

SQUiD: Ultra-Secure Storage and Analysis of Genetic Data for the Advancement of Precision Medicine

Jacob Blindenbach, Jiayi Kang, Seungwan Hong, Caline Karam, Thomas Lehner, Gamze Gürsoy

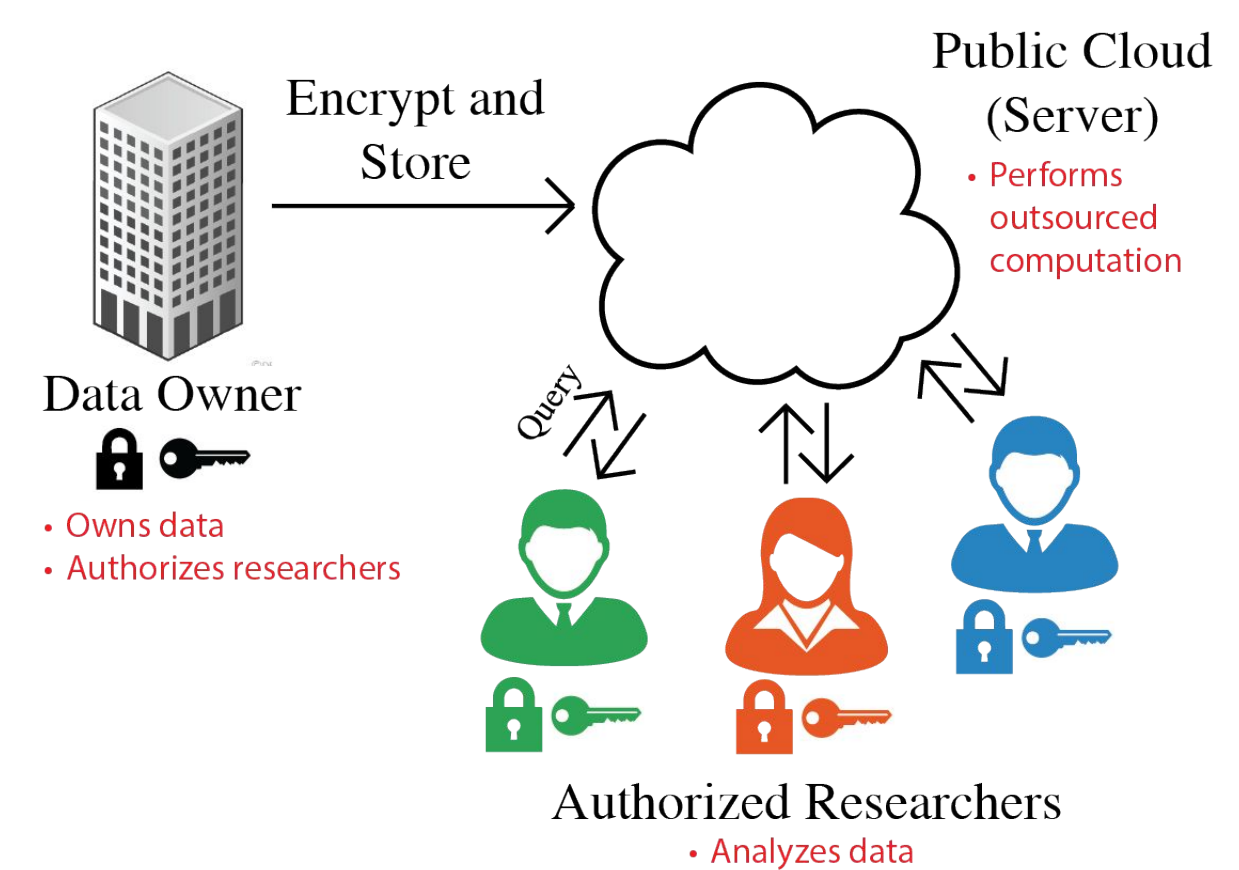
Problem Statement

- Databases are needed for storing extensive sensitive patient disease and genetic information.
- The large amount of data necessitates the use of cloud storage, coupled with strong security measures due to its sensitive nature.

