

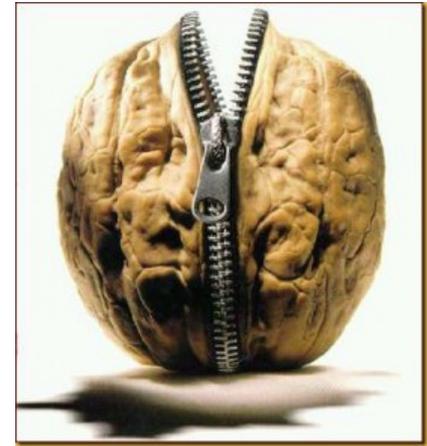
The background is a composite image. On the left, there is a semi-transparent image of a man wearing a dark hat and a dark jacket with a white stripe on the sleeve. On the right, there is a semi-transparent image of a DNA double helix structure. The text 'THE GOOD BAD UGLY' is repeated in a large, bold, grey font across the top and sides of the image.

The Genomics Revolution: The Good, The Bad, and The Ugly

(A Privacy Researcher's Perspective)

**Emiliano De Cristofaro
University College London
<https://emilianodc.com>**

This Talk In a...



The Good

Revolution in medicine and healthcare
Genetic testing for the masses

The Bad

Collection of highly sensitive data
Very hard to anonymize / de-identify

The Ugly

Greater good vs privacy
Encryption might not be the answer

History

1970s: DNA sequencing starts

1990: The “Human Genome Project” starts

2003: First human genome fully sequenced

2012: UK announces sequencing of 100K genomes

2015: USA announces sequencing of 1M genomes

\$\$\$

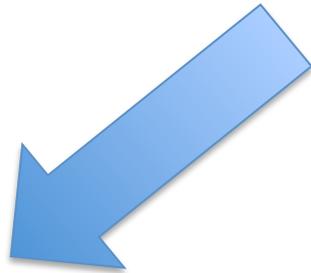
\$3B: Human Genome Project

\$250K: Illumina (2008)

\$5K: Complete Genomics (2009), Illumina (2011)

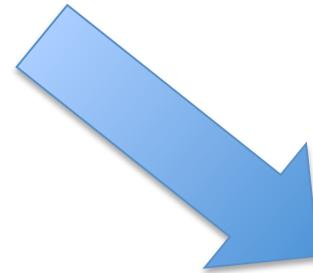
\$1K: Illumina (2014)

How to read the genome?



Genotyping

Testing for genetic differences using a set of markers



Sequencing

Determining the full nucleotide order of an organism's genome

1/05/2011 @ 4:57PM | 30,076 views

The First Child Saved By DNA Sequencing

+ Comment Now + Follow Comments



In Treatment for Leukemia, Glimpses of the Future



LETTER

doi:10.1038/nature13394

Genome sequencing identifies major causes of severe intellectual disability

Christian Gilissen^{1*}, Jayne Y. Hehir-Kwa^{1*}, Djie Tjwan Thung¹, Maartje van de Vorst¹, Bregje W. M. van Bon¹, Marjolein H. Willemsen¹, Michael Kwint¹, Irene M. Janssen¹, Alexander Hoischen¹, Annette Schenck¹, Richard Leach², Robert Klein², Rick Tearle², Tan Bo^{1,3}, Rolph Pfundt¹, Helger G. Yntema¹, Bert B. A. de Vries¹, Tjitske Kleefstra¹, Han G. Brunner^{1,4*}, Lisenka E. L. M. Vissers^{1*} & Joris A. Veltman^{1,4*}

TIME

THE ANGELINA EFFECT

Angelina Jolie's double mastectomy puts genetic testing in the spotlight. What her choice reveals about calculating risk, cost and peace of mind

BY JEFFREY KLUGER & ALICE PARK



Show results for

[See new and recently updated reports >](#)

23andMe Discoveries were made possible by 23andMe members who took surveys.

