



Genome Privacy: A Computer Scientist's Perspective

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

Genomics: A CS Perspective

[illegible]

Genomics: A CS Perspective

Once sequenced... a genome becomes an (annotated) file

Needs to be stored somewhere

Can be queried/searched/tested/etc

But... not all data are
created equal!

Security Researcher's Perspective

Ultimate identifier

Hard to anonymize / de-identify

Once leaked, we cannot “revoke” it

Extremely sensitive information

Ethnic heritage, predisposition to diseases

Leaking one's genome \approx leaking **relatives' genome***

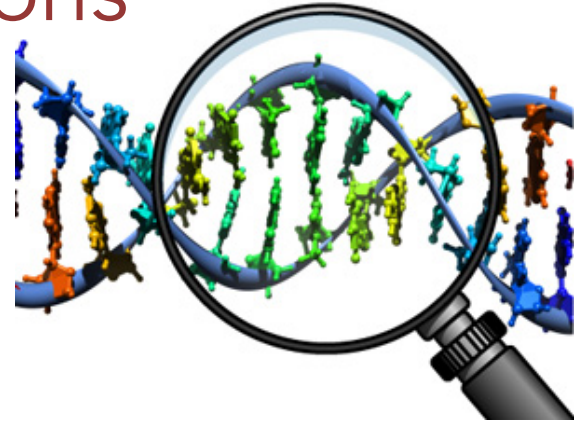
Sensitivity of the genome is (almost) perpetual

Long after owner's death

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

The rise of a new research community

Studying the privacy implications



Exploring techniques to protect privacy



Studying Privacy

Re-identification of anonymous DNA donors

Infer surnames using (public) information available from popular genealogy sites*

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

Studying Privacy

OK... anonymization doesn't really work.
What about aggregation?

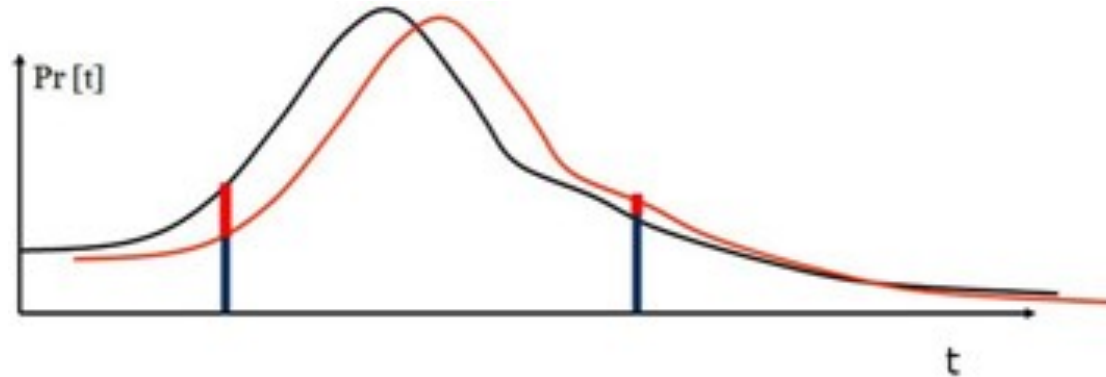
Even statistics from allele frequencies can be used to identify genetic trial participants

Rui Wang et al. "Learning Your Identity and Disease from Research Papers: Information Leaks in Genome Wide Association Study." Proceedings of ACM CCS, 2009

Routes for breaching privacy

Y. Erlich and A. Narayanan. "Routes for Breaching and Protecting Genetic Privacy." Nature Review Genetics, Vol. 15, No. 6, 2014