Goal: Ensure the data is secure from both cloud vulnerabilities and unauthorized users, yet accessible for authorized researchers to safely perform queries and analyses.

Approach / Scenario

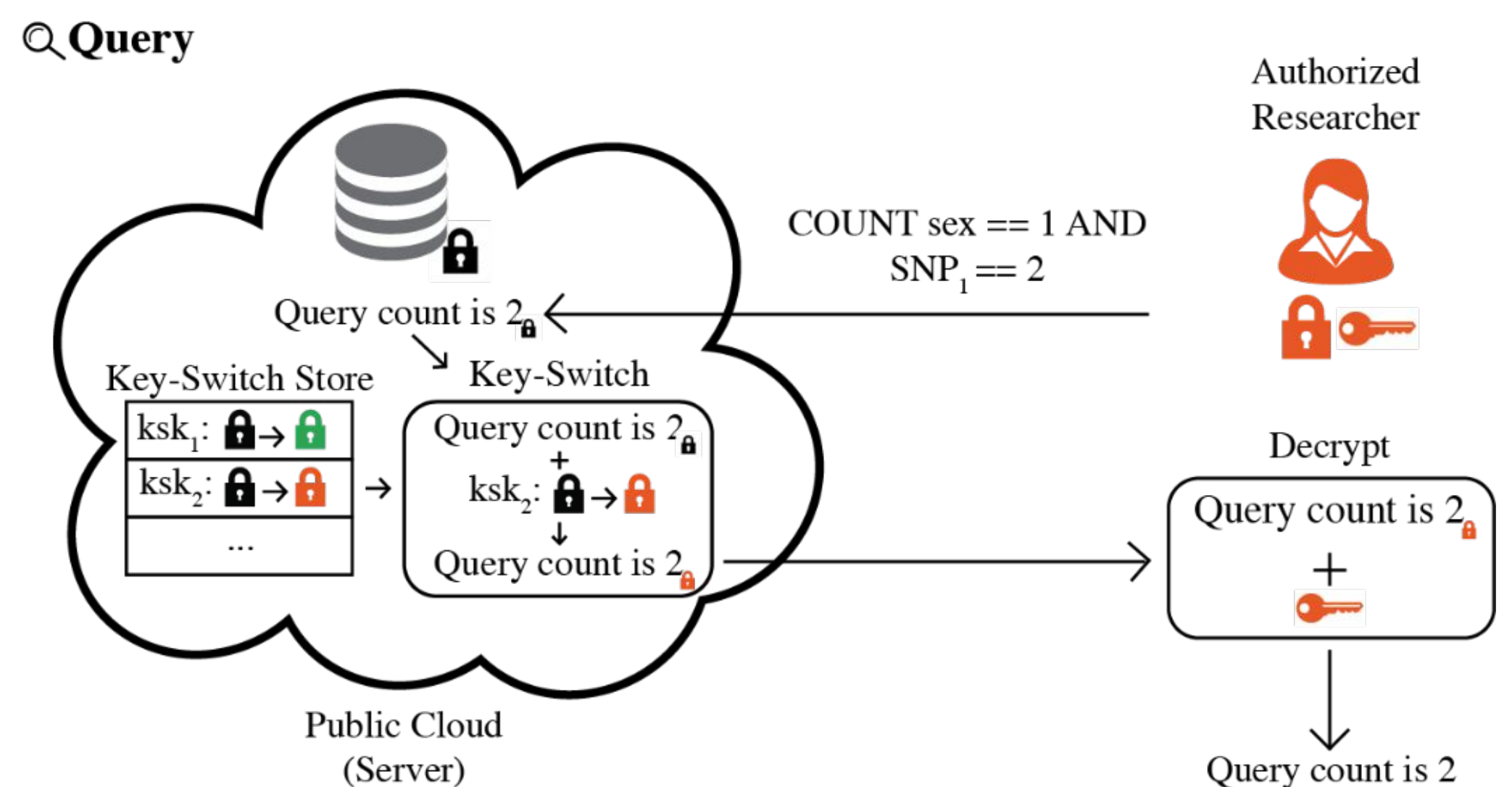
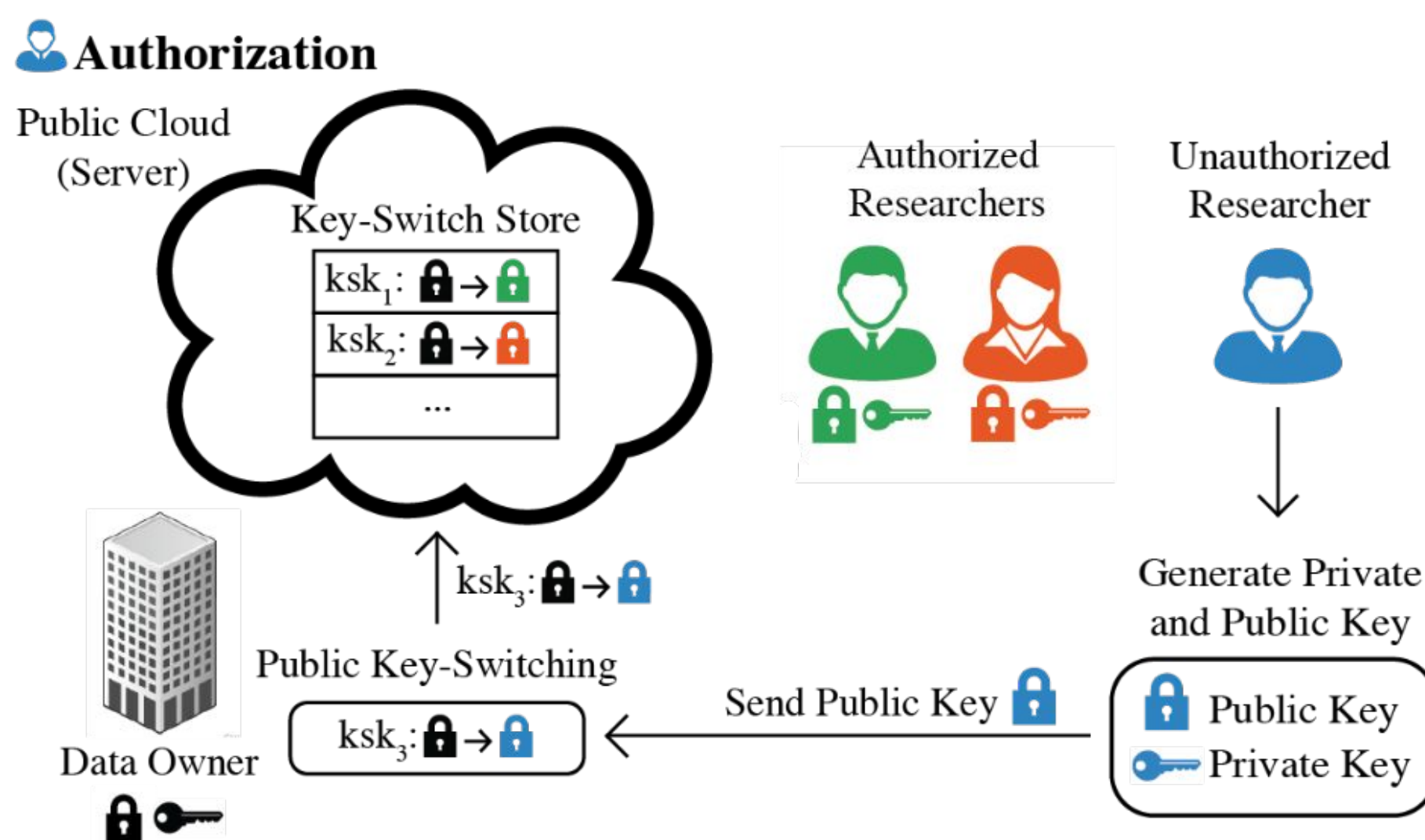
SQUiD (Secure QUeryable Database) is designed for a multiparty setting with a data owner, a public cloud, and multiple researchers.



