

RESEARCH

Open Access

# Privacy-preserving search for chemical compound databases

Kana Shimizu<sup>1\*</sup>, Koji Nuida<sup>2,3</sup>, Hiromi Arai<sup>4</sup>, Shigeo Mitsunari<sup>5</sup>, Nuttapong Attrapadung<sup>3</sup>, Michiaki Hamada<sup>6</sup>, Koji Tsuda<sup>7</sup>, Takatsugu Hirokawa<sup>8</sup>, Jun Sakuma<sup>9</sup>, Goichiro Hanaoka<sup>3</sup>, Kiyoshi Asai<sup>7</sup>

From Joint 26th Genome Informatics Workshop and Asia Pacific Bioinformatics Network (APBioNet) 14th International Conference on Bioinformatics (GIW/InCoB2015) Tokyo, Japan. 9-11 September 2015

## Abstract

**Background:** Searching for similar compounds in a database is the most important process for in-silico drug screening. Since a query compound is an important starting point for the new drug, a query holder, who is afraid of the query being monitored by the database server, usually downloads all the records in the database and uses them in a closed network. However, a serious dilemma arises when the database holder also wants to output no information except for the search results, and such a dilemma prevents the use of many important data resources.

**Results:** In order to overcome this dilemma, we developed a novel cryptographic protocol that enables database searching while keeping both the query holder's privacy and database holder's privacy. Generally, the application of cryptographic techniques to practical problems is difficult because versatile techniques are computationally expensive while computationally inexpensive techniques can perform only trivial computation tasks. In this study, our protocol is successfully built only from an additive-homomorphic cryptosystem, which allows only addition performed on encrypted values but is computationally efficient compared with versatile techniques such as general purpose multi-party computation. In an experiment searching ChEMBL, which consists of more than 1,200,000 compounds, the proposed method was 36,900 times faster in CPU time and 12,000 times as efficient in communication size compared with general purpose multi-party computation.

**Conclusion:** We proposed a novel privacy-preserving protocol for searching chemical compound databases. The proposed method, easily scaling for large-scale databases, may help to accelerate drug discovery research by making full use of unused but valuable data that includes sensitive information.

## Introduction

In recent years, the increasing cost of drug development and decreasing number of new chemical entities have become growing concerns [1]. One of the most popular approaches for overcoming these problems is searching for similar compounds in databases [2]. In order to improve the efficiency of this task, it is important to utilize as many data resources as possible. However, the following dilemma prevents the use of many existing data resources.

Unpublished experimental results have been accumulated at many research sites, and such data has scientific value [3]. Since data holders are usually afraid of sensitive information leaking from the data resources, they do not want to release the full data, but they might allow authorized users to search the data as long as the users obtain only search results from which they cannot infer sensitive information. Likewise, private databases of industrial research might be made available if the sensitive information were sufficiently protected. On the other hand, query compounds are also sensitive information for the users, and thus the users usually avoid sending queries and want to download all of the data in order to conduct search tasks on their local computers. In short, we cannot utilize

\* Correspondence: shimizu-kana@aist.go.jp

<sup>1</sup>Biotechnology Research Institute for Drug Discovery, National Institute of Advanced Industrial Science and Technology (AIST), 2-4-7 Aomi Koto-ku, Tokyo, Japan

Full list of author information is available at the end of the article

important data resources because both the data holder and the data user insist on their privacy. Therefore, an emerging issue is to develop novel technology that enables privacy-preserving similarity searches. We show several use cases in the next section.

Let us start by clarifying privacy problems in database searches. In a database search, two types of privacy are of concern: “*user privacy*” (also known as input privacy) and “*database privacy*” (also known as output privacy). The first is equal to protecting the user’s query from being leaked to others including the database holder. The second is equal to protecting the database contents from being leaked to others including the database user, except for the search results held by the user. Here we firstly consider the case of using no privacy-preserving techniques; namely, the user sends a plain query to the server and the server sends the search result. In this case, the user’s query is fully obtained by the server. On the database side, the server’s data is not directly leaked to the user. However, there is a potential risk that the user may infer the database contents from the search results. To protect user privacy, a scheme called single-database private information retrieval (PIR) has been proposed [4]. The simplest method for achieving PIR is that the user downloads all the contents of the database and searches on his/her local computer. Since this naive approach needs a huge communication size, several cryptographic techniques have been developed, in which the query is safely encrypted/randomized in the user’s computer and the database conducts the search without seeing the query. Although PIR is useful for searching public databases, it does not suit the purpose of searching private databases because of the lack of database privacy. Likewise, similarity evaluation protocols keep user privacy [5-7] but they do not sufficiently protect database privacy because the server directly outputs similarity scores that become important hints for inferring database contents.

Generally speaking, it is very difficult to keep both user privacy and database privacy, because the database side must prevent various attacks without seeing the user’s query. Among them, the following two attacks are major concerns.

- Regression attack

Given one data point, the similarity between a target and the data point becomes a strong hint for detecting the target. The accuracy of the detection increases as the number of given data points becomes larger. In fact, a protocol that is not suitably designed may lead to even a small number of queries enabling the database user to detect the target. For example, when the server returns the exact distance between a query and a database

entry, the range of the entry is rapidly narrowed as the number of queries increases, and the entry is finally detected uniquely by only almost the same number of queries as the dimension of the entry (see Figure 1 for a detailed explanation). For example, in the case of using the MACCS keys, which are 166 bit structural key descriptors and often used for representing chemical compounds, a database entry is detected by sending only 166 queries. Therefore, it is necessary for the server to return the minimum information that is sufficient for the purpose of the search. In the Thresholding largely improves database privacy section, we will compare success probability of the regression attack for the case when the server returns minimum information (which our protocol aims for) and the other case when the server returns more information (which the previous method aims for).

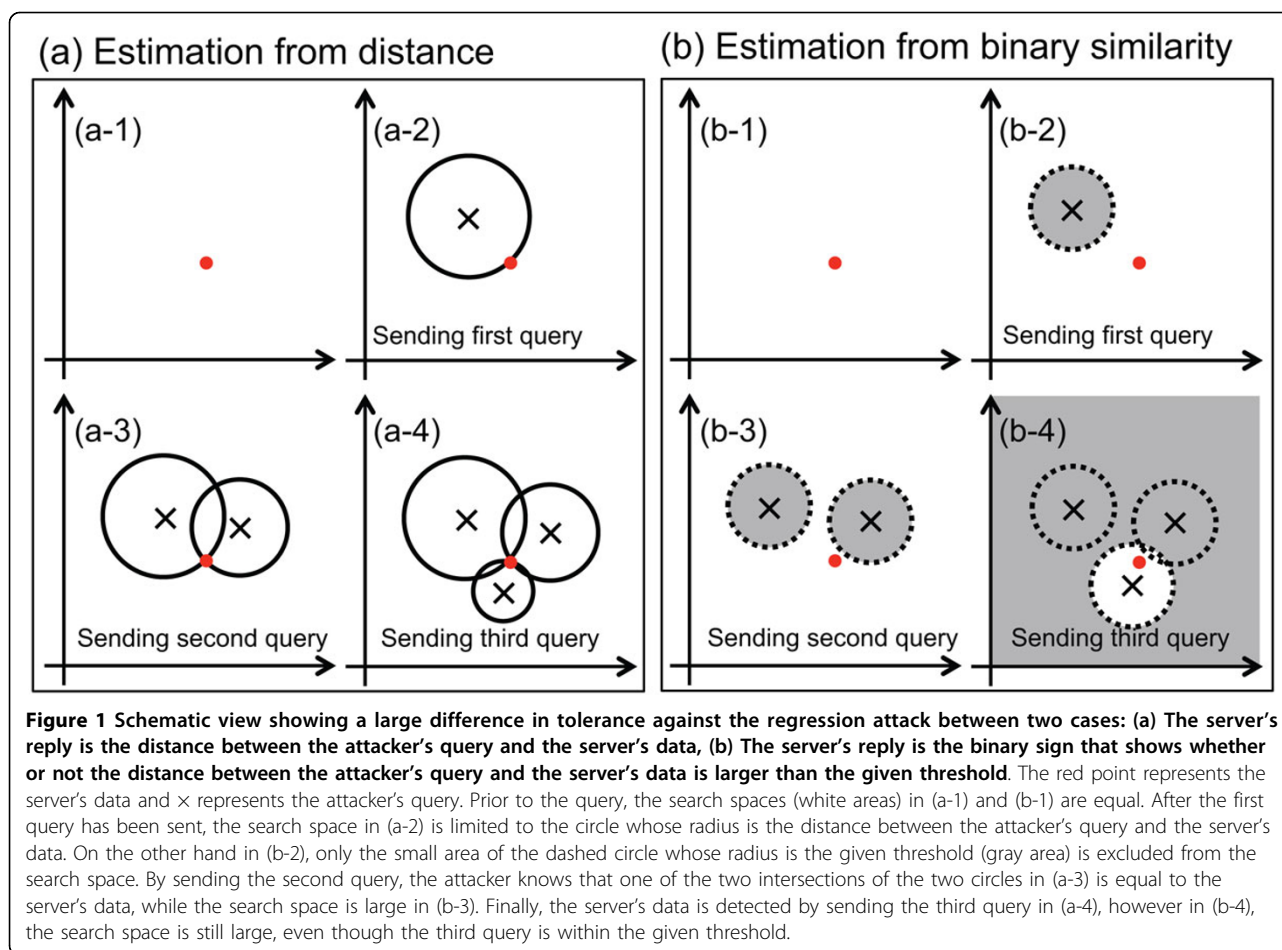
- Illegal query attack

Searching with an illegal query often causes unexpected server behaviour. In such a case, the server might return unexpected results that include important server information. To prevent this, the server should ensure the correctness of the user’s query.