Disease Risks (114, 2 locked reports)

Elevated Risks	Your Risk	Average Risk
Psoriasis	22.4%	11.4%
Celiac Disease	0.5%	0.1%
Bipolar Disorder	0.2%	0.1%
Primary Biliary Cirrhosis	0.10%	0.08%
Scleroderma (Limited Cutaneous Type)	0.06%	0.07%

[See all 114 risk reports...](#)

Carrier Status (27, 1 locked report)

Hemochromatosis	Variant Present
Alpha-1 Antitrypsin Deficiency	Variant Absent
Bloom's Syndrome	Variant Absent
Canavan Disease	Variant Absent
Congenital Disorder of Glycosylation Type 1a (PMM2-CDG) 	Variant Absent
Cystic Fibrosis	Variant Absent
Familial Dysautonomia	Variant Absent
Factor XI Deficiency	Variant Absent

[See all 27 carrier status...](#)

Traits (52)

Alcohol Flush Reaction	Does Not Flush
Bitter Taste Perception	Can Taste
Earwax Type	Wet
Eye Color	Likely Blue
Hair Curl 	Slightly Curlier Hair on Average

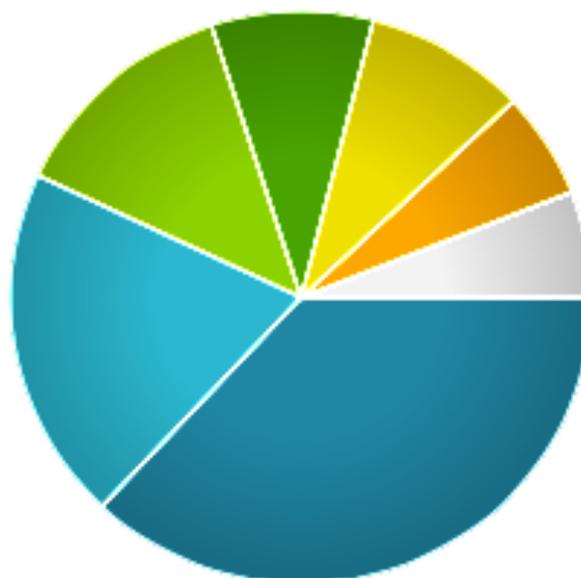
[See all 52 traits...](#)

Drug Response (20)

Warfarin (Coumadin®) Sensitivity	Increased
Abacavir Hypersensitivity	Typical
Alcohol Consumption, Smoking and Risk of Esophageal Cancer	Typical
Clopidogrel (Plavix®) Efficacy	Typical
Fluorouracil Toxicity	Typical

[See all 20 drug response...](#)

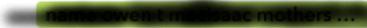
Genetic Ethnicity



	Southern European	37%
	West African	20%
	British Isles	13%
	Native South American	9%
	Finnish/Volga-Ural	9%
	Eastern European	6%
	Uncertain	6%

List View Map View Surname View

search matches Show: both sides Sort: relationship 25 per page 1 - 25 of 424

	Male	You		UPDATE YOUR PROFILE
	Female	2nd to 3rd Cousin 1.68% shared, 5 segments	J2a2	Send an Introduction
	Female	3rd to 4th Cousin 1.30% shared, 3 segments	United States Alsace-Lorraine (Strasbourg), Fr... Paternal  Senape 5 more U5b2	Public Match Send a Message
	Male	3rd to 4th Cousin 1.03% shared, 2 segments	H13a1a R1b1b2	Send an Introduction
	Female	3rd to 5th Cousin 0.45% shared, 2 segments	H7	Send an Introduction
	Female	3rd to 5th Cousin 0.42% shared, 2 segments	H1	Send an Introduction
	Male	3rd to 5th Cousin 0.40% shared, 2 segments	United States Reno, Nevada San Diego, California Tucker Littlefield Warga 4 more H1c G2a	Public Match Send a Message
	Male	3rd to 5th Cousin 0.37% shared, 2 segments	United States fathers father prince Edward isla...   K1a1b R1b1b2a1a	Public Match Send a Message
	Male, b. 1978	3rd to 6th Cousin 0.40% shared, 1 segment	United States New Jersey Utah California Northern Europe U3b1 T	Send an Introduction

Privacy Researcher's Perspective

Treasure trove of **sensitive** information

Ethnic heritage, predisposition to diseases

Genome = the ultimate **identifier**

Hard to anonymize / de-identify

Sensitivity is **perpetual**

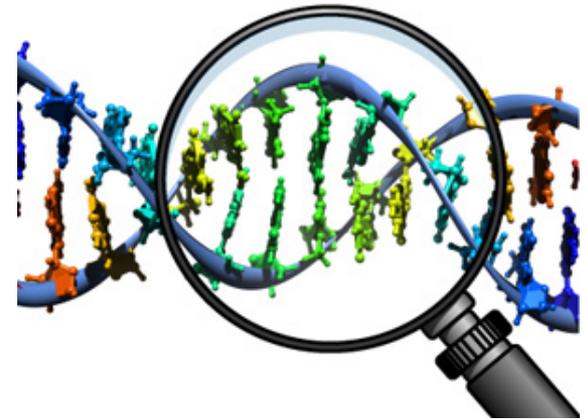
Cannot be “revoked”

Leaking one's genome \approx leaking relatives' genome

The Greater Good
vs
Privacy?

