

The Genomics Revolution: Innovation Dream or Privacy Nightmare?

Emiliano De Cristofaro
<https://emilianodc.com>

TL;DR

Progress in Genomics:

- Enables advances in medicine and healthcare

- Genetic testing for the masses

- Prompts a greater good vs privacy tension

Genomic Data is:

- Sensitive

- Hard to anonymize / de-identify



WGS Progress

Some dates

1970s: DNA sequencing starts

1990: The “Human Genome Project” starts

2003: First human genome fully sequenced

2012: UK announces sequencing of 100K genomes

Some numbers

\$3B: Human Genome Project

\$250K: Illumina (2008)

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

\$1K: Illumina (2014)

How to read the genome?



Sequencing

Determining the full nucleotide order of an organism's genome



Genotyping

Testing for genetic differences using a set of markers

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

Genetic Gamble

New Approaches to Fighting Cancer

PART ONE
A Race to Leukemia's
Source

PART TWO
Promise and
Heartbreak

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*}

MAY 27, 2013

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

time.com

Time



health overview

[Print my health overview](#) | [Share my health results](#)

Show results for [redacted]

[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) new	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...](#)

The genotyping services of 23andMe are performed in LabCorp's CLIA-certified laboratory. The tests have not been cleared or approved by the FDA but have been analytically validated according to CLIA standards. The information on this page is intended for research and educational purposes only, and is not for diagnostic use.

Genetic Ethnicity



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

1 - 25 of 424



You

[UPDATE YOUR PROFILE](#)

1.68% shared, 5 segments

J2a2

[Send an Introduction](#)

1.30% shared, 3 segments

Paternal

[Send a Message](#)

Cousin
1.03% shared, 2
segments

R1b1b2

[Send an Introduction](#)

Cousin
0.45% shared, 2
segments

H7

[Send an Introduction](#)

Cousin
0.42% shared, 2
segments

H1

[Send an Introduction](#)

Cousin
0.40% shared, 2
segments

San Diego, California

U1 = 62 =

[Send a Message](#)

0.37% shared, 2 segments

fathers father prince Edward isla...

name owen t m msaac mothers ...	msaac	K1a1b
---------------------------------	-------	-------

R1b1b2a1a

[Send a Message](#)

Cousin
0.40% shared, 1
segment

Utah

California

U3b1 T

T

[Send an Introduction](#)



```
ex1.sam
ex1.sam > No Selection
1 @HD VN:1.0 SO:coordinate
2 @SQ SN:seq1 LN:5000
3 @SQ SN:seq2 LN:5000
4 @CO Example of SAM/BAM file format.
5 B7_591:4:96:693:509 73 seq1 1 99 36M * 0 0 CACTAGTGGCTCATTGTAAATGTGTGGTTTAACTCG
6 EAS54_65:7:152:368:113 73 seq1 3 99 35M * 0 0 CTAGTGGCTCATTGTAAATGTGTGGTTTAACTCGT
7 EAS51_64:8:5:734:57 137 seq1 5 99 35M * 0 0 AGTGGCTCATTGTAAATGTGTGGTTTAACTCGTCC
8 B7_591:1:289:587:906 137 seq1 6 63 36M * 0 0 GTGGCTCATTGTAATTTTTGTTTTAACTCTTCTCT
9 EAS56_59:8:3:671:758 137 seq1 9 99 35M * 0 0 GGTCAATGTAAATGTGTGGTTTAACTCGTCCATGG
10 EAS56_61:6:8:467:211 73 seq1 13 99 35M * 0 0 ATTGTAAATGTGTGGTTTAACTCGTCCATGGCCCA
11 EAS114_28:5:296:340:699 137 seq1 13 99 36M * 0 0 ATTGTAAATGTGTGGTTTAACTCGTCCATGGCCAG
12 B7_597:6:194:894:408 73 seq1 13 99 35M * 0 0 GTAAATGTGTGGTTTAACTCGTCCATGGCCAGC
13 EAS188_4:8:12:628:973 89 seq1 18 75 35M * 0 0 TAAATGTGTGGTTTAACTCGTCCATGGCCAGCATT
14 EAS51_66:7:68:402:50 137 seq1 22 99 35M * 0 0 GTGTGGTTTAACTCGTCCATGGCCAGCATTGGG
15 EAS114_30:6:298:115:564 137 seq1 22 99 35M * 0 0 GTGTGGTTTAACTCGTCCATGGCCAGCATTAGGG
16 B7_591:3:188:662:155 73 seq1 24 99 36M * 0 0 GTGGTTTAACTCGTCCATGGCCAGCATTAGGGAGC
17 EAS56_59:2:225:608:291 73 seq1 28 99 35M * 0 0 TTTAACTCGTCCATGGCCAGCATTAGGGATCTGT
18 EAS51_66:7:328:397:316 73 seq1 29 99 35M * 0 0 TTAACCTCGTCCATGGCCAGCATTAGGGAGCTGTG
19 EAS51_64:5:257:960:682 73 seq1 31 75 35M * 0 0 AACTCGTCCATGGCCAGCATTAGGGAGCTGTGGA
20 EAS54_61:4:143:69:578 99 seq1 36 98 35M = 185 184 GTACATGGCCAGCATTAGGGAGCTGTGGACCCCG
```