A schematic view of the privacy-preserving database search problems discussed here is shown in Figure 2.

In the field of cryptography, there have been studies of versatile techniques such as general purpose multi-party computation (GP-MPC) [8] and fully homomorphic encryption (FHE) [9], which enable the design of systems that maintain both user privacy and database privacy. However, these techniques require huge computational costs as well as intensive communications between the parties (see the recent performance evaluation of FHE [10]), so they are scarcely used in practical applications. In order to avoid using such techniques, a similarity search protocol using a trusted third party [11] and a privacy preserving SQL database using a trusted proxy server [12] have been proposed, but those methods assure privacy only when the third party does not collude with the user or the server, which is not convenient for many real problems. As far as we know, no practical method has been proposed despite the great importance of privacy-preserving similarity searching. To overcome this lack, we propose a novel privacy-preserving similarity search method that can strongly protect database privacy as well as user privacy while keeping a significantly low computational cost and small communication size.

The rest of this paper is organized as follows. In the next section, we summarize our achievements in this study. This is followed by the Cryptographic background section and the Method section, where we define the problem and introduce details of the proposed protocol.



**Figure 1** Schematic view showing a large difference in tolerance against the regression attack between two cases: (a) The server's reply is the distance between the attacker's query and the server's data, (b) The server's reply is the binary sign that shows whether or not the distance between the attacker's query and the server's data is larger than the given threshold. The red point represents the server's data and x represents the attacker's query. Prior to the query, the search spaces (white areas) in (a-1) and (b-1) are equal. After the first query has been sent, the search space in (a-2) is limited to the circle whose radius is the distance between the attacker's query and the server's data. On the other hand in (b-2), only the small area of the dashed circle whose radius is the given threshold (gray area) is excluded from the search space. By sending the second query, the attacker knows that one of the two intersections of the two circles in (a-3) is equal to the server's data, while the search space is large in (b-3). Finally, the server's data is detected by sending the third query in (a-4), however in (b-4), the search space is still large, even though the third query is within the given threshold.

In the Security analyses section, both the user privacy and database privacy of the proposed protocol are discussed in detail. In the Performance evaluation section, the central processing unit (CPU) time and communication size of the proposed protocol are evaluated for two datasets extracted from ChEMBL. Finally, we present our conclusions for this study in the Conclusion section.

### Our Achievements

Here we focus on similarity search with the Tversky index of fingerprints, which is the most popular approach for chemical compound searches [13] and is used for various search problems in bioinformatics. To provide a concrete application, we address the problem of counting the number of similar compounds in a database, which solves various problems appearing in chemical compound searches. The following model describes the proposed method.

**Model 1** *The user is a private chemical compound holder, and the server is a private database holder. The user learns nothing but the number of similar compounds in the server's database, and the server learns nothing about the user's query compound.*

Here we introduce only a small fraction of the many scientific or industrial problems solved by Model 1.

1 Secure pre-purchase inspection service for chemical compound.

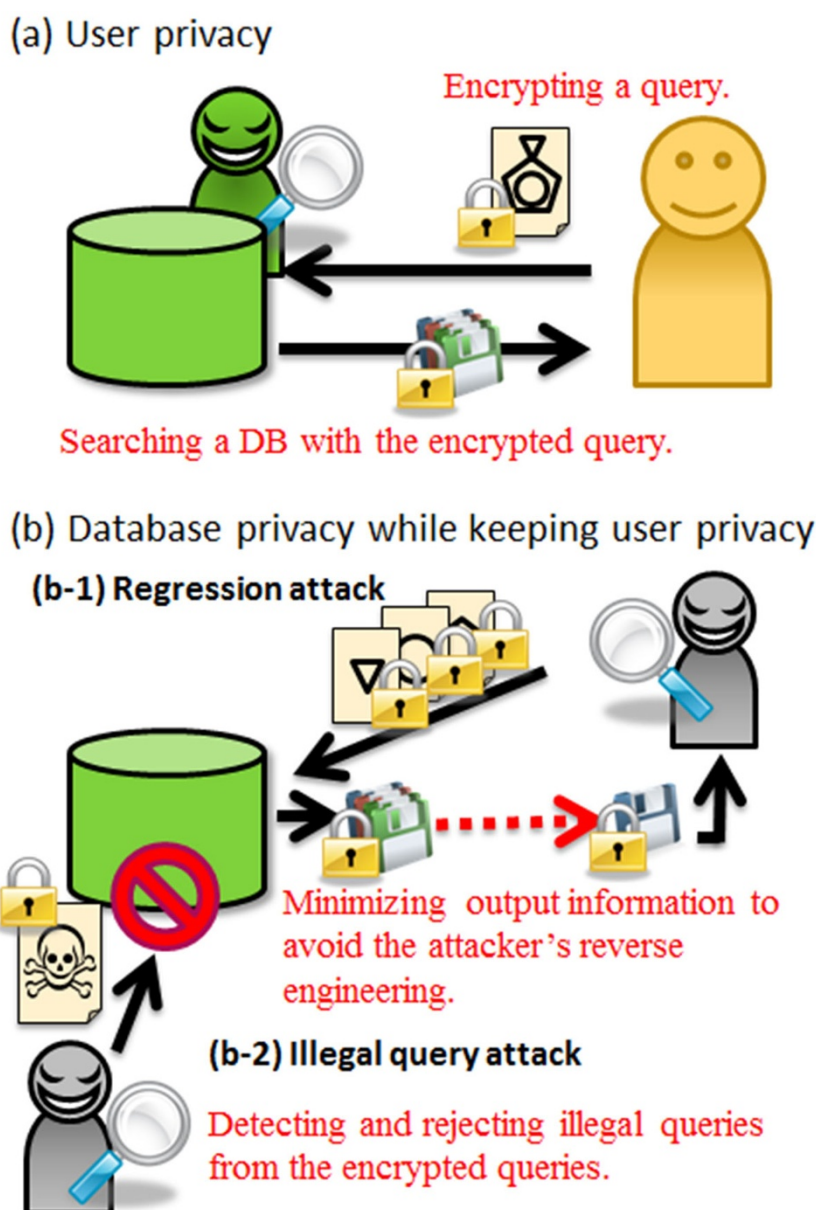
When a client considers the purchase of a commercial database such as a focused library [14], he/she wants to check whether the database includes a sufficient number of similar compounds, without sending his/her private query, but the server does not allow downloading of the database.

2 Secure patent compound filtering.

When a client finds a new compound, he/she usually wants to know whether it infringes on competitors' patents by searching the database of patent-protected compounds maintained by third parties. The same problem occurs when the client wants to check whether or not the compound is undesirable.

3 Secure negative results check.

It is a common perception that current scientific publication is strongly biased against negative results [3], although a recent study showed statistically that negative results brought meaningful benefit [15]. Since researchers are reluctant to provide negative results, which often



**Figure 2 Schematic view of protection of (a) user privacy and (b) database privacy while keeping user privacy.** For user privacy, the user's query and the search result which includes the query information must be invisible to the database side during the search task. For database privacy, the server minimizes output information for preventing regression attacks (b-1), and also detects and rejects illegal queries that might cause unexpected information leakage (b-2). These server's tasks must be carried out with the encrypted queries in order to keep user privacy.

include sensitive information, a privacy-preserving system for sharing those results would greatly contribute to reducing redundant efforts for similar research topics. For example, it would be useful to have a system that allows a user to check whether the query is similar to failed compounds that have previously been examined in other laboratories.

In this study, we propose a novel protocol called the **secure similar compounds counter (SSCC)** which

achieves Model 1. The first main achievement of this study is that SSCC is remarkably tolerant against regression attacks compared with existing protocols which directly output the similarity score. Moreover, we propose an efficient method for protecting the database from illegal query attacks. These points are discussed in the Security analyses section.

The second main achievement is that SSCC is significantly efficient both in computational cost and

communication size. We carefully designed the protocol such that it uses only an additive-homomorphic cryptosystem, which is computationally efficient, and does not rely on any time-consuming cryptographic methods such as GP-MPC or FHE. Hence the performance of the protocol is sufficiently high for a large-scale database such as ChEMBL [16], as is shown in the Performance evaluation section.

## Cryptographic background

### Additively homomorphic encryption scheme

In this paper, we use an additive-homomorphic cryptosystem to design our protocol. The key feature of the additive-homomorphic cryptosystem is that it enables to perform additive operations on *encrypted* values. Therefore, intuitively, any standard computation algorithm can be converted into the privacy-preserving computation algorithm, if operations used in the standard algorithm can be replaced by additions.

More formally, we use a public-key encryption scheme (KeyGen; Enc; Dec), which is semantically secure; that is, an encryption result (ciphertext) leaks no information about the original message (plaintext) [17]. Here, KeyGen is a key generation algorithm for selecting a pair (pk, sk) of a public key pk and a secret key sk; Enc( $m$ ) denotes a ciphertext obtained by encrypting message  $m$  under the given pk; and Dec( $c$ ) denotes the decryption result of ciphertext  $c$  under the given sk. We also require the following additive-homomorphic properties:

- Given two ciphertexts Enc( $m_1$ ) and Enc( $m_2$ ) of messages  $m_1$  and  $m_2$ , Enc( $m_1 + m_2$ ) can be computed without knowing  $m_1$ ,  $m_2$  and the secret key (denoted by Enc( $m_1$ )  $\oplus$  Enc( $m_2$ )).
- Given a ciphertext Enc( $m$ ) of a message  $m$  and an integer  $e$ , Enc( $e \oplus m$ ) can be computed without knowing  $m$  and the secret key (denoted by  $e \otimes$  Enc( $m$ )).