The rise of a new research community

Studying privacy issues



Exploring techniques to protect privacy



De-Anonymization

TECH 4/25/2013 @ 3:47PM | 17,111 views

Harvard Professor Re-Identifies Anonymous Volunteers In DNA Study

+ Comment Now + Follow Comments

A Harvard professor has re-identified the names of more than 40% of a sample of anonymous participants in a high-profile DNA study, highlighting the dangers that ever greater amounts of personal data available in the Internet era could unravel personal secrets.



Harvard Professor Latanya Sweeney

From the onset, the Personal Genome Project,

Melissa Gymrek et al. *"Identifying Personal Genomes by Surname Inference."* Science Vol. 339, No. 6117, 2013

Aggregation

Re-identification of aggregated data

Statistics from allele frequencies can be used to identify genetic trial participants [1]

Presence of an individual in a group can be determined by using allele frequencies and his DNA profile [2]

[1] R. Wang et al. "Learning Your Identity and Disease from Research Papers: Information Leaks in Genome Wide Association Study." CCS, 2009

[2] N. Homer et al. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays.

PLoS Genetics, 2008

Kin Privacy

Quantifying how much privacy do relatives lose when one's genome is leaked?



Also read: “Routes for breaching genetic privacy”
Y. Erlich and A. Narayanan,
Nature Review Genetics
Vol. 15, No. 6, 2014

M. Humbert et al., “Addressing the Concerns of the Lacks Family: Quantification of Kin Genomic Privacy.” Proceedings of ACM CCS, 2013

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

SHARE

+



Most Read

1

Read the Iranian foreign minister's passive aggressive response to Tom

2

Where the world's migrants go, in

3

Why there's a roaring controversy over Hillary Clinton's "homebrewed"

4

A new theory for why the bees are v

5

Human Aspects of Genome Privacy

Dynamic Consent:

Patients electronically control consent through time and receive information about the uses of their data

Jane Kaye's work at Oxford

Understanding “fears” and “reactions”, including:

Survivor's guilt

Freedom to withdraw is crucial but poorly understood

Insurance carriers and big corporations most distrusted

Ethnographic Studies in WGS

Semi-structured interviews with 16 participants

Assessing perception of genetic tests, attitude toward WGS programs, as well as perception of privacy/ethical issues

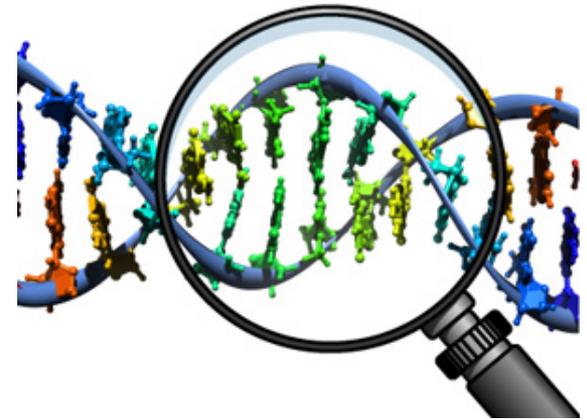
(Some) Preliminary results

1. Preferred method is through doctors not companies (trust)
2. Labor/healthcare discrimination top concerns
3. Differences in correlation with income and education