Results

To implement an encrypted database using homomorphic encryption (HE), several challenges need to be addressed.

HE is designed for two parties: Using public key-switching, multiple researchers can access the database without sharing private keys. Furthermore, the data owner can control researcher access.



HE has limited practical functionality: SQUiD provides four useful queries that analyze genotype phenotype data.

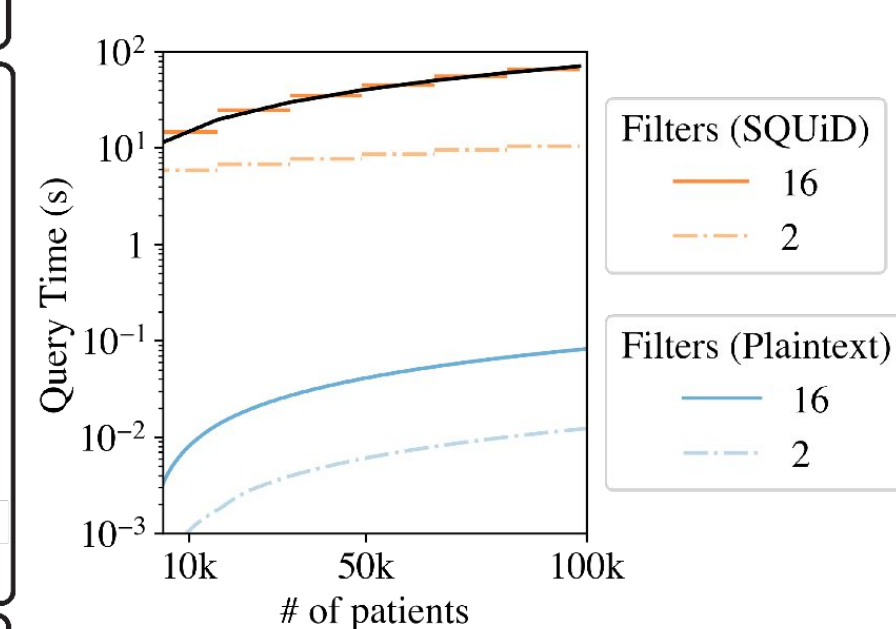
A - Count Query

Query: `COUNT WHERE SNP2 == 1 AND sex == 1`

Patients ID	SNP ₁	SNP ₂	SNP ₃	...	SNP _k	sex	smoker	disease
0	0	1	1	...	1	0	0	1
1	0	1	0	...	0	1	1	0
2	0	0	0	...	0	1	0	1
3	1	1	0	...	0	0	0	1
4	0	1	0	...	0	1	1	1
5	2	2	0	...	0	1	0	1
6	2	0	1	...	0	1	0	0
7	0	0	0	...	2	0	1	0

Compute: `COUNT = 2` → Key Switch

Return: `COUNT = 2`



B - MAF Query

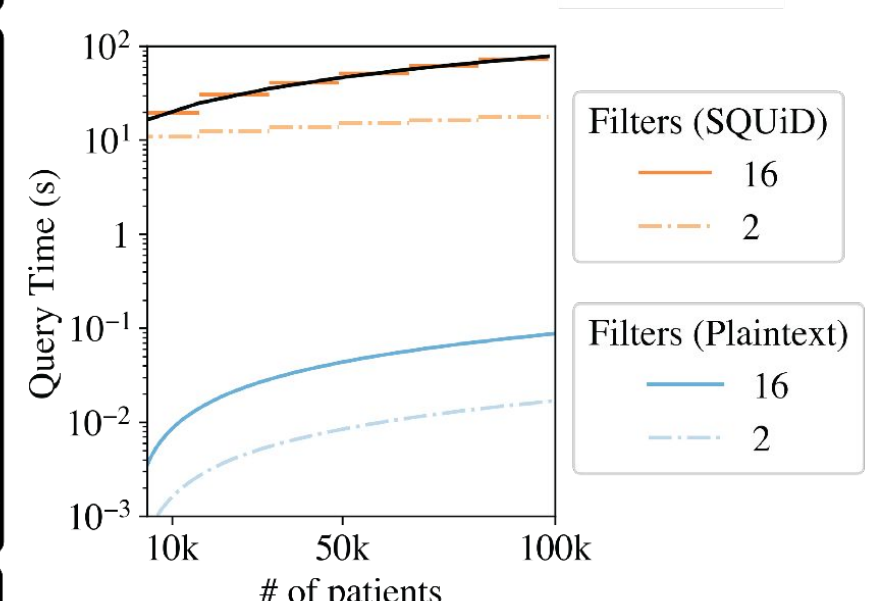
Query: `MAF FOR SNP1 WHERE SNP2 == 1 OR sex == 1`