Security Researcher's Perspective

Genome = the ultimate **identifier**

Hard to anonymize / de-identify

Treasure trove of **sensitive** information

Ethnic heritage, predisposition to diseases

Sensitivity is **perpetual**

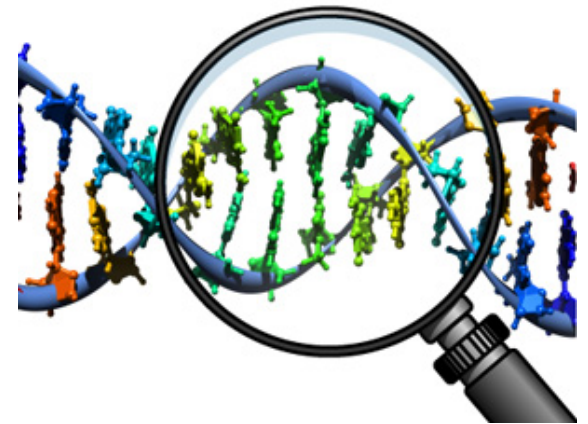
Cannot be “revoked”

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



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.” ACM 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, 4, Aug. 2008

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.

From the onset, the Personal Genome Project,
a project to create a public database of



Harvard Professor Latanya Sweeney

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

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



Most Read

1

Read the Iranian foreign minister's passive aggressive response to Tony Blair

2

Where the world's migrants go, in 10 minutes

3

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

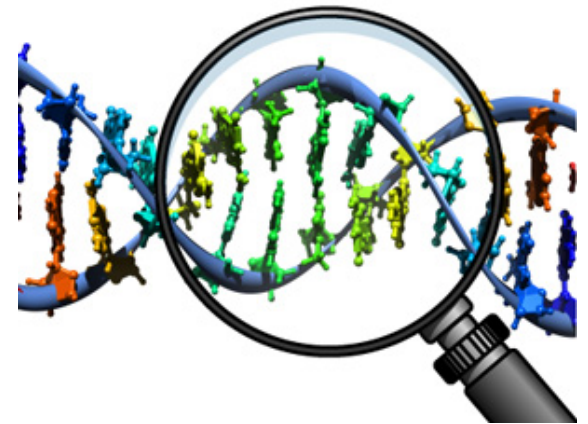
4

A new theory for why the bees are vanishing

5

The rise of a new research community

Studying privacy issues

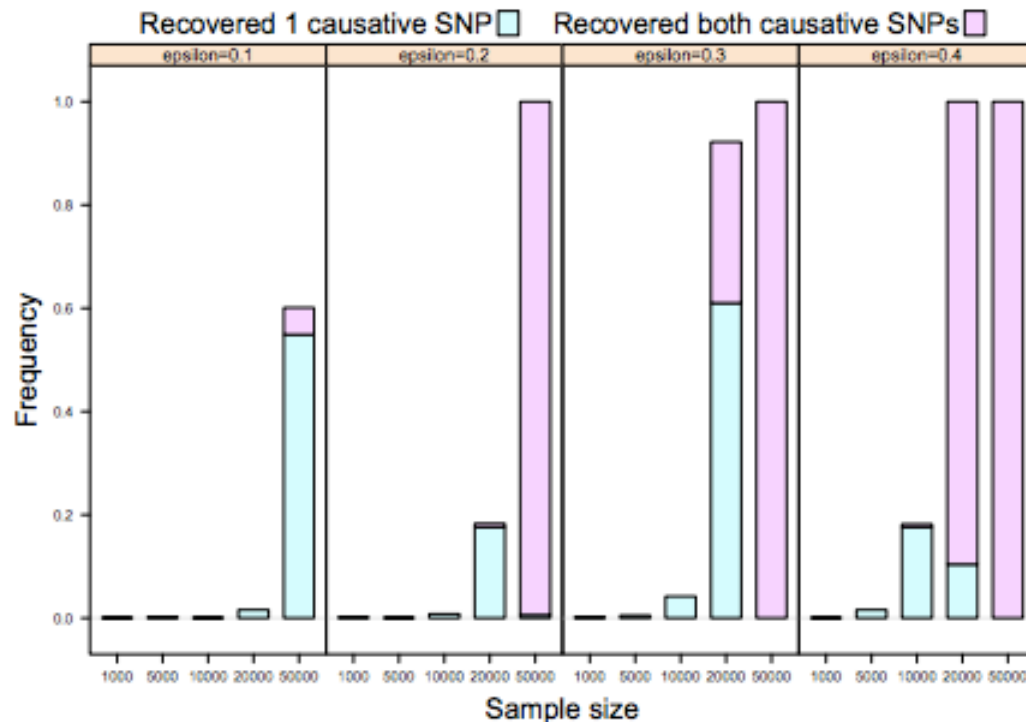


Exploring techniques to protect privacy



Differential Privacy

Privacy in 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

Privacy-Preserving Genomic Tests

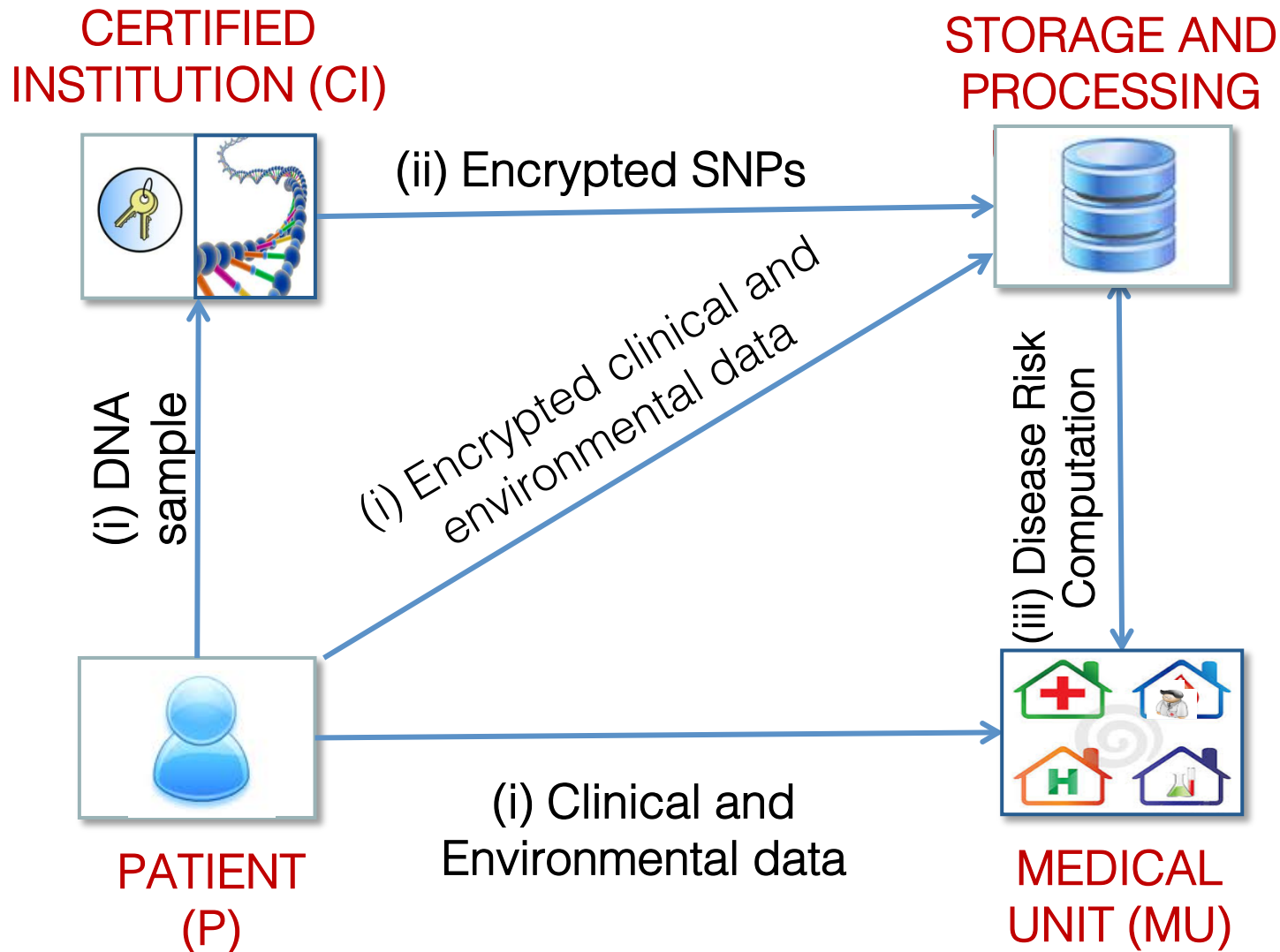
Individuals retain **control** of their sequenced genome