E. De Cristofaro. *“Users' Attitudes, Perception, and Concerns in the Era of Whole Genome Sequencing.”* (USEC 2014)

The rise of a new research community

Studying privacy issues

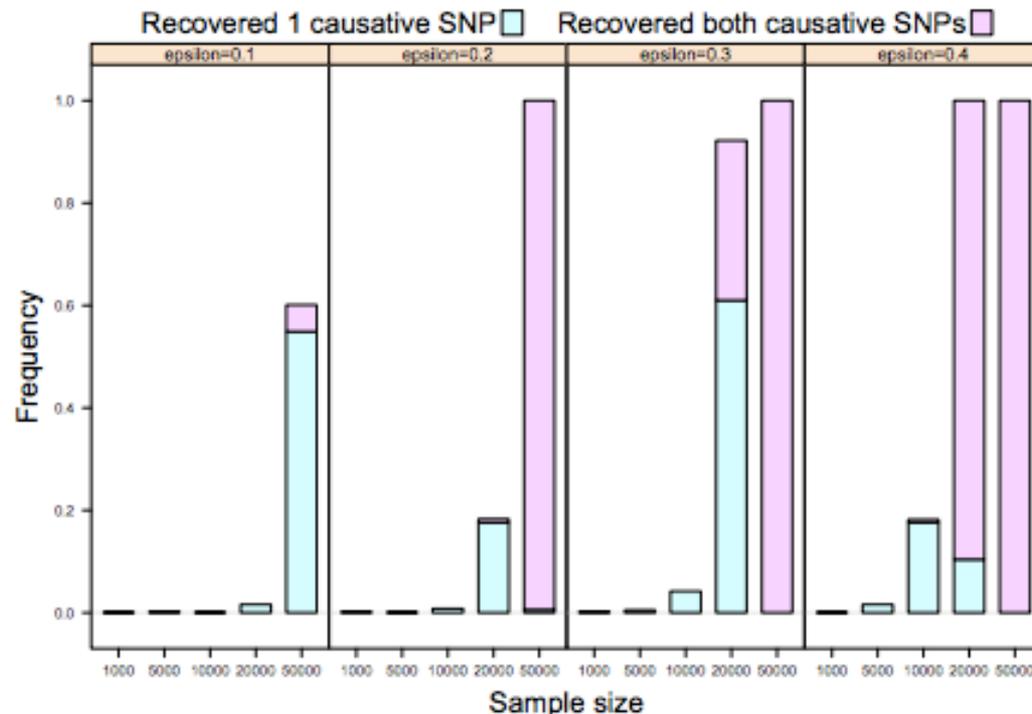


Exploring techniques to protect privacy



Differential Privacy

Genome Wide Association Studies (GWAS)



Computing number/location of SNPs associated to disease
Significance/correlation between a SNP and a disease

A. Johnson and V. Shmatikov. "Privacy-Preserving Data Exploration in Genome-Wide Association Studies." Proceedings of KDD, 2013

Computing on Encrypted Genomes

Genomic datasets often used for association studies

Encrypt data & outsource to the cloud

- Perform private computation over encrypted data

- Using partial & fully homomorphic encryption

Examples:

