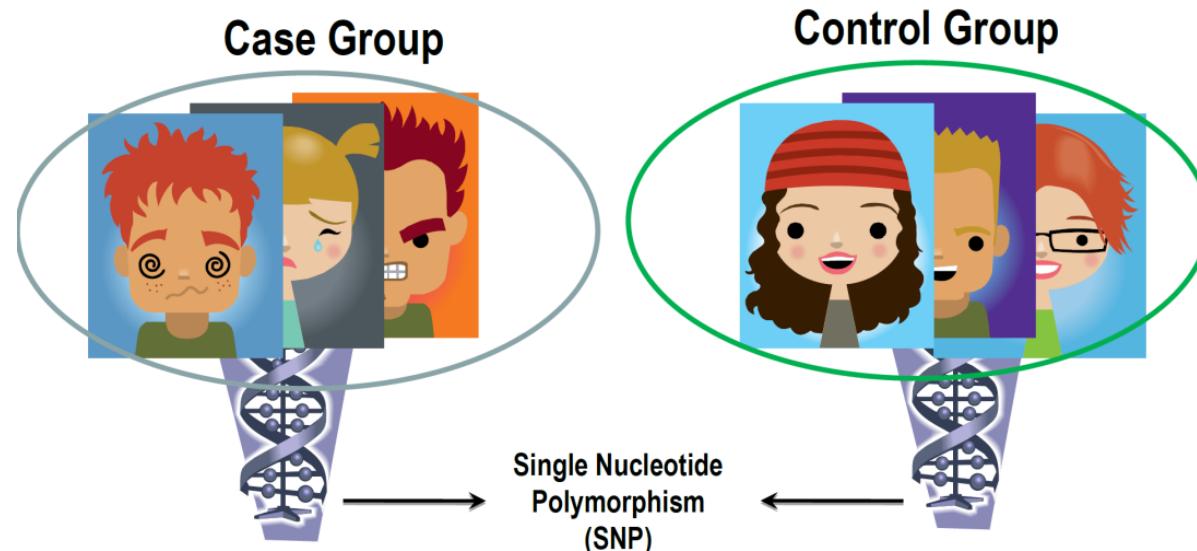
sources due to sample contamination. These findings also suggest that composite statistics across cohorts, such as allele frequency or genotype counts, do not mask identity within genome-wide association studies. The implications of these findings are discussed.

Many other subsequent studies extended the range of vulnerabilities for summary statistics:

[Jacobs et al. *Nature Genet.* '09], [Vissecher and Hill *PLoS Genet.* '09], [Sankararaman et al. *Nature Genet.* '09], [Wang et al. *CCS*'09], [Clayton *Biostatistics* '10], [Im et al. *Am. J. Hum. Genet.* '12], ...

Homer's Attack

- Adversary has access to a known participant's genome
- Goal: determine if the target individual is in the case group
- Uses simple correlation in the genome (linkage disequilibrium)
- Attack later improved by Wang et al.



N. Homer, S. Szelinger, M. Redman, D. Duggan, and W. Tembe. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. PLoS Genetics, 4, Aug. 2008.