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

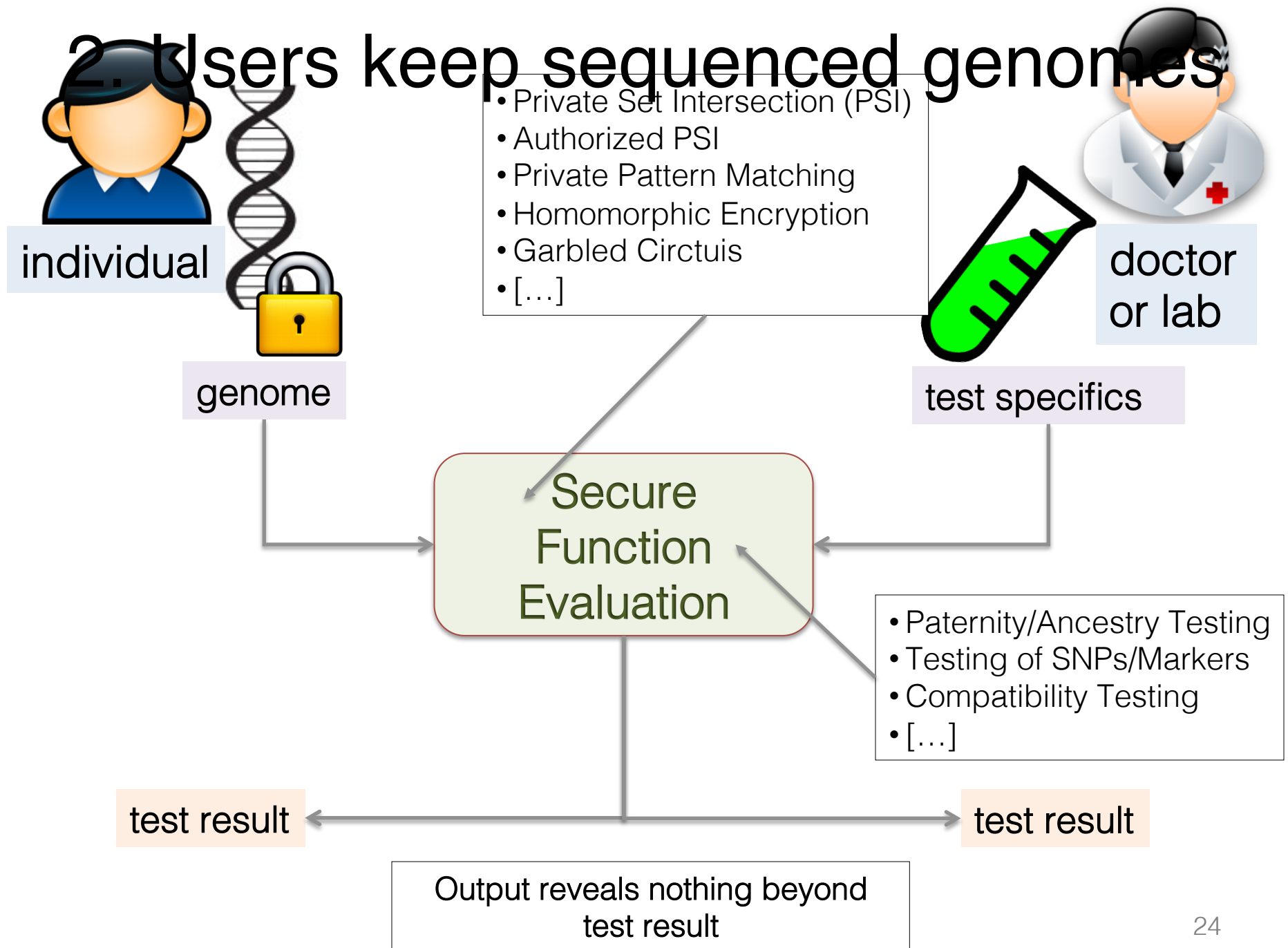
Ayday et al. (DPM'13)