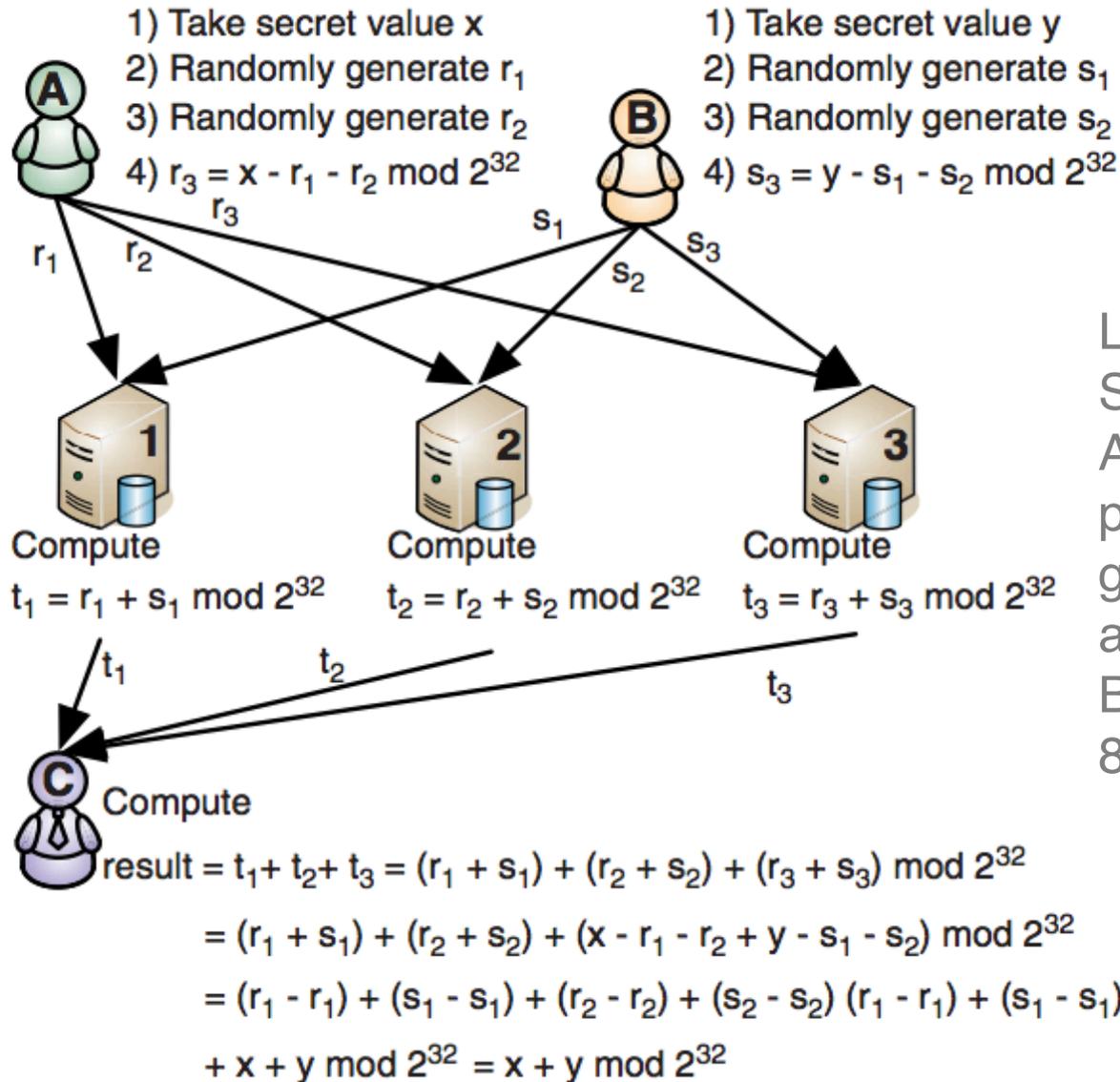
- Pearson Goodness-of-Fit test, linkage disequilibrium

- Estimation Maximization, Cochran-Armitage TT, etc.

K. Lauter, A. Lopez-Alt, M. Naehrig.

Private Computation on Encrypted Genomic Data

Computing on Encrypted Genomes



L. Kamm, D. Bogdanov, S. Laur, J. Vilo.
 A new way to protect privacy in large-scale genome-wide association studies.
 Bioinformatics 29 (7): 886-893, 2013.