For example, we can use either the Paillier cryptosystem [18] or the “lifted” version of the ElGamal cryptosystem [19] as such an encryption scheme; now the second operation  $\otimes$  can be achieved by repeating the first operation  $\oplus$ . We notice that the range of plaintexts for those cryptosystems can be naturally set as an integer interval  $[-N_1, N_2]$  for some sufficiently large  $N_1, N_2 > 0$ ; therefore, the plaintexts are divided into positive ones, negative ones, and zero.

### Non-interactive zero-knowledge proof

Below, we discuss the following situation: A user (a prover) wants to make a server (a verifier) convinced that a ciphertext  $c$  generated by the user corresponds to a message  $m$  in  $\{0, 1\}$ , but does not want to reveal any information about which of 0 and 1 is  $m$ . This can be achieved by

using a cryptographic tool called *non-interactive zero-knowledge (NIZK) proof*. In the present case, it enables the user to generate a “proof” associated with  $c$ , so that:

- If  $m$  is indeed in  $\{0, 1\}$ , then the server can verify this fact by testing the proof (without knowing  $m$  itself).
- If  $m \notin \{0, 1\}$ , then the user cannot generate a proof that passes the server’s test.
- The server cannot obtain any information about  $m$  from the proof, except for the fact that  $m \in \{0, 1\}$ .

(See [20] for a general formulation.) Besides the existing general-purpose NIZK proofs, Sakai et al. [21] proposed an efficient scheme specific to the “lifted” ElGamal cryptosystem, which we use below. (See Section 1 of Additional File 1 in which we give the brief description of the NIZK proofs [21].)

## Method

The goal of this study is to design a protocol between a user and a server that enables the user to obtain the number of compounds in the server’s database that are similar to the user’s target compound. Here, a fingerprint of compound is modeled as  $\vec{p} \in \{0, 1\}^\ell$  (i.e., a bit string of length  $\ell$ ). An equivalent way to refer to  $\vec{p}$  is the set of all indices  $i$  where  $p_i = 1$ . We denote such a set by  $\mathbf{p}$ . The similarity of two compounds  $\mathbf{p}, \mathbf{q}$  is then measured by *Tversky index* which is parameterized by  $\alpha, \beta > 0$  and is defined as:

$$TI_{\alpha, \beta}(\mathbf{p}, \mathbf{q}) = \frac{|\mathbf{p} \cap \mathbf{q}|}{|\mathbf{p} \cap \mathbf{q}| + \alpha |\mathbf{p} \setminus \mathbf{q}| + \beta |\mathbf{q} \setminus \mathbf{p}|}.$$

Tversky index is useful since it includes several important similarity measurements such as Jaccard Index (JI, which is exactly  $TI_{1,1}$  and also known as Tanimoto Index) and Dice index (which is exactly  $TI_{1/2,1/2}$ ) [22]. First, we introduce the basic idea and two efficient techniques for improving database privacy. Then, we describe our full proposed protocol.

### Basic idea

We firstly consider the simplest case that the user has (the fingerprint of) a target compound  $\mathbf{q}$  as a query and the server’s database consists of only a single fingerprint  $\mathbf{p}$ . The case of a larger database is discussed later. The goal here is to detect whether or not the Tversky index of  $\mathbf{p}$  and  $\mathbf{q}$  is larger than a given threshold  $1 \geq \theta > 0$ . The main idea of our approach is to calculate the score

$$\theta^{-1}(|\mathbf{p} \cap \mathbf{q}|) - (|\mathbf{p} \cap \mathbf{q}| + \alpha |\mathbf{p} \setminus \mathbf{q}| + \beta |\mathbf{q} \setminus \mathbf{p}|) \quad (1)$$

from *encrypted* fingerprints  $\mathbf{p}$  and  $\mathbf{q}$  by an additive-homomorphic cryptosystem. The score is non-negative if

and only if the Tversky index of  $\mathbf{p}$  and  $\mathbf{q}$  is at least  $\theta$ . Now since  $|\mathbf{p} \setminus \mathbf{q}| = |\mathbf{p}| - |\mathbf{p} \cap \mathbf{q}|$  and a similar relation holds for  $|\mathbf{q} \setminus \mathbf{p}|$ , the score (1) is positively proportional to

$$\lambda_1 |\mathbf{p} \cap \mathbf{q}| - \lambda_2 |\mathbf{p}| - \lambda_3 |\mathbf{q}|,$$

where  $\lambda_1 = c(\theta^{-1} - 1 + \alpha + \beta)$ ,  $\lambda_2 = c\alpha$ ,  $\lambda_3 = c\beta$  and any positive value  $c$ . We assume that the parameters and the threshold for the Tversky index are rational numbers denoted by  $\alpha = \mu_a/\gamma$ ,  $\beta = \mu_b/\gamma$  and  $\theta = \theta_n/\theta_d$ , where  $\mu_a$ ,  $\mu_b$ ,  $\gamma$ ,  $\theta_n$  and  $\theta_d$  are non-negative integers. By using  $c = \gamma\theta_n g^{-1}$  under this assumption,  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  become non-negative integers where  $g$  is the greatest common divisor of  $\gamma(\theta_d - \theta_n) + \theta_n(\mu_a + \mu_b)$ ,  $\theta_n\mu_a$  and  $\theta_n\mu_b$ .

Motivated by this observation, we define the following modified score, called the **threshold Tversky index**:

**Definition 1** Given parameters  $\alpha$  and  $\beta$  and a threshold  $\theta$  for the Tversky index which are rational numbers denoted by  $\alpha = \mu_a/\gamma$ ,  $\beta = \mu_b/\gamma$  and  $\theta = \theta_n/\theta_d$  where  $\mu_a$ ,  $\mu_b$ ,  $\gamma$ ,  $\theta_n$  and  $\theta_d$  are non-negative integers, then the threshold Tversky index  $\overline{\text{TI}}_{\alpha,\beta,\theta} = \overline{\text{TI}}_{\alpha,\beta,\theta}(\mathbf{p} \cap \mathbf{q})$  for fingerprints  $\mathbf{p}$  and  $\mathbf{q}$  is defined by

$$\overline{\text{TI}}_{\alpha,\beta,\theta} := \lambda_1 |\mathbf{p} \cap \mathbf{q}| - \lambda_2 |\mathbf{p}| - \lambda_3 |\mathbf{q}|,$$

and non-negative integer parameters  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are defined by

$$\begin{aligned} \lambda_1 &= \gamma\theta_n g^{-1}(\theta^{-1} - 1 + \alpha + \beta), \\ \lambda_2 &= \gamma\theta_n g^{-1}\alpha, \\ \lambda_3 &= \gamma\theta_n g^{-1}\beta, \end{aligned}$$

where  $g$  is the greatest common divisor of  $\gamma(\theta_d - \theta_n) + \theta_n(\mu_a + \mu_b)$ ,  $\theta_n\mu_a$  and  $\theta_n\mu_b$ .

By the above argument, we have  $\text{TI}_{\alpha,\beta}(\mathbf{p}, \mathbf{q}) \geq \theta$  if and only if  $\overline{\text{TI}}_{\alpha,\beta,\theta}(\mathbf{p}, \mathbf{q}) \geq 0$ . Therefore, the user can know whether or not his/her target compound  $\mathbf{q}$  is similar (i.e.,  $\text{TI}_{\alpha,\beta}(\mathbf{p}, \mathbf{q}) \geq \theta$ ) to the fingerprint  $\mathbf{p}$  in the database, by obtaining only the value  $\overline{\text{TI}}_{\alpha,\beta,\theta}(\mathbf{p}, \mathbf{q})$ .

In the protocol, the bits of the user's target fingerprint  $\mathbf{q}$  and the value  $|\mathbf{p}|$  held by the server are both encrypted using the user's public key. Since  $\overline{\text{TI}}_{\alpha,\beta,\theta}(\mathbf{p}, \mathbf{q})$  can be computed by the addition of these values and multiplication by integers, the protocol can calculate (without the secret key) a ciphertext of  $\overline{\text{TI}}_{\alpha,\beta,\theta}(\mathbf{p}, \mathbf{q})$ , which is then decrypted by the user. For simplicity, we will abuse the notation and write  $\text{TI}(\mathbf{p}, \mathbf{q})$ ,  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$  without subscripts  $\alpha$ ,  $\beta$ ,  $\theta$  when the context is clear.

We emphasize that our protocol does not use time-consuming cryptographic methods such as GP-MPC and FHE, and data transfer occurs only twice during an execution of the protocol. Hence, our protocol is efficient enough to scale to large databases.