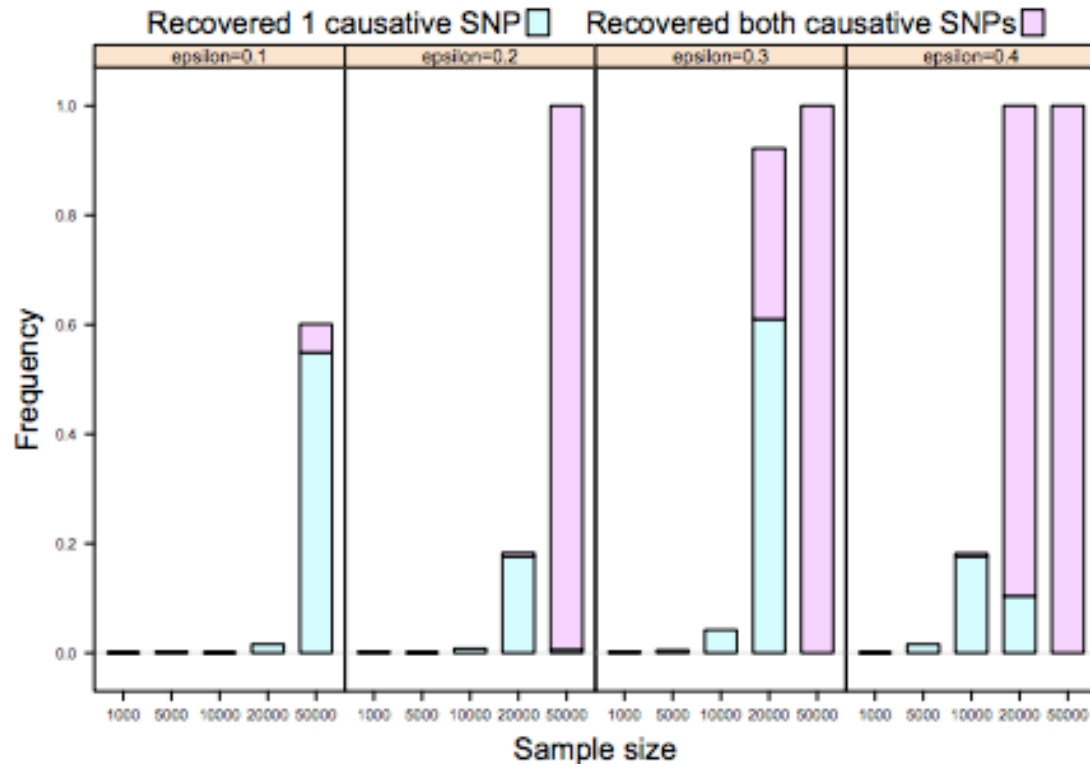
Differential Privacy



Maximizing the **accuracy** of queries from statistical databases

Minimizing the chances of **identifying** its records

Differential Privacy

