# Getting Real: Merging Bioinformatics Standards and Crypto-State-of-the-Art

Workshop Omics in Medical Research – TMF, Berlin, 2017

Prof.Dr. Kay Hamacher

TU Darmstadt
http://www.kay-hamacher.de

## Examples (1): Oblivious RAM (ORAM)

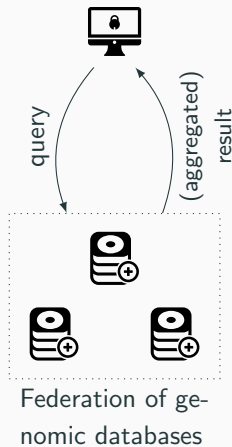write several times; shuffle data chunks; hide query $\Rightarrow$ protect queries



- pattern matching $\checkmark$
- DNA fingerprint(s) $\equiv$ description vector $\checkmark$
- Bayesian learning/updates $\checkmark$

Karvelas, Peter, Katzenbeisser, Tews, Hamacher *"Privacy Preserving Whole Genome Sequence Processing through Proxy-Aided ORAM"* Workshop on Privacy in the Electronic Society (WPES2014):1-10

## Examples (2): Privacy-Preserving VCF Queries

- Genomic data is sensitive data
- Medical treatment and research profits from access to large genomic datasets
- Trivially aggregating multiple genomic databases . . .
    - allows for more detailed analysis
    - but risks privacy of each dataset and entry
    - might not be desirable by the database provider
- **Solution:** Secure Multi-Party Computation (SMPC) for private aggregation and database queries



Federation of genomic databases

# Variant Call Format (VCF)



**Figure 1:** VCF format example – from vcftools website.

Beacon Network (GA4GH Project) established to evaluate eagerness of institutions world-wide to engage in distributed variant query service.

**"Do you have a genome with mutation X in your database?"**

**"Yes"** / **"No"**

## Beacon Privacy Issues

Recent work (Thenen, N., Ayday, E., C. Ercument; bioRxiv; 200147; Sep. 2017) strongest challenge to privacy of Beacon networks yet

- Need only 5 queries to re-identify member in 65-individuals beacon (95 % confidence)
    - even when single-nucleotide polymorphisms (SNPs) with minor allele frequency (MAF) < 5% are hidden
    - exploiting linkage disequilibrium
    - and high-order Markov chains

Our work mitigates this risk by

- hiding which beacon contributed to total
- incurring a threshold on count-queries
- enabling a reverse follow-up scenario

## Variant Query Format

- Compression/Encoding of variant data into $(\kappa = 32, \psi = 16)$ bit (key, value) space
- Key holds the position in a reference genome $(\log_2(3.2 \cdot 10^9) \approx 31.6 < 32)$
- Small value space must incur variant information loss
- But domain-specific knowledge makes a virtue of this necessity to let similar variants match
- Allow same key multiple times to store e.g. diploid variations or upper and lower length of deletion/copy number variation (CNV)/inversion

## VQF Translation

| Variant | Stored Information | #Values |
|---------|---------------------|---------|
| SNP/SNV | Alternative nucleotide. | 4 |
| Deletion/CNV /Inversion | Up and down rounded logarithm (base $b = 2$) of length, up to log-length $s$. Plus frameshift bit for deletions. | $s = 16$ |
| Insertion | Up to $s_{\text{ins}} = 7$ inserted nucleotides and a frameshift bit. | $2 \cdot 4^{s_{\text{ins}} = 7}$ |
| *Other* | Only flag as *other* variation. Captures more complex variations. | 1 |

## Privacy-Preserving Genomic Queries

Supported queries are of the form:

```
SELECT f(*) FROM Variants
WHERE ((loc_1, var_1), ..., (loc_m, var_m)) IN Genome
AND cancer = X AND ... AND age_min ≤ age ≤ age_max
```

`f(*)` can be an arbitrary function (e.g., aggregation, threshold).

D. Demmler, K. Hamacher, T. Schneider, S. Stammler.
*Privacy-Preserving Whole-Genome Variant Queries* 16th Int. Conf.
Cryptology And Network Security (CANS) accepted, 2017.