Homer's Attack

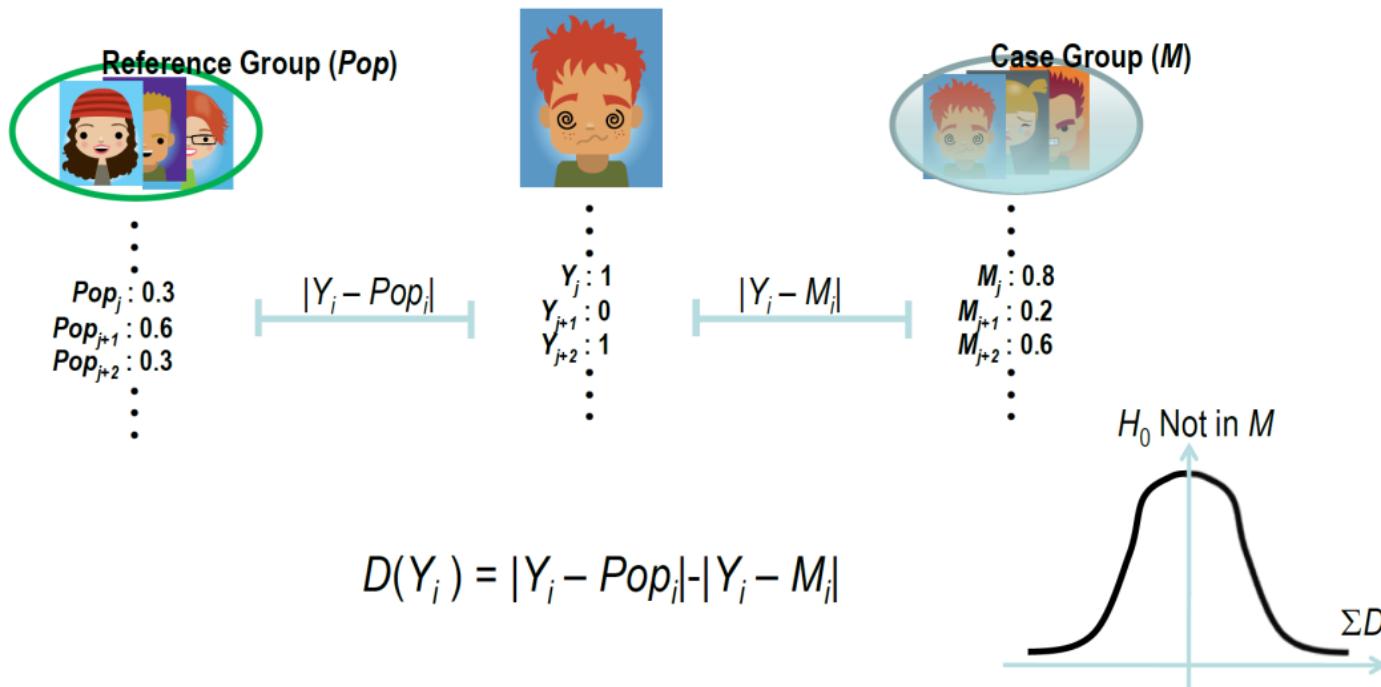


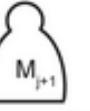
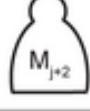
Figure from: Wang et al. Learning Your Identity and Disease from Research Papers: Information Leaks in Genome-Wide Association Study

N. Homer, S. Szelinger, M. Redman, D. Duggan, and W. Tembe. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. PLoS Genetics, 4, Aug. 2008.

Homer's attack in a nutshell

The attacker knows:

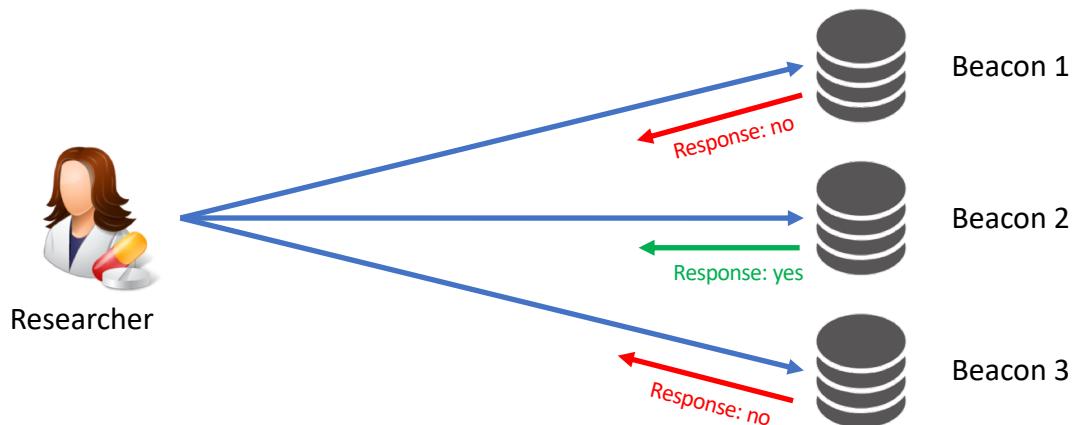
- The genome of the victim (her set of variants)
- The size of the mixture he's attacking
- Population allele frequencies