Private *Personal* Genomic Tests

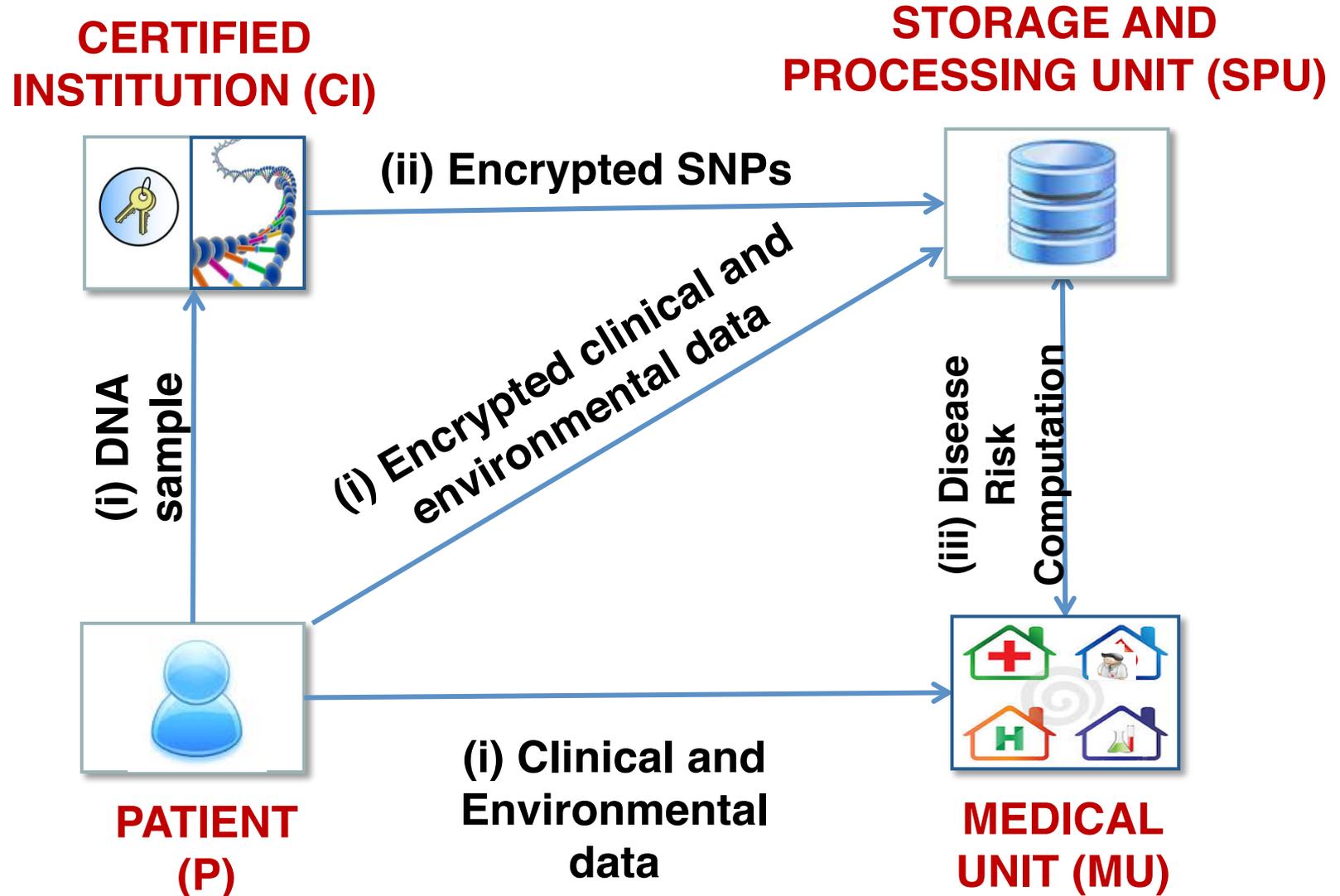
Individuals retain **control** of their sequenced genome

Allow doctors/labs to run genetics tests, but:

1. Genome never disclosed, only test output is
2. Pharmas can keep test specifics confidential

... two main approaches ...

1. Using Semi-Trusted Parties



1. Using Semi-Trusted Parties

Ayday et al. (WPES'13)

Data is encrypted and stored at a “Storage Process Unit”
Disease susceptibility testing