## Schematic Overview



Independent database providers

DB$_1$ · DB$_2$ · … · DB$_n$

$D'_1$, $D''_1$, …, $D'_n$, $D''_n$

SMPC proxies — $D'$, $D''$, SMPC

Client $C$, querying databases

# The ABY Framework

- ABY (Demmler, D., Schneider, T., Zohner, M. (NDSS'2015)) is a `C++` framework for secure two-party computation
- 3 supported protocols:
    - Arithmetic sharing:
      $+$ secret-sharing
    - Boolean sharing (GMW):
      $\oplus$ secret-sharing
    - Yao's garbled circuits
- Open source on github:
  `http://encrypto.de/code/ABY`

**Figure 2:** Online runtime in ms.

A typical query from 100 000 variants on 5 positions with $\kappa = 32$ and $\psi = 16$ needs 178 ms runtime and 11 MiB communication in the online phase. (Setup phase: 3.4 s, 733 MiB)

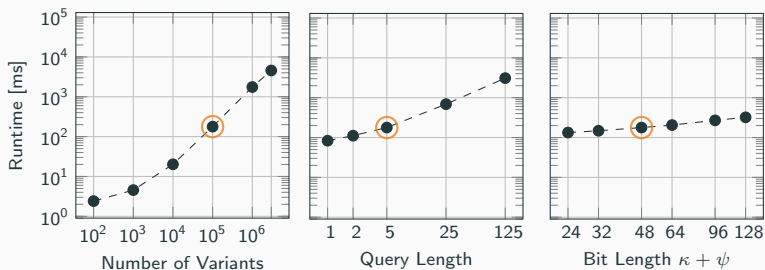# Summary & Acknowledgements

- Trust not essential: computational trust vs. social/legal construct

- work in security / protocol design / cryptography mature and applicable

- loss of (computational) efficiency $\checkmark$ but trade-off might be worth it

- framework(s) do exist; but you <u>need</u> experts in crypto & bioinfo
  *"Anyone, from the most clueless amateur to the best cryptographer, can create an algorithm that he himself can't break."* — Bruce Schneier

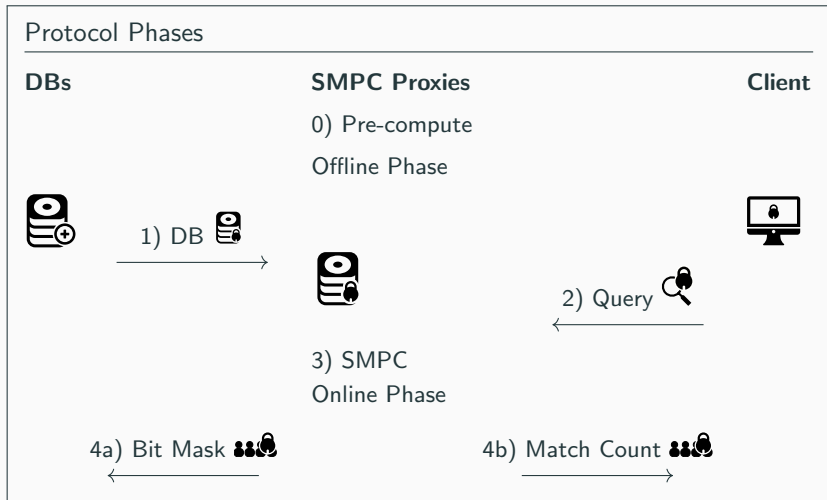http://www.kay-hamacher.de

http://www.kay-hamacher.de

kay.hamacher@cysec.de

# Protocol

## Protocol Phases

**DBs**  **SMPC Proxies**  **Client**

0) Pre-compute

Offline Phase

1) DB →

2) Query ←

3) SMPC

Online Phase

4a) Bit Mask ←  4b) Match Count →

## Circuit Design

- GMW protocol for inherent secret-sharing
- Efficient SIMD operations
- Variant matching: array of equality gates combined with OR tree
- Auxiliary information: combined with AND tree

## Homomorphic Encryption (HE)

Starting with
C. Gentry. *"Fully Homomorphic Encryption Using Ideal Lattices"* 41$^{\text{st}}$ ACM Symposium on Theory of Computing (STOC), 2009

$$
\begin{array}{ccc}
A, B & \xrightarrow{\quad + \quad} & A + B \\[2mm]
\downarrow {\scriptstyle \mathrm{Enc}} & & \uparrow {\scriptstyle \mathrm{Dec}} \\[2mm]
\mathrm{Enc}(A), \mathrm{Enc}(B) & \xrightarrow{\quad \oplus \quad} & \mathrm{Enc}(A + B)
\end{array}
$$

Several framworks, e.g., Damgård, et.al., *"Efficient and secure comparison for online-auctions"*, ACSIP2007, Lecture Notes in Computer Science (4586)