### Database security enhancement techniques against regression attack

As discussed in Introduction section, the server needs to minimize returned information in order to minimize the success ratio of the regression attack. That is, the ideal situation for the server is that the user learns only the similarity/non-similarity property of fingerprints  $\mathbf{p}$  and  $\mathbf{q}$ , without knowing any other information about the secret fingerprint  $\mathbf{p}$ . This means that only the sign of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$  should be known by the user. However, in our basic protocol, the value of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$  is fully obtained by the user; Database privacy is not protected from regression attacks. (See the Security analyses section for details.) In order to send only the sign of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$ , we firstly considered using a bit-wise decomposition protocol [23] for extracting and sending only the sign bit of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$ . Although this approach is ideal in terms of security, the protocol requires more than 30 rounds of communications, which is much more efficient than using GP-MPC or FHE, but rather time-consuming for large-scale databases. Therefore, here we propose the novel technique of using dummy replies, which requires only one round of communication while sufficiently minimizing information leakage of  $\mathbf{p}$ . In the proposed technique, besides its original reply  $t = \text{Enc}(\overline{\text{TI}}(\mathbf{p}, \mathbf{q}))$ , the server also chooses random integers  $\phi_1, \dots, \phi_n$  from a suitable interval and encrypts those values under the user's public key  $\text{pk}$ . Then the server sends the user a collection of ciphertexts  $t, \text{Enc}(\phi_1), \dots, \text{Enc}(\phi_n)$  that are shuffled to conceal the true ciphertext  $t$ , as well as the number  $s_d$  of dummy values  $\phi_k$  with  $\phi_k \geq 0$ . The user decrypts the received  $n + 1$  ciphertexts, counts the number  $s_c$  of non-negative values among the decryption results, and compares  $s_c$  to  $s_d$ . Now we have  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q}) \geq 0$  if and only if  $s_c - s_d = 1$ ; therefore, the user can still learn the sign of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$ , while the actual value of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$  is concealed by the dummies. We have confirmed that the information leakage of  $\mathbf{p}$  approaches zero as the number of dummies becomes large; see the Security analyses for pudding dummies section for more detailed discussion. (We have also developed another security enhancement technique using sign-preserving randomization of  $\overline{\text{TI}}(\mathbf{p}, \mathbf{q})$ ; see Section 2 of Additional File 1 for details.)

Database security enhancement technique against illegal query attack Illegal query attacks can be prevented if the server can detect whether or not the user's query is valid. To keep user privacy, the server must conduct this task without obtaining more information than the validity/invalidity of the query. In fact, this functionality can be implemented by using the NIZK proof by Sakai et al. [21] mentioned in the Non-interactive zero-knowledge proof

section. The improved protocol requires the user to send the server a proof associated with the encrypted fingerprint bits  $q_i$ , from which the server can check whether  $q$  is indeed a valid fingerprint (without obtaining any other information about  $q$ ); the server aborts the protocol if  $q$  is invalid. Here we use the “lifted” ElGamal cryptosystem as our basic encryption scheme to apply Sakai’s scheme. (We note that if we require the user to send  $\text{Enc}(-|q|)$  used by server’s computation, then another NIZK proof is necessary to guarantee the validity of the additional ciphertext, which decreases the communication efficiency of our protocol. Hence our protocol requires the server to calculate  $\text{Enc}(-|q|)$  by itself.) The formal definition of the valid query is given in the Database privacy in malicious model section.

### Secure similar compounds counter

For the general case that the database consists of more than one fingerprint  $p$ , we propose the protocol shown in Algorithm 1 to count the number of fingerprints  $p$  similar to the target fingerprint  $q$ . In the protocol, the server simply calculates the encryption of the threshold Tversky indices for all database entries and, as discussed above, replies with a shuffled collection of these true ciphertexts and dummy ciphertexts, as well as the number  $s_d$  of non-negative dummy values. Then the value  $s_c - s_d$  finally obtained by the user is equal to the number of similar fingerprints  $p$  in the database.

**Algorithm 1** The secure similar compounds counter (SSCC)

- Public input: Length of fingerprints  $\ell$  and parameters for the Tversky index  $\theta = \theta_n/\theta_d$ ,  $\alpha = \mu_a/\gamma$ ,  $\beta = \mu_b/\gamma$
- Private input of a user: Target fingerprint  $q$
- Private input of a server: Set of fingerprints  $P = \{p^{(1)}, \dots, p^{(M)}\}$

1 (Key setup of cryptosystem) The user generates a key pair  $(pk, sk)$  by the key generation algorithm KeyGen for the additive-homomorphic cryptosystem and sends public key  $pk$  to the server (the user and the server share public key  $pk$  and only the user knows secret key  $sk$ ).

2 (Initialization) The user encrypts his/her fingerprint  $q$  as a vector of ciphertexts:  $\text{Enc}(q_k) := (\text{Enc}(q_1), \dots, \text{Enc}(q_\ell))$ . He/she also generates  $v$  as a vector of proofs. Each proof  $v_i$  is associated with  $\text{Enc}(q_i)$ .

3 (Query of entry) The user sends the vector of ciphertexts  $\text{Enc}(q_k)$  and the vector of proofs  $v$  to the server as a query.

4 (Query validity verification) The server verifies the validity of  $\text{Enc}(q_k)$  by testing the vector of proof  $v$ . If  $v$  does not pass the server’s test, the user cannot move on to the next step.

5 (Calculation of threshold Tversky index)

(a) The server calculates the greatest common divisor of  $\gamma(\theta_d - \theta_n) + \theta_n(\mu_a + \mu_b)$ ,  $\theta_n\mu_a$  and  $\theta_n\mu_b$  as  $g$ , and calculates  $\lambda_1 = \gamma\theta_n g^{-1}(\theta^{-1} - 1 + \alpha + \beta)$ ,  $\lambda_2 = \gamma\theta_n g^{-1}\alpha$ , and  $\lambda_3 = \gamma\theta_n g^{-1}\beta$ .

(b) The server calculates  $\text{Enc}(-|q|) = \text{Enc}\left(-\sum_{i=1}^{\ell} q_i\right)$

from  $\text{Enc}(q_k) : \text{Enc}(-|q|) = -1 \otimes \oplus_{i=1}^{\ell} \text{Enc}(q_i)$ .

(c) for  $j = 1$  to  $M$  do

i. The server calculates  $-|p^{(j)}| = -\sum_{i=1}^{\ell} p_i^{(j)}$  and encrypts it to obtain a ciphertext  $\text{Enc}(-|p^{(j)}|)$ .

ii. The server calculates a ciphertext  $t_j$  of threshold Tversky index  $\overline{\text{TI}}(p^{(j)}, q)$ .

$c \leftarrow \text{Enc}(0)$

for  $k = 1$  to  $\ell$  do

if  $p_k^{(j)} = 1$

$c \leftarrow c \oplus \text{Enc}(q_k) \triangleright$  Computing  $\text{Enc}(|p^{(j)} \cap q|)$

end if

end for

$t_j \leftarrow \lambda_1 \otimes c \oplus \lambda_2 \text{Enc}(-|p^{(j)}|) \oplus \lambda_3 \otimes \text{Enc}(-|q|)$

end for

6 (Padding of dummies)

(a) The server generates a set of dummy values  $\{\phi_1, \dots, \phi_n\}$  and counts the number  $s_d$  of non-negative dummies  $\phi_i \geq 0$ .

(b) The server encrypts  $\phi_i$  to obtain a ciphertext  $\text{Enc}(\phi_i)$  for  $i = 1, \dots, n$ .

(c) The server shuffles the contents of the set  $T = \{t_1, \dots, t_M, \text{Enc}(\phi_1), \dots, \text{Enc}(\phi_n)\}$ .

7 (Return of matching results) The server sends  $T$  and  $s_d$  to the user.

8 (Decryption and counting) The user decrypts the contents of  $T$  and counts the number  $s_c$  of non-negative values.

9 (Evaluation) The user obtains  $s_c - s_d$  as the number of similar fingerprints in the database.

### Parameter settings of the protocol

Decrypting an encryption of too large value needs huge computation cost if the lifted-ElGamal cryptosystem is used. Therefore, in order to keep the consistency and efficiency of the protocol, the range of  $\overline{\text{TI}}(p, q)$  should not be too large. i.e., the integer parameters  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  in the threshold Tversky index should not be too large. In fact, this will not cause a problem in practice; For example, the parameters become  $\lambda_1 = 9$ ,  $\lambda_2 = \lambda_3 = 4$  for computing  $\overline{\text{TI}}_{1, 1, 0.8}$  which is a typical setting of a chemical compound search. In this case, a minimum value and a maximum value of  $\overline{\text{TI}}(p, q)$  is -664 and 166 for 166 MACCS keys, which is a sufficiently small range. (See Section 3 of Additional File 1 for details.)

## Security analyses

In this section, we evaluate security of SSCC by several approaches.

In the area of cryptology, the following two standard security models for two-party computation have been considered:

- *Semi-honest model* : Both parties follow the protocol, but an adversarial one attempts to infer additional information about the other party's secret input from the legally obtained information.
- *Malicious model* : An adversarial party cheats even in the protocol (e.g., by inputting maliciously chosen invalid values) in order to illegally obtaining additional information about the secret.

We analyze user privacy and database privacy in both the semi-honest and malicious models. For the database privacy, we firstly compare attack success ratios for the case of using our method which aims to output a binary sign and the other case of using the previous methods which aim to output a similarity score, and show that outputting a binary sign improves database privacy. We also evaluate security strength of our method against a regression attack by comparing attack success ratios for the case of using dummies and the ideal case that uses a versatile technique (such as GP-MPC and FHE) to output a binary sign, and show that the security strength for the case of using dummies is almost the same as the ideal case under realistic settings.

### User privacy

The semantic security of the encryption scheme used in the protocol (see the Additively homomorphic encryption scheme section) implies immediately that the server cannot infer any information about the user's target fingerprint  $q$  during the protocol. This holds in both the semi-honest and malicious models.

### Thresholding largely improves database privacy

We mentioned in the introduction section that minimizing information returned from the server reduces success ratio of regression attack. Therefore, SSCC aims for "ideal" case in which the user learns only the sign of  $\overline{\text{TI}}(p, q)$  during the protocol. The previous methods that compute Jaccard Index aim for the "plain" case, in which the user fully learns the value  $\text{TI}(p, q)$ . Here we evaluate the efficiency of the thresholding by comparing success probabilities of regression attack for those two cases. We consider the general case in which the user is allowed to send more than one query and those queries are searched by Jaccard Index. We also suppose that the database consists of a single fingerprint  $p$  in order to clarify the effect of thresholding.