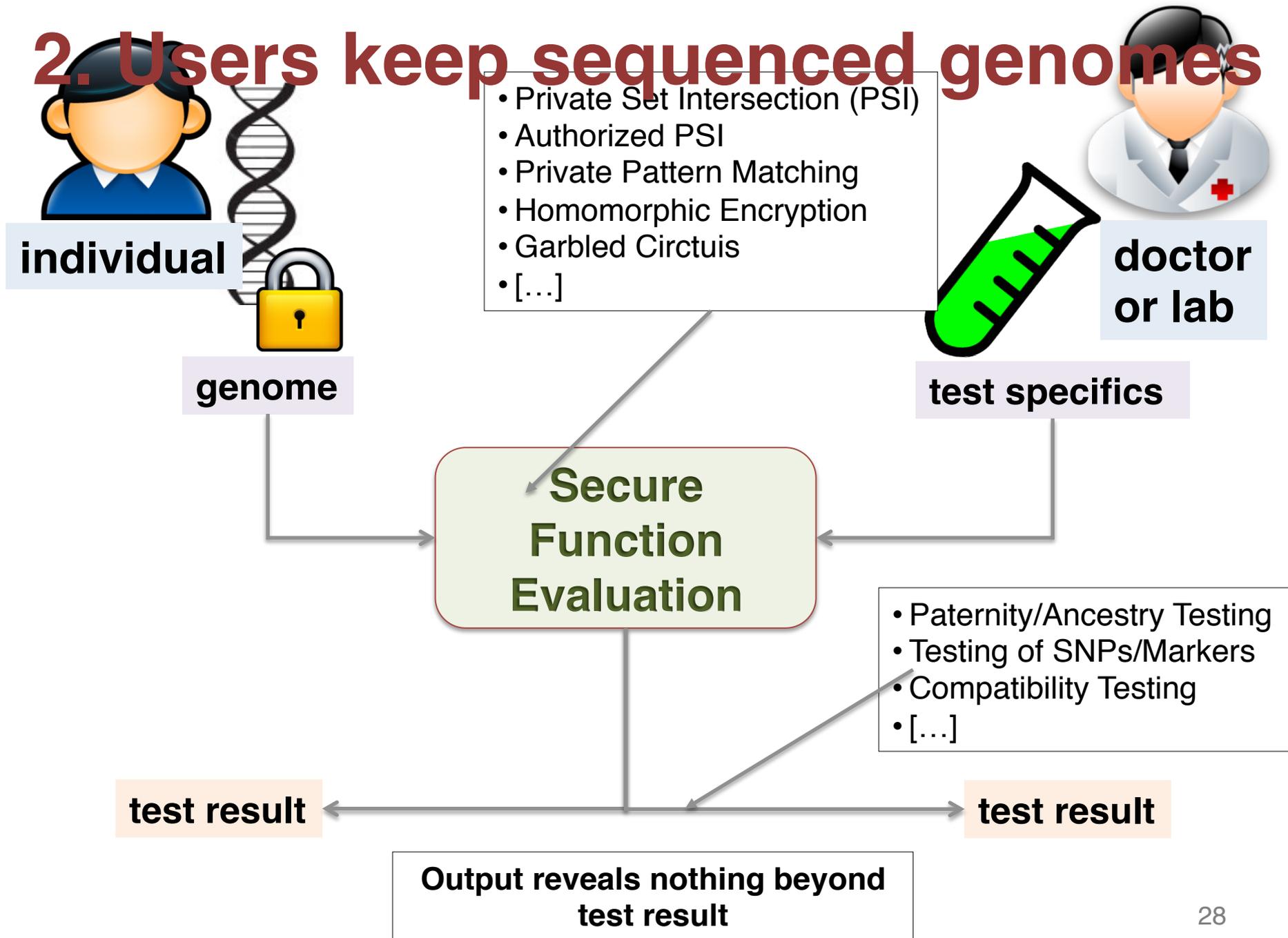
Ayday et al. (DPM'13)

Encrypting raw genomic data (short reads)
Allowing medical unit to privately retrieve them

Danezis and De Cristofaro (WPES'14)

Regression for disease susceptibility

2. Users keep sequenced genomes



2. Users keep sequenced genomes

Baldi et al. (CCS'11)

Privacy-preserving version of a few genetic tests, based on private set operations

Paternity test, Personalized Medicine, Compatibility Tests
(First work to consider fully sequenced genomes)

De Cristofaro et al. (WPES'12), extends the above

Framework and prototype deployment on **Android**

Adds Ancestry/Genealogy Testing

Open Problems

Where do we store genomes?

Encryption can't guarantee **security** past 30-50 yrs

Reliability and **availability** issues?

Cryptography

Efficiency overhead

Data representation **assumptions**

How much understanding required from **users**?

Why do we even care about genome privacy?

We all leave biological cells behind...

Hair, saliva, etc., can be collected and sequenced?

Compare this “attack” to re-identifying millions of DNA donors or hacking into 23andme...

The former: expensive, prone to mistakes, only works against a handful of targeted victims

The latter: very “scalable”

Epilogue

Whole Genome Sequencing

A revolution in healthcare

Raises worrisome privacy/ethical concerns

The Genomic Privacy research community

Understanding the privacy issues

Privacy-preserving testing on whole genomes

Possible using efficient crypto protocols & cross-discipline collaboration

A number of open research issues...

For more info:

<http://genomeprivacy.org>

Also:

E. Ayday, E. De Cristofaro, J.P. Hubaux, G. Tsudik.

“Whole Genome Sequencing: Revolutionary
Medicine or Privacy Nightmare?”

IEEE Computer Magazine



Thank you!

Special thanks to

E. Ayday, P. Baldi, R. Baronio, G. Danezis, S. Faber,
P. Gasti, J-P. Hubaux, B. Malin, G. Tsudik