Snp	Allele Frequency ($Y_{i,j}$)	Distance Measure	Interpretation at the given SNP
	0.0 0.25 0.50 0.75 1.0	$D(Y_{i,j}) = Y_{i,j} - Pop_j - Y_{i,j} - M_j $	
j	  	= 1.0 - 0.25 - 1.0 - 0.75 = 0.75 - 0.25 = 0.50	most likely to be in the Mixture
j+1	  	= 0.50 - 0.25 - 0.50 - 0.75 = 0.25 - 0.25 = 0.00	equally likely to be in the Mixture and in the Reference Population
j+2	  	= 0.00 - 0.25 - 0.00 - 0.75 = 0.25 - 0.75 = - 0.50	most likely to be in the Reference Population

Person of Interest (Y_i) Reference Population (Pop) Mixture (M)

Figure taken from: Homer N, Szelinger S, Redman M, Duggan D, Tembe W, et al. (2008) Resolving Individuals Contributing Trace Amounts of DNA to Highly Complex Mixtures Using High-Density SNP Genotyping Microarrays. PLoS Genet 4(8): e1000167. doi:10.1371/journal.pgen.1000167

GA4GH Beacon Project



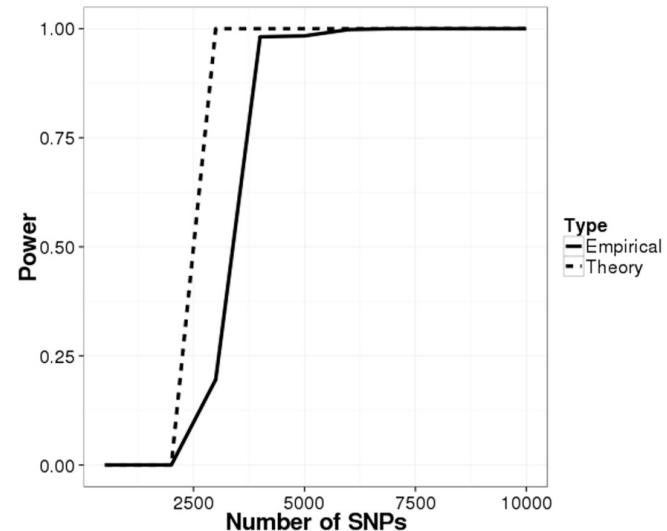
Main features:

- Allows researchers to quickly query multiple database to find the sample they need
- Encourages cross-borders collaboration among researchers
- Only provides minimal responses back in order to mitigate privacy concerns

Beacon used as an oracle: the SB attack



Shringarpure SS, Bustamante CD. Privacy risks from genomic data-sharing beacons. The American Journal of Human Genetics. 2015 Nov 5;97(5):631-46.



- The attack relies on the assumption that the adversary knows:
 - The set of variants (VCF file) of the target individual
 - The size of the beacon
- The attack is based on a likelihood ratio test where the adversary repeatedly queries the beacon in order to re-identify the individual
- The attack can be **extremely dangerous** if the beacon is associated with a sensitive phenotype (e.g., cancer)

SB Attack: Likelihood Ratio Test

- H₀: the target individual is not in the beacon
- H₁: the target individual is in the beacon
- $R = \{x_1, \dots, x_n\}$ is the set of beacon responses
- δ probability of sequencing errors
- D_N^i denote the probability that none of the N other genomes in the beacon have an alternate allele at position i

$$L_{H_0}(R) = \sum_{i=1}^n x_i \log(1 - D_N^i) + (1 - x_i) \log(D_N^i),$$

$$L_{H_1}(R) = \sum_{i=1}^n x_i \log(1 - \delta D_{N-1}^i) + (1 - x_i) \log(\delta D_{N-1}^i).$$

$$\Lambda = L_{H_0}(R) - L_{H_1}(R)$$

“Optimal” re-identification attack with real allele frequencies

- **Smarter adversary** who makes use of publicly available information
- **Idea:** rare alleles have higher re-identification power !!