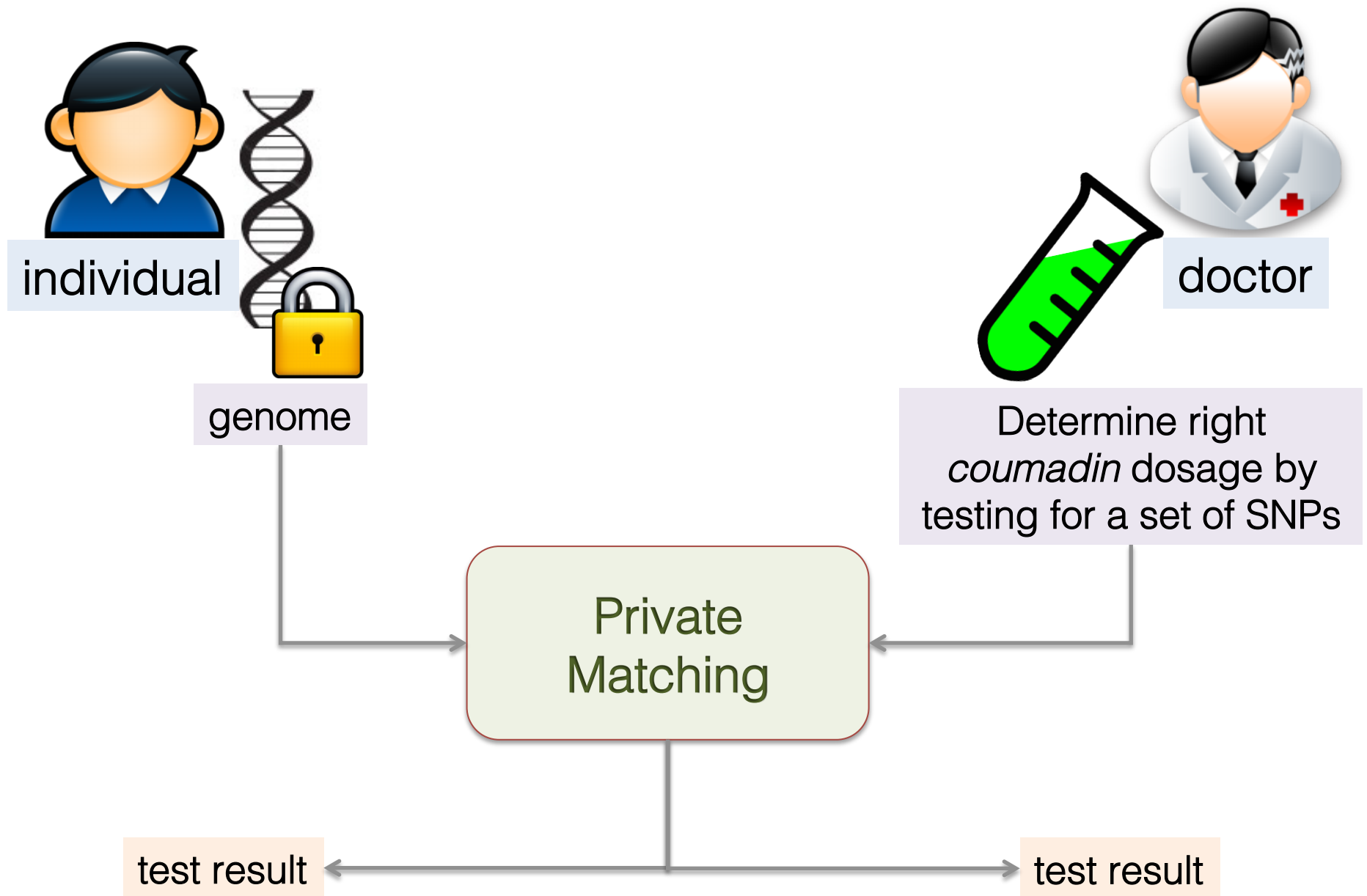
Encrypting raw genomic data (short reads)
Allowing medical unit to privately retrieve them

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 databases

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

The latter: very “scalable”

Thank you!

Special thanks to

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