Patients ID	SNP ₁	SNP ₂	SNP ₃	...	SNP _k	sex	smoker	disease
0	0	1	1	...	1	0	0	1
1	0	1	0	...	0	1	1	1
2	0	0	0	...	0	1	0	0
3	1	1	0	...	0	0	0	1
4	0	1	0	...	0	1	1	1
5	2	2	0	...	0	1	0	1
6	2	0	1	...	0	1	0	0
7	0	0	0	...	2	0	1	0

Compute: Allele Count = 5, 2 + # of filtered patients = 12 → Key Switch

Return: Allele Count = 5, 2 + # of filtered patients = 12

Finalize: `MAF = 5 / 12 = 0.4167`



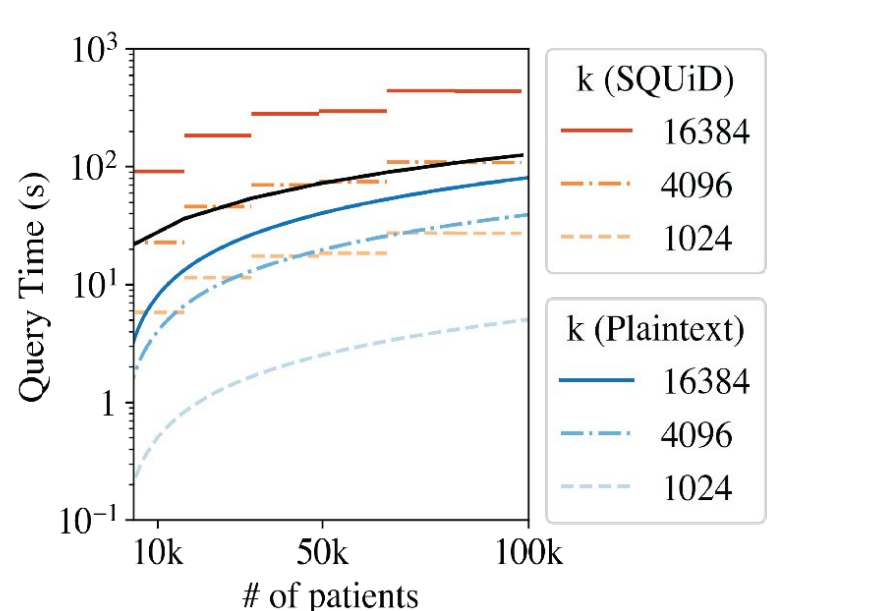
C - PRS Query

Query: `PRS SNPs: [SNP1, SNP2, ..., SNPk] Effect Size: [β1, β2, ..., βk]`

Patients ID	SNP ₁	SNP ₂	SNP ₃	...	SNP _k	sex	smoker	disease
0	0	1	1	...	1	0	0	1
1	0	1	0	...	0	1	1	1
2	0	0	0	...	0	1	0	0
3	1	1	0	...	0	0	0	1
4	0	0	0	...	0	0	1	1
5	2	2	0	...	0	1	0	1
6	2	0	1	...	0	1	0	0
7	0	0	0	...	2	0	1	0