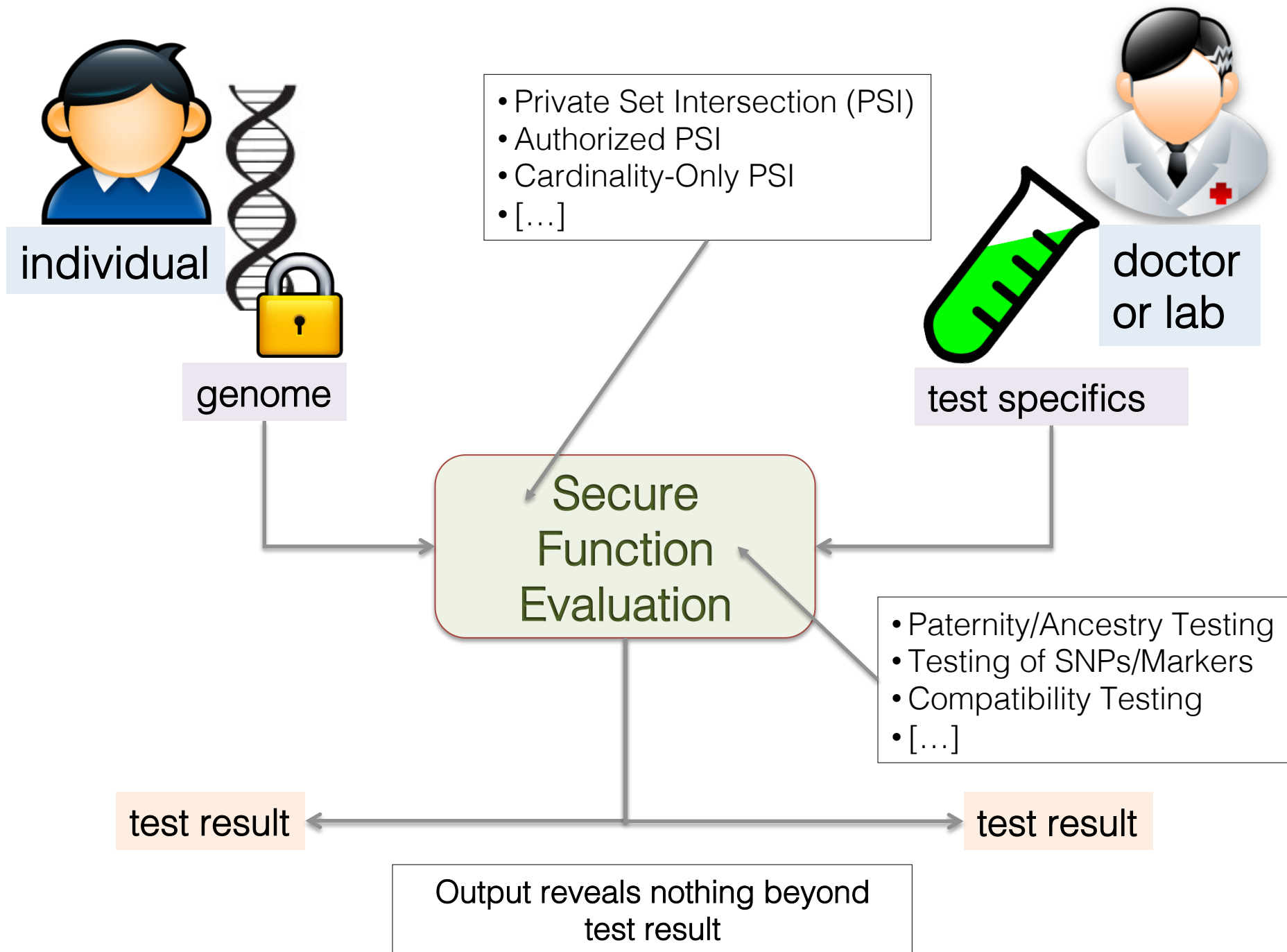


Supporting Genome Wide Association Studies (GWAS)

Computing number and location of SNPs associated to disease
Test significance, correlation, etc. between a SNP and a disease