The goal of an attacker is to reveal  $p$  by analysing the results returned from the server. It is generally effective for the attacker to exploit the difference between the two outputs obtained by sending two different queries. In fact, when the server returns  $\overline{\text{TI}}$ ,  $\overline{\text{TI}}(p, q) - \overline{\text{TI}}(p, 0)$  becomes positive if and only if  $p_i = 1$ , where  $q = (0, \dots, q_i = 1, \dots, 0)$  and  $0 = (0, \dots, 0)$ . This means that the attacker can reveal any bit in  $p$  by sending the single query after sending the first query  $0$ . Therefore,  $p$  can be fully revealed by sending only  $\ell + 1$  queries. On the other hand, there is no deterministic attack for revealing  $p$  from only the sign of  $\overline{\text{TI}}$ , because two different inputs do not always lead to different outputs. Since we know of a linear algorithm that fully reveals  $p$  in response to at most  $2\ell$  queries after making a "hit" query  $q$  such that  $\overline{\text{TI}}(p, q) > 0$ , here we evaluate database privacy by the probability of making at least one hit query when the user is allowed to send  $x$  queries. (See Section 4 of Additional File 1 for details.) This probability is denoted as

$$\sum_p \Pr(X = p) \cdot (1 - (1 - f_p)^x), \quad (2)$$

where  $f_p$ , defined as follows, is the probability that the user makes one hit query with a single trial when  $p$  is given.

$$f_p := \sum_q \Pr(Y = p) \cdot \Pr(\overline{\text{TI}}(p, Y) > 0 | Y = q).$$

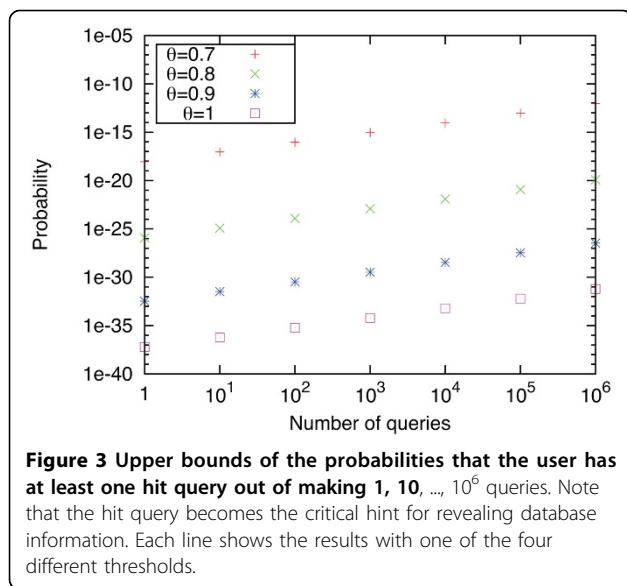
For ease of calculation, we computed the upper bound of equation (2) for  $x = 1, 10, 10^2, \dots, 10^6$  and  $\theta = 0.7, 0.8, 0.9, 1.0$ . (See Section 5 of Additional File 1 for details.) Since publicly available 166 MACCS keys are the most popular fingerprint for chemical compound searches, we set  $\ell$  to 166. From the results shown in Figure 3, we can see that the probability of making a hit query is sufficiently small for practical use even though the user is allowed to send a million of queries. Considering that the user learns  $p$  by using no more than  $\ell + 1$  queries when he/she learns  $\overline{\text{TI}}$ , we can conclude that database privacy is dramatically improved by thresholding. In other words, the proposed protocol, which aims to output only the sign of the similarity score, has stronger security than other previous methods, which directly output similarity scores.

### Security analyses for padding dummies

We showed that the output privacy in the "ideal" case is significantly improved from the "plain" case. Here we experimentally evaluate how the actual situation of our proposed protocol is close to the "ideal" case.

Before going into detail analyses, let us discuss how to generate dummies. It is ideal for the server privacy to generate a dummy according to the same distribution where  $\overline{\text{TI}}(p, q)$  is generated from. However, this is not realistic





because  $\overline{\Pi}(p, q)$  is determined by both  $p$  and  $q$  which is user's private information. Therefore, in our analyses, we assume that a dummy is generated from uniform distribution over possible values of  $\overline{\Pi}(p, q)$ . For example, if possible values of  $\overline{\Pi}(p, q)$  is  $\{1, 2, 3, 4, 5\}$ , dummies are randomly selected from any one of them. The purpose of padding dummies is to mitigate the risk of leaking  $\overline{\Pi}(p, q)$ . In order to clarify the effect of the use of dummy values, we concentrate on the basic case; the database contains a single  $p$ , and there exist  $k$  possible values of  $\overline{\Pi}(p, q)$ .  $i$ -th value of the  $k$  possible values arises as the true  $\overline{\Pi}(p, q)$  according to the probability  $w_i$ . Namely, true  $\overline{\Pi}(p, q)$  is generated from the multinomial distribution with  $k$  different probabilities  $w = w_1, \dots, w_k$ , while dummies are generated from the multinomial distribution with equal probability  $1/k$ . To conduct stringent analyses, we assume that the user knows  $w$ , and he/she also knows that dummies are uniformly distributed over  $k$  possible  $\overline{\Pi}(p, q)$ .

The security provided by our protocol can be formalized in the following manner. First we recall that, in our protocol, the server computes encryption of  $\overline{\Pi}(p, q)$  and encryption of dummy values  $\phi_1, \dots, \phi_m$  and then sends the user the  $n+1$  encrypted values as well as the number of positive dummy values in  $\phi_1, \dots, \phi_n$ . For the purpose of formalizing the security, we introduce a "fictional" server that performs the following: It first receives the encrypted values  $\overline{\Pi}(p, q), \phi_1, \dots, \phi_n$  from the real server. Secondly, it gets the sign of  $\overline{\Pi}(p, q)$ . (We note that a real server cannot do it since it requires unrealistic computational power that breaks the security of the encryption scheme, so this is just fictional for the sake of mathematical definition.)

Thirdly, it generates another dummy value  $\overline{\Pi}'$  randomly, and independently of the values of  $\overline{\Pi}(p, q), \phi_1, \dots, \phi_n$  (except for the sign of  $\overline{\Pi}(p, q)$ ), in the following manner:

- If  $\overline{\Pi}(p, q)$  is positive, then  $\overline{\Pi}'$  is chosen randomly from positive values.
- If  $\overline{\Pi}(p, q)$  is negative, then  $\overline{\Pi}'$  is chosen randomly from negative values.

Finally, the fictional server sends the user an encryption of  $\overline{\Pi}'$  (instead of  $\overline{\Pi}(p, q)$ ) as well as the encrypted  $\phi_1, \dots, \phi_n$  and the number of positive values in  $\phi_1, \dots, \phi_n$ . We note that, when the user receives a reply from the fictional server, the user can know the sign of  $\overline{\Pi}(p, q)$  which is the same as that of  $\overline{\Pi}'$ , but cannot know any other information on  $\overline{\Pi}(p, q)$  since  $\overline{\Pi}'$  is independent of  $\overline{\Pi}(p, q)$ . In the setting, the following property can be proven:

**Theorem 1** Suppose that the user cannot distinguish, within computational time TIME, the sets of decrypted values of ciphertexts involved in outputs of the real server and of the fictional server. Then any information computable within computational time TIME from the decryption results for output of the real server is equivalent to information computable within computational time TIME' from the sign of  $\overline{\Pi}(p, q)$  only, where TIME' is a value which is close to TIME.

*Proof* Let  $\mathcal{A}$  be an algorithm, with running time TIME, which outputs some information on the decrypted values for an output of the real server. We construct an algorithm  $\mathcal{A}'$  which computes, from the sign of  $\overline{\Pi}(p, q)$  only, an information equivalent to the information computed by  $\mathcal{A}$ . The construction is as follows; from the sign of  $\overline{\Pi}(p, q)$ ,  $\mathcal{A}'$  generates dummy values by mimicking the behavior of the fictional server, and then  $\mathcal{A}'$  inputs these dummy values to a copy of  $\mathcal{A}$ , say  $\mathcal{A}^*$ , and gets the output of  $\mathcal{A}^*$ . Now if the output of  $\mathcal{A}'$  is not equivalent to the output of  $\mathcal{A}$ , then the definition of  $\mathcal{A}'$  implies that the probability distributions of the outputs of  $\mathcal{A}$  with inputs given by the decrypted values for outputs of the real server and of the fictional server are significantly different (since  $\mathcal{A}^*$  used in  $\mathcal{A}'$  is a copy of  $\mathcal{A}$ ); it enables the user to distinguish the two possibilities of his/her received values by observing the output of  $\mathcal{A}$ , but this contradicts the assumption of the theorem. Therefore, the output of  $\mathcal{A}'$  is equivalent to the output of  $\mathcal{A}$  as claimed. Moreover, the computational overhead of  $\mathcal{A}'$  compared to  $\mathcal{A}$  is just the process of generating dummy values by mimicking the behavior of the fictional server; it is not large (i. e., TIME' is close to TIME as claimed) since the server-side computation of our proposed protocol is efficient. Hence, the theorem holds.

Roughly rephrasing, if the assumption of the theorem is true for a larger TIME, then the actual situation of our proposed protocol becomes closer to the “ideal” case provided we focus on any information available from efficient computation. As a first step to evaluate how the assumption is plausible (i.e., how the value TIME in the assumption can be large), we performed computer experiments to show that some natural attempts to distinguish the actual and the fictional cases do not succeed, as explained below.