Compute: $\beta_1 \text{SNP}_1 + \beta_2 \text{SNP}_2 + \dots + \beta_k \text{SNP}_k = \text{PRS Scores}$ → Key Switch

Return: PRS Scores



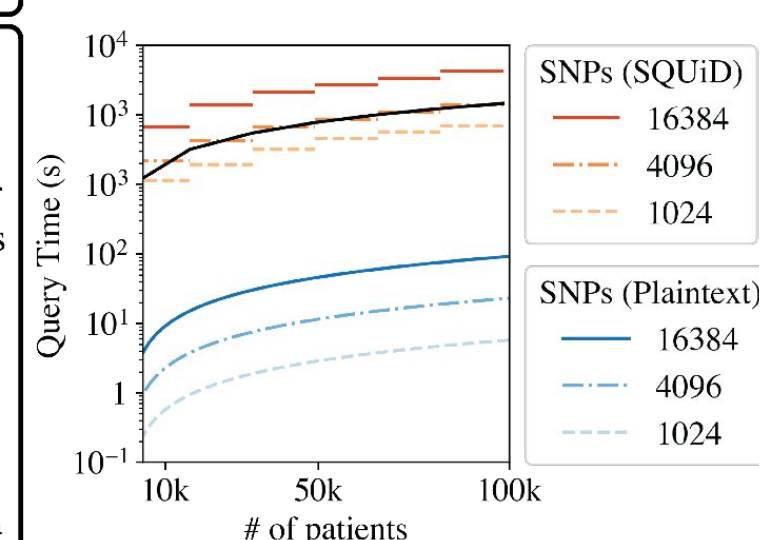
D - Similarity Query

Query Patient: `SNP1 SNP2 SNP3 ... SNPk`

Patients ID	SNP ₁	SNP ₂	SNP ₃	...	SNP _k	sex	smoker	disease
0	0	1	1	...	1	0	0	1
1	0	1	0	...	0	1	1	1
2	0	0	0	...	0	1	0	0
3	1	1	0	...	0	0	0	1
4	0	0	0	...	0	0	1	1
5	2	2	0	...	0	1	0	1
6	2	0	1	...	0	1	0	0
7	0	0	0	...	2	0	1	0

Compute: Count with disease = 4, Count without disease = 1 → Key Switch

Return: Count with disease = 4, Count without disease = 1





SQUiD: Ultra-Secure Storage and Analysis of Genetic Data for the Advancement of Precision Medicine

Jacob Blindenbach, Jiayi Kang, Seungwan Hong, Caline Karam, Thomas Lehner, Gamze Gürsoy

Abstract

Cloud computing provides the opportunity to store the ever-growing genotype-phenotype data sets needed to achieve the full potential of precision medicine. However, due to the sensitive nature of this data and the patchwork of data privacy laws, additional security protections are proving necessary to ensure data privacy and security. Here we present SQUiD, a Secure QUeryable Database for storing and analyzing genotype-phenotype data. With SQUiD, genotype-phenotype data can be stored in a low-security, low-cost public cloud in the encrypted form, which researchers can securely query without the public cloud ever being able to decrypt the data.

SQUiD utilizes homomorphic encryption to enable direct computations on encrypted data, introducing improvements that bolster security while also increasing practicality. Firstly, we developed a new cryptographic primitive, public key-switching, to allow for multi-client interactions with the database without ever sharing sensitive private keys. Secondly, we used ciphertext packing to decrease the storage and query runtime by orders of magnitude. Thirdly, we developed practical protocols to enable four queries commonly used with genotype-phenotype datasets. Specifically, SQUiD can count the number of patients that pass a filter, compute the minor allele frequency (MAF) for SNPs in a filtered patient cohort, compute polygenic risk scores for patients in the database, and perform patient similarity analysis. We demonstrate the usability of SQUiD by replicating various commonly used calculations such as cohort creation for GWAS, MAF filtering, and patient similarity analysis both on synthetic and UK Biobank data.