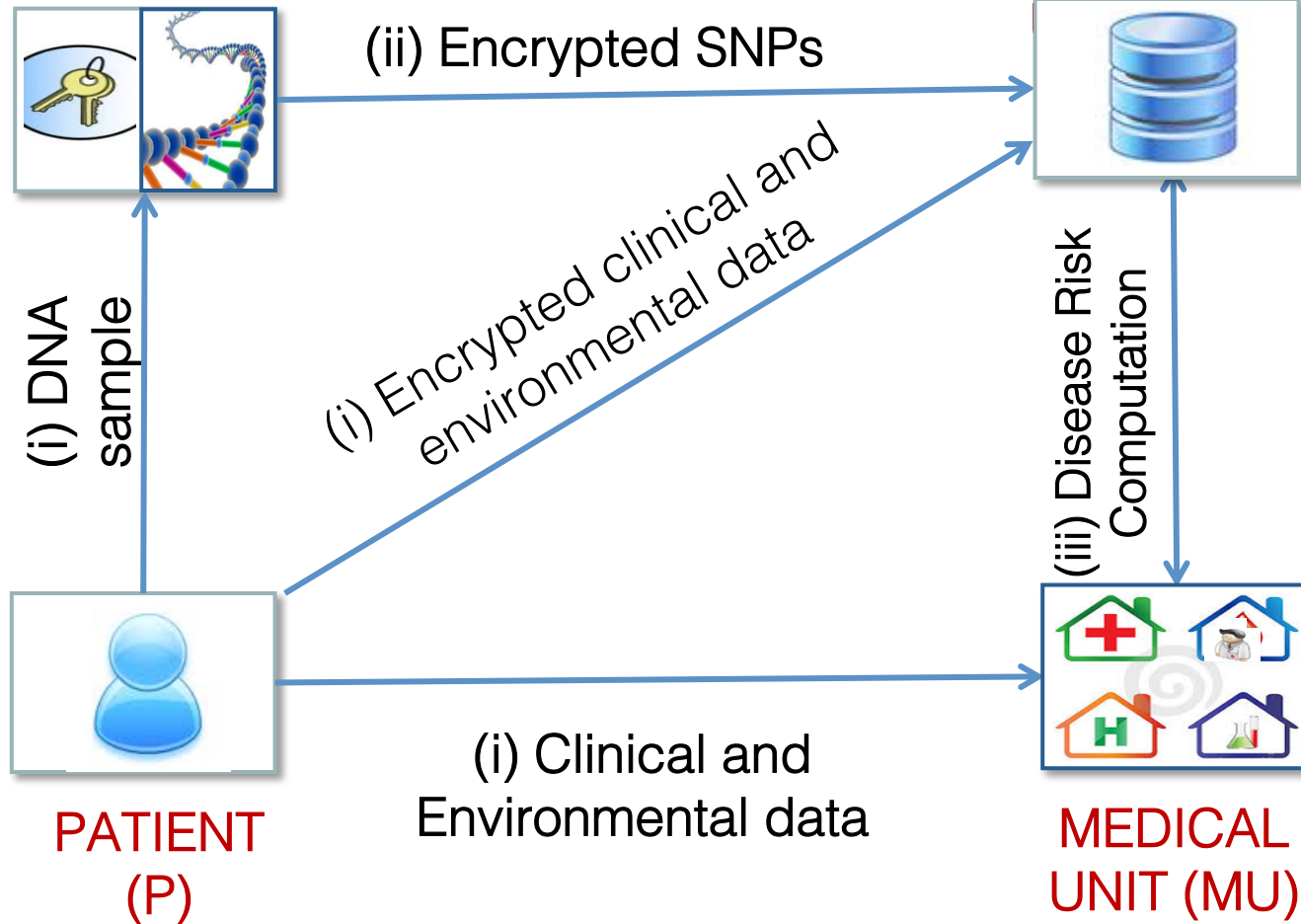
A. Johnson and V. Shmatikov. "Privacy-Preserving Data Exploration in Genome-Wide Association Studies." Proceedings of KDD, 2013

Privacy-Friendly Personal Genomics



CERTIFIED
INSTITUTION (CI)

STORAGE AND
PROCESSING



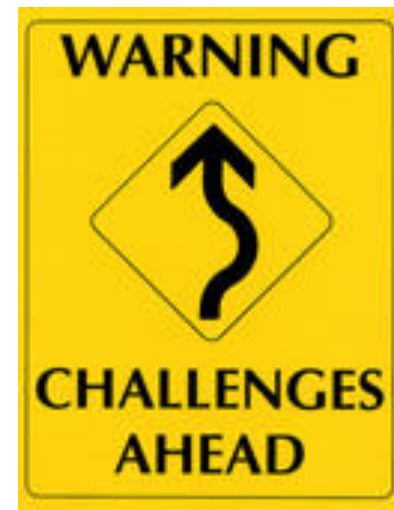
Open Problems

Guarantee security past 30 years?

User's involvement? Usability?

Computation overhead?

Data representation assumptions?



2nd International Workshop on
Genome Privacy and Security
(GenoPri 2015)

May 21, San Jose, CA
<http://genopri.org>

Submission deadline: January 25

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

Question?

Why do we care about genome privacy???

We all leave biological cells behind...

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

But... collecting and sequencing samples is
expensive, illegal, prone to mistakes

Different scale of attacks!