In this experiment, we evaluate the security of our protocol by comparing the probabilities that the user correctly guesses the value  $\overline{\text{TI}}(p, q)$  in two cases: The case in which the user makes a guess based only on a prior knowledge  $w$ , and the other case in which the user makes a guess based on the observation of the search result under the condition that he/she knows  $w$ .

For the first case, the user’s best strategy for guessing  $\overline{\text{TI}}(p, q)$  is to choose the  $i_0$ -th possible value, where

$$i_0 = \arg \max_{1 \leq i \leq k} w_i. \quad (3)$$

In this case, the success probability of the guess is  $w_{i_0}$ .

Let us consider the best strategy for the second case. As described above, we consider a practical case that  $n$  dummy values  $\phi_1, \dots, \phi_n$  chosen from the  $k$  possible values uniformly at random, and the user makes a guess from the received  $n + 1$  shuffled values  $\phi_1, \dots, \phi_n, \overline{\text{TI}}(p, q)$ . Now suppose that the user received the  $i$ -th possible value  $a_i$  times for each  $1 \leq i \leq k$  (hence  $\sum_{i=1}^k a_i = n + 1$ ). Since the choices of  $\phi_1, \dots, \phi_n$  are independent of  $\overline{\text{TI}}(p, q)$ , the probability that the user received  $i$ -th possible value  $a_i$  times for each  $1 \leq i \leq k$  and that  $\overline{\text{TI}}(p, q)$  is  $i_0$ -th possible value is

$$\binom{n}{a_1, \dots, a_{i_0} - 1, \dots, a_k} \left(\frac{1}{k}\right)^n \cdot w_{i_0} = a_{i_0} \cdot w_{i_0} \frac{n!}{a_1! \dots a_k! k^n}.$$

Therefore, the conditional probability that  $\overline{\text{TI}}(p, q)$  is the  $i_0$ -th possible value, conditioned on the set of the user’s received values, is

$$\frac{\binom{n}{a_1, \dots, a_{i_0} - 1, \dots, a_k} \left(\frac{1}{k}\right)^n \cdot w_{i_0}}{\sum_{i=1}^k \binom{n}{a_1, \dots, a_i - 1, \dots, a_k} \left(\frac{1}{k}\right)^n \cdot w_i} = \frac{a_{i_0} \cdot w_{i_0}}{\sum_{i=1}^k a_i \cdot w_i}.$$

This implies that the user’s best strategy is to guess that  $\overline{\text{TI}}(p, q)$  is the  $i_0$ -th possible value, where

$$i_0 = \arg \max_{1 \leq i \leq k} a_i \cdot w_i. \quad (4)$$

We estimated success probabilities of user’s guess for the both cases by simulation experiments. Here we

assumed typical case when  $\text{TI}_{1,1,0.8}$  and 166 MACCS keys are used. In this case,  $k = 831$  and we performed the experiments for  $n = 831 \times 10^0, 831 \times 10^1, \dots, 831 \times 10^4$  on three different distributions of  $\overline{\text{TI}}(p, q)$  which were obtained by the following schemes:

1 We randomly selected one fingerprint  $q$  from ChEMBL and calculated  $\overline{\text{TI}}(p, q)$  for all the entries in ChEMBL and used the observed distribution as  $w$ . In our experiment, 177159-th fingerprint was selected as  $q$  (referred as  $w^{\text{ChEMBL-177159}}$ ).

2 The same scheme as 1) was used when  $q$  was 265935-th fingerprint (referred as  $w^{\text{ChEMBL-265935}}$ ).

3 We randomly selected a value from  $1, \dots, k$  for  $m$  times and count frequency of  $i$  as  $h_i$  and set  $w_i = h_i/m$  (referred as  $w^{\text{random}}$ ). We used  $k \times 5$  as  $m$ .

All the distributions used here are shown in Section 6 of Additional File 1.

We performed 100, 000 trials for each experiment. Each trial consisted of choosing  $\phi_1, \dots, \phi_n$  uniformly at random; choosing  $\overline{\text{TI}}(p, q)$  according to  $w$ ; deciding the user’s guess  $i_0$  by formula (3) and formula (4) respectively (we adopted a uniformly random choice if there were more than one such  $i_0$ ); and checking whether or not  $\overline{\text{TI}}(p, q)$  was the  $i_0$ -th possible value for both rules (i.e., the user’s guess succeeded). The results of the experiment are given in Table 1; they show that the user’s attack success probability became significantly close to the ideal case when a sufficiently large number of dummies were used; therefore, our technique of using dummies indeed improves the output privacy.

Security analyses for padding dummies for the case when the user is allowed to send more than one query

One might suspect that the attacker can detect the true  $\overline{\text{TI}}(p, q)$  by sending the same query twice and finding the value which is appeared in both results. However, this attack does not easily succeed if  $n$  is sufficiently larger than  $k$  (i.e., ideally, all possible values of  $\overline{\text{TI}}$  are covered by sufficient number of dummies), and we consider that  $k$  is not too large in practice as we discussed in Parameters settings of the protocol section.

In order to evaluate the security achieved by the padding dummies for the case when the user is allowed to submit  $L$  queries, we performed following analyses. Here we evaluate the security achievement by comparing the case of using our protocol based on the padding dummy and the ideal case of returning only the sign of  $\overline{\text{TI}}(p, q)$ . In order to perform rigorous analyses, we assume the most severe case in which the attacker keeps sending the same query  $L$  times. For this case, the probability that  $\overline{\text{TI}}(p, q)$  is the  $i_0$ -th possible value after sending  $L$

**Table 1 The experimental success ratios of the user's guess based on the server's return and the prior distribution of true value ( $n = 813$ ,**

	$n = 831$	$n = 831 \times 10^1$	$n = 831 \times 10^2$	$n = 831 \times 10^3$	$n = 831 \times 10^4$	Ideal value
$w^{\text{ChEMBL-177159}}$	0.03552	0.01738	0.01101	0.01009	0.00977	0.00981
$w^{\text{ChEMBL-265935}}$	0.02991	0.01337	0.00903	0.00798	0.00784	0.00807
$w^{\text{rand}}$	0.00914	0.0041	0.00309	0.00279	0.00305	0.00289

$\overline{\Pi}_{1, 1, 0.8}$  ( $k = 831$ ) is assumed and results are calculated for five different numbers of dummies ( $n = 831, 831 \times 10^1, 831 \times 10^2, 831 \times 10^3, 831 \times 10^4$ ) are used for three different distributions:  $w^{\text{ChEMBL-177159}}$  and  $w^{\text{ChEMBL-265935}}$  are actual distributions of  $\overline{\Pi}_{1, 1, 0.8}$  on ChEMBL obtained by querying two randomly selected fingerprints from ChEMBL,  $w^{\text{rand}}$  is obtained by randomly selecting a value from  $1, \dots, k$  for  $m = 5 \times 831$  times and dividing each observed frequency by  $m$ .

queries on condition that frequency of  $i$ -th possible value of  $j$ -th query  $a_i^{(j)}$  for  $j = 1, \dots, L$  is

$$\prod_{j=1}^L \frac{a_i^{(j)} \cdot w_i}{\sum_{h=1}^k a_h^{(j)} \cdot w_h} \quad (5)$$

This implies that the user's best strategy is to choose  $i$ -th possible value which maximizes equation (5). As mentioned above, we compared the success ratio of the attack based on the above strategy and the ideal success ratio when the user makes the guess only from the given distribution  $w$ . We also assumed more realistic case that user did not know the exact distribution of dummy but knew the distribution that was similar to the actual distribution the server used. For the evaluation of this case, we generated dummies from the distribution  $u$ , which was slightly different from uniform distribution, while the user assumed that dummies were generated from uniform distribution.  $u$  was generated as follows:

$$u_i = r \cdot 1/k, \text{ where } r \sim N(1, \delta^2).$$

We performed the experiment for  $L = 1, 10, 10^2, \dots, 10^5$ ,  $n = 831 \times 10, 831 \times 50, 831 \times 10^2$  and  $\delta = 0, 0.05, 0.1, 0.15, 0.2$  based on the same approach used in the evaluation of single query security. i.e., for each trial,  $n$  dummies were randomly chosen according to  $u$  (note that  $u$  was equal to uniform distribution when  $\delta = 0$ ), true value  $\overline{\Pi}(p, q)$  was selected according to  $w$  and the attacker's guess was made based on the equation (5). We performed 10, 000 trials for each triplet of  $L, n$  and  $\delta$ . Those experiments were conducted for the same three distributions:  $w^{\text{ChEMBL-177159}}$ ,  $w^{\text{ChEMBL-265935}}$  and  $w^{\text{ChEMBL-random}}$ . We compared the success ratio of the attack and the ideal success ratio when the user made the guess without seeing search results. The results are shown in Figure 4. The success ratio of user's attack decreased as the number of dummies increased and became closer to the ideal value when the sufficient number of dummies are given, even for the case that a large number of queries were sent. Although an efficient method for dummy generation remains as a future task, the results also show that hiding the distribution of dummy is significantly effective for

protecting database privacy and the user has to know it with high accuracy in order to steal extra information from the server.

### Database privacy in malicious model

For our protocol, the difference between the malicious and semi-honest models is that in the malicious model the user may use an invalid input  $q$  whose components  $q_i$  are not necessarily in  $\{0, 1\}$ . If the user chooses  $q$  in such a way that some component  $q_i$  is extremely large and the remaining  $\ell - 1$  components are all zero, then  $\overline{\Pi}(p, q)$  will also be an extreme value (distinguishable from the dummy values) and depend dominantly on the bit  $p_i$ ; therefore, the user can almost surely guess the secret bit  $p_i$ . Since our protocol detects whether or not  $q_i$  is a bit value without invading user privacy, it can safely reject illegal queries and prevent any illegal query attacks, including above case.