➔ The attacker queries rarest alleles first $f_1 \leq f_2 \leq \dots \leq f_M$ ↵

$$D_N^i = \Pr(\text{none of the other } N \text{ genomes have an alternate allele at position } i)$$

Original attack [SB15]

$$D_N^i = \mathbb{E}[(1 - f_i)^{2N}],$$

Where:

- $f_i \sim \text{beta}(a, b)$

“Optimal” attack [R et al. 16]

$$D_N^i = (1 - f_i)^{2N}$$

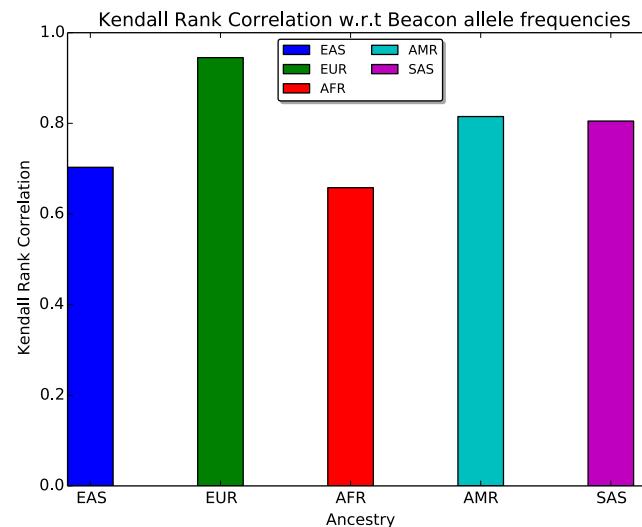
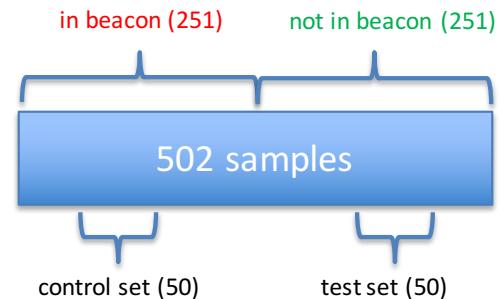
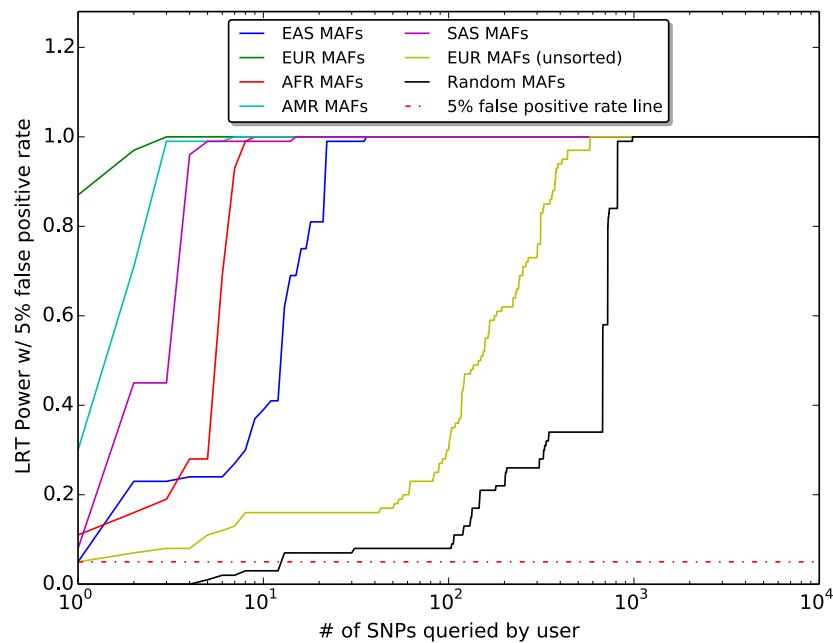
Where:

- f_i is different at every position i

“Optimal” re-identification attack in single population beacon

- **Experimental setup:**

- 502 EUR samples from 1000 Genomes Project Phase 3 (chromosome 10)
- False positive rate of 5%



Kin Genomic Privacy

He