### Performance evaluation

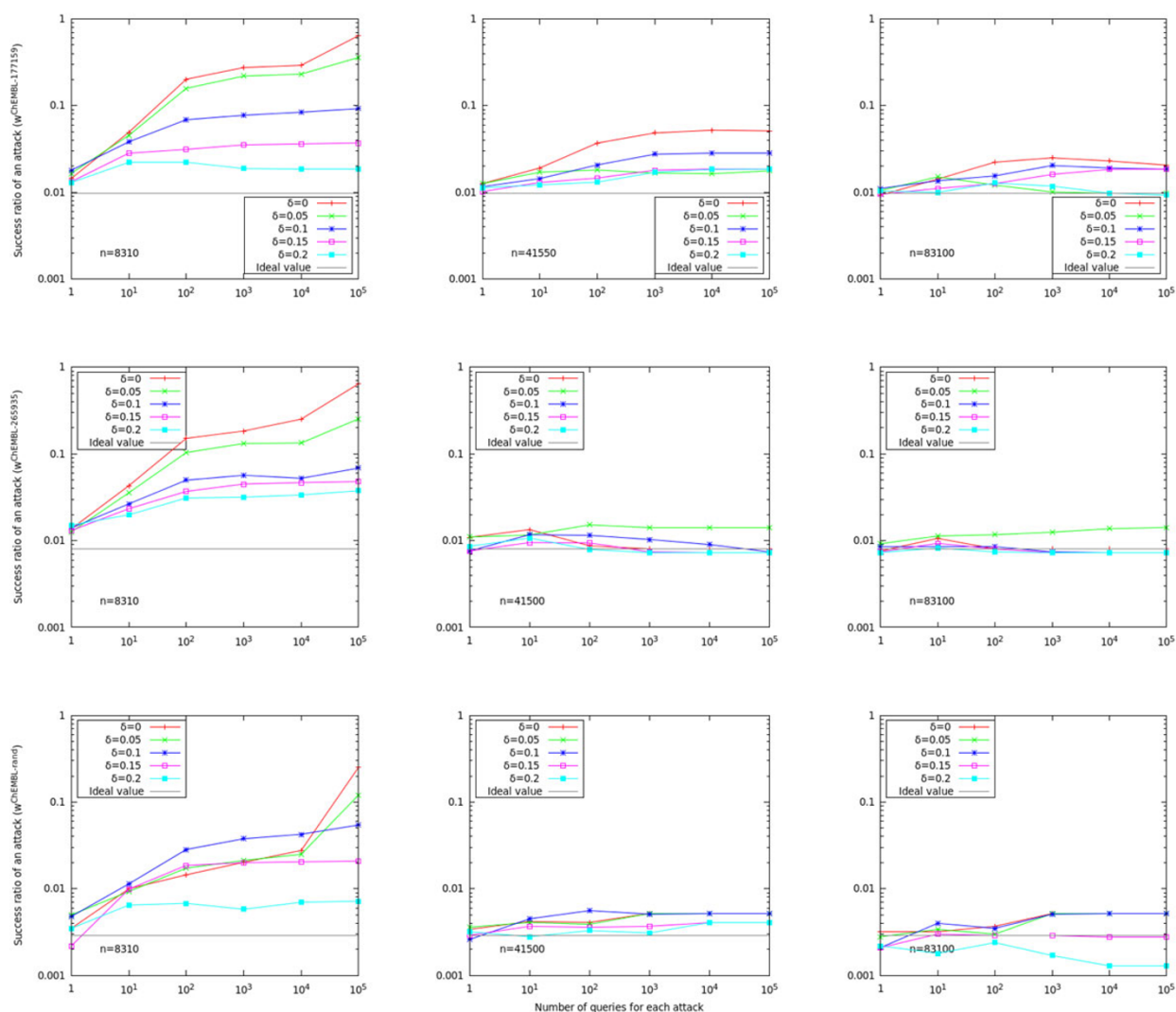
In this section, we evaluate the performance of the proposed method on two datasets created from ChEMBL.

### Implementation

We implemented the proposed protocol using the C++ library of elliptic curve ElGamal encryption [24], in which the NIZK proposed in the previous study [21] is also implemented.

For the implementation, we used parameters called secp192k1, as recommended by SECG (The Standards for Efficient Cryptography Group). These parameters are considered to be more secure than 1024-bit RSA encryption, which is the most commonly used public-key cryptosystem. The implementation of

Owing to the limitation of the range of plaintext, the implementation here does not include sign-preserving randomization. For the purpose of comparison, we also implemented a GP-MPC protocol by using Fairplay [25]. In order to reduce the circuit size of the GP-MPC, we implemented a simple task that computes the sign of Tversky index between a query and a fingerprint in the database, and repeated the task for all the fingerprints in



**Figure 4** The comparison of the experimental success ratios of the user's guess based on the server's return as well as the prior distribution of true value when the user sends many queries ( $\delta = 0, \dots, 0.2$ ), and success probability based only on a guess from the prior distribution (ideal value).  $\overline{\Pi}_{1, 1, 0.8}$  ( $k = 831$ ) is assumed and results are calculated for three different numbers of dummies ( $n = 831 \times 10, 831 \times 50, 831 \times 10^2$ ) when the user sends  $L = 1, 10, \dots, 10^5$  queries and three different distributions:  $w^{\text{CHEMBL-177159}}$  and  $w^{\text{CHEMBL-265935}}$  are actual distributions of  $\overline{\Pi}_{1, 1, 0.8}$  on ChEMBL obtained by querying two randomly selected fingerprints from ChEMBL,  $w^{\text{rand}}$  is obtained by randomly selecting a value from  $1, \dots, k$  for  $m = 5 \times 831$  times and dividing each observed frequency by  $m$ .

the database. Thus the CPU time and data transfer size of the implementation is linear to the size of database.

### Experimental setup

The Jaccard index along with the threshold  $\theta = 0.8$  were used for both protocols. For SSCC, we used 10,000 dummies. These two implementations were tested on two datasets: one, referred to as ChEMBL 1000, was the first 1000 fingerprints stored in ChEMBL, and the other, referred to as ChEMBL Full, was 1,292,344 fingerprints in the latest version of ChEMBL. All the programs were run on a single core of an Intel Xeon 2.9 GHz on the

same machine equipped with 64 GB memory. To avoid environmental effects, we repeated the same experiment five times and calculated average values.

### Results

The results are shown in Table 2. Despite the proposed method including elaborate calculation like the NIZK proof, we can see from the results that both the CPU time and communication size of the proposed method are significantly smaller than those of the GP-MPC protocol. Furthermore, it is clear that SSCC provides industrial-strength performance, considering that it works,

**Table 2 CPU time and communication size of secure similar compounds counter (SSCC) and those of general-purpose multi-party computation (GP-MPC).**

	ChEMBL_1000	ChEMBL_Full
CPU time (s)		
SSCC (server)	0.69	167.19
SSCC (client)	1.53	172.37
GP-MPC (server)	4, 075.15	–
GP-MPC (client)	4, 366.18	–
Communication size (MB)		
SSCC (server → client)	2.24	265.33
SSCC (client → server)	0.03	0.03
GP-MPC (server → client)	42.50	–
GP-MPC (client → server)	2, 128.00	–

The experiment on ChEMBL Full by GP-MPC did not finish within 24 hours.

even on a huge database like ChEMBL Full, taking no more than 167 s and 173 s for the server and client respectively.

The experiment on ChEMBL Full by using GP-MPC did not finish within 24 hours. Since both CPU time and communication size are exactly linear to the size of database for the GP-MPC protocol, the results of ChEMBL Full for GP-MPC are estimated to be more than 1600 hours for both sides and 3 Gbyte data transfer from client to server, considering the results of ChEMBL 1000.

By using simple data parallelization, the computational speed will be improved linearly with the number of CPUs. Since all the programs were run on the same machine there was almost no latency for the communication between the two parties in these experiments. Therefore, GP-MPC, whose communication size is huge, is expected to require far more time when it runs on an actual network that is not always in a good condition. The other important point is that SSCC requires only two data transfers, which enables data transfer after off-line calculation. On the other hand, GP-MPC must keep online during the search because of the high communication frequency. We also note that it took less than 100 MB to compile SSCC, while GP-MPC required more than 16 GB. Considering these observations, SSCC is efficient for practical use. It is known that several techniques improve the performance of GP-MPC and the previous work by Pinkas et al. [26] reported that Free XOR [27] and Garbled Row Reduction [26], which are commonly used in state-of-the-art GP-MPC methods [28-31], reduced running time and communication size by factors of 1.8 and 6.3 respectively when a circuit computing an encryption of AES was evaluated. Though these techniques are not implemented in Fairplay, we consider that GP-MPC is yet far less practical for the large-scale chemical compound search problem compared to our method which improved running time and communication size by factors of 36, 900 and 12, 000.

## Conclusion

In this study, we proposed a novel privacy-preserving protocol for searching chemical compound databases. To our knowledge, this is the first practical study for privacy-preserving (for both user and database sides) similarity searching in the fields of bioinformatics and cheminformatics. Moreover, the proposed method could be applied to a wide range of life science problems such as searching for similar single-nucleotide polymorphism (SNP) patterns in a personal genome database. While the protocol proposed here focuses on searching for a number of similar compounds, we are examining further improvements of the protocol such as the client being able to download similar compounds; we expect this on-going study to further contribute to the drug screening process. In recent years, open innovation has been attracting attention as a promising approach for speeding up the process of new drug discovery [32]. For example, research on neglected tropical diseases including malaria has been promoted by the recent attempt to share chemical compound libraries in the research community. In spite of high expectations, such an approach is still limited to economically less important problems on account of privacy problems [33]. Therefore, privacy-preserving data mining technology is expected to be the breakthrough promoting open innovation and we believe that our study will play an important role.

## Additional material

Additional File 1: Supplementary text (pdf).

### Competing interests

The authors declare that they have no competing interests.

### Authors' contributions

KS, HA, KN and KA designed the protocol inspired by the discussion with TH. KS, KN, NA and GH conducted security analyses. KS and SM implemented the protocol. KS evaluated performance of the protocol. All authors wrote the manuscript. All authors read, commented and approved the final manuscript.

### Acknowledgements

KS thanks Yusuke Sakai and Takahiro Matsuda for fruitful discussions.

### Declarations

This work was supported by the Japan-Finland Cooperative Scientific Research Program of JST/AMED (to KS) and JSPS KAKENHI Grant Number 25540131 (to KS and MH). JS and KT are supported by JST CREST. KT is supported by JST ERATO, RIKEN PostK, NIMS MI2I, JSPS KAKENHI Nanostructure and JSPS KAKENHI Grant Number 15H05711. JS is supported by JSPS KAKENHI Grant Number 24680015. This article has been published as part of *BMC Bioinformatics* Volume 16 Supplement 18, 2015: Joint 26th Genome Informatics Workshop and 14th International Conference on Bioinformatics: Bioinformatics. The full contents of the supplement are available online at <http://www.biomedcentral.com/bmcbioinformatics/supplements/16/S18>.

#### Authors' details

<sup>1</sup>Biotechnology Research Institute for Drug Discovery, National Institute of Advanced Industrial Science and Technology (AIST), 2-4-7 Aomi Koto-ku, Tokyo, Japan. <sup>2</sup>Information Technology Research Institute, National Institute of Advanced Industrial Science and Technology (AIST), 2-4-7 Aomi Koto-ku, Tokyo, Japan. <sup>3</sup>Japan Science and Technology Agency (JST) PRESTO Researcher, Tokyo, Japan. <sup>4</sup>Information Technology Center, The University of Tokyo 7-3-1 Hongo Bunkyo-ku, Tokyo, Japan. <sup>5</sup>Cybozu Labs 12F, Koraku Mori Bldg. 1-4-14, Tokyo, Japan. <sup>6</sup>Faculty of Science and Engineering, Waseda University 3-4-1 Okubo Shinjuku-ku, Tokyo, Japan. <sup>7</sup>Graduate School of Frontier Sciences, The University of Tokyo 5-1-5, Chiba, Japan. <sup>8</sup>Molecular Profiling Research Center for Drug Discovery, National Institute of Advanced Industrial Science and Technology (AIST), 2-4-7 Aomi Koto-ku, Tokyo, Japan. <sup>9</sup>Graduate School of SIE, University of Tsukuba 1-1-1 Tennodai Tsukuba, Ibaraki, Japan.

Published: 9 December 2015

#### References

- Subbaraman N: **Flawed arithmetic on drug development costs.** *Nature Biotechnology* 2011, **29**(5):381-381.
- Miller Ma: **Chemical database techniques in drug discovery.** *Nature Reviews Drug Discovery* 2002, **1**(3):220-7.
- Schooler J: **Unpublished results hide the decline effect.** *Nature* 2011, **470**:437.
- Ostrovsky R, Skeith WE III: **A survey of single-database private information retrieval: techniques and applications.** *Proceedings of the 10th International Conference on Practice and Theory in Public-key Cryptography PKC'07* 2007, 393-411.
- Goethals B, Laur S, Lipmaa H, Mielik'ainen T: **On private scalar product computation for privacy-preserving data mining.** *Proceedings of the 7th Annual International Conference on Information Security and Cryptology ICISC 2004* 2004, 104-120.
- Blundo C, Cristofaro ED, Gasti P: **ESPRESSo : Efficient Privacy-Preserving Evaluation of Sample Set Similarity.** *Proceedings of Data Privacy Management and Autonomous Spontaneous Security: 7th International Workshop, DPM 2009 and 5th International Workshop, SETOP 2012 DMP/SETOP 2012* 2012, 89-103.
- Murugesan M, Jiang W, Clifton C, Si L, Vaidya J: **Efficient privacy-preserving similar document detection.** *The VLDB Journal* 2010, **19**(4):457-475.
- Yao ACC: **How to generate and exchange secrets.** *Proceedings of the 27th Annual Symposium on Foundations of Computer Science SFCS '86* 1986, 162-167.
- Gentry C: **Fully homomorphic encryption using ideal lattices.** *Proceedings of the 41st Annual ACM Symposium on Theory of Computing STOC '09* 2009, 169-178.
- Togan M, Plesca C: **Comparison-based computations over fully homomorphic encrypted data.** *Communications (COMM), 2014 10th International Conference* 2014, 1-6, doi:10.1109/ICComm.2014.6866760.
- Laur S, Lipmaa H: **On private similarity search protocols.** *Proceedings of the 9th Nordic Workshop on Secure IT Systems NordSec 2004*, 73-77.
- Popa RA, Redfield CMS, Zeldovich N, Balakrishnan H: *Proceedings of the 23rd ACM Symposium on Operating Systems Principles SOSP 11* 85-100.
- Martin YC, Kofron JL, Traphagen LM: **Do structurally similar molecules have similar biological activity?** *Journal of Medicinal Chemistry* 2002, **45**(19):4350-4358.
- Miller JL: **Recent developments in focused library design: targeting gene-families.** *Current Topics in Medicinal Chemistry* 2006, **6**(1):19-29.
- Curtis R, Tang J: **Someone's loss might be your gain: A case of negative results publications in science.** *Proceedings of the American Society for Information Science and Technology ASIST 2012*, 49.
- Gaulton A, Bellis LJ, Bento AP, Chambers J, Davies M, Hersey A, Light Y, McGlinchey S, Michalovich D, Al-Lazikani B, Overington JP: **ChEMBL: a large-scale bioactivity database for drug discovery.** *Nucleic Acids Research* 2012, **40**(Database):1100-1107.
- Goldwasser S, Micali S: **Probabilistic encryption.** *J Comput Syst Sci* 1984, **28**(2):270-299.
- Paillier P: **Public-key cryptosystems based on composite degree residuosity classes.** *Proceedings of the 17th International Conference on Theory and Application of Cryptographic Techniques EUROCRYPT'99* 1999, 223-238.
- ElGamal T: **A public key cryptosystem and a signature scheme based on discrete logarithms.** *IEEE Transactions on Information Theory* 1985, **31**(4):469-472.
- Goldreich O: *Foundations of Cryptography: Volume 1.* Cambridge University Press; 2001.
- Sakai Y, Emura K, Hanaoka G, Kawai Y, Omote K: **Methods for restricting message space in public-key encryption.** *IEICE Transactions* 2013, **96-A**(6):1156-1168.
- Tversky A: **Features of similarity.** *Psychological Review* 1977, **84**(4):327-352.
- Damgård I, Fitz M, Kiltz E, Nielsen JB, Toft T: **Unconditionally secure constant-rounds multi-party computation for equality, comparison, bits and exponentiation.** *Proceedings of the 3rd Theory of Cryptography Conference TCC 2006* 2006, 285-304.
- C++ **Library implementing elliptic curve ElGamal crypto system [19].** 2015 [https://github.com/aistcrypt/Lifted-ElGamal], URL accessed April 13, 2015.
- Ben-david A, Nisan N, Pinkas B: **Fairplaymp: A system for secure multi-party computation.** *Proceedings of ACM Conference on Computer and Communications Security CCS 2008* 2008, 17-21.
- Pinkas B, Schneider T, Smart NP, Williams SC: **Secure two-party computation is practical.** *Proceedings of the 15th International Conference on the Theory and Application of Cryptology and Information Security ASIACRYPT 2009* 2009, 250-267.
- Kolesnikov V, Schneider T: **Improved garbled circuit: Free XOR gates and applications.** *Proceedings of the 35th International Colloquium on Automata, Languages and Programming ICALP 2008* 2008, 486-498.
- Henecka W, Kögl S, Sadeghi A, Schneider T, Wehrensberg I: **TASTY: tool for automating secure two-party computations.** *Proceedings of the 17th ACM Conference on Computer and Communications Security CCS 2010* 2010, 451-462.
- Huang Y, Evans D, Katz J, Malka L: **Faster secure two-party computation using garbled circuits.** *Proceedings of the 20th USENIX Security Symposium USENIX 2011* 2011.
- Huang Y, Shen CH, Evans D, Katz J, Shelat A: **Efficient secure computation with garbled circuits.** *Proceedings of the 7th International Conference on Information Systems Security ICIS 2011*, 28-48.
- Kreuter B, Shelat A, Shen C: **Billion-gate secure computation with malicious adversaries.** *Proceedings of the 21th USENIX Security Symposium USENIX Security 2012* 2012, 285-300.
- Williams AJ, Harland L, Groth P, Pettifer S, Chichester C, Willighagen EL, Evelo CT, Blomberg N, Ecker G, Goble C, Mons B: **Open PHACTS: semantic interoperability for drug discovery.** *Drug Discovery Today* 2012, **17**(21-22):1188-1198.
- Hunter J, Stephens S: **Is open innovation the way forward for big pharma?** *Nature Reviews Drug Discovery* 2010, **9**(2):87-88.

doi:10.1186/1471-2105-16-S18-S6

Cite this article as: Shimizu et al.: Privacy-preserving search for chemical compound databases. *BMC Bioinformatics* 2015 **16**(Suppl 18):S6.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit

