

Efficient Privacy-Preserving Variable-Length Substring Match for Genome Sequence

Yoshiki Nakagawa ✉

Department of Computer Science and Engineering, Waseda University, Tokyo, Japan

Satsuya Ohata ✉

Digital Garage, Inc., Tokyo, Japan

Kana Shimizu¹ ✉

Department of Computer Science and Engineering, Waseda University, Tokyo, Japan

National Institute of Advanced Industrial Science and Technology, Tokyo, Japan

Abstract

Finding a similar substring that commonly appears in query and database sequences is an essential task for genome data analysis. This study proposes a secure two-party variable-length string search protocol based on secret sharing. The unique feature of our protocol is that time, communication, and round complexities are not dependent on the database length N , after the query input. This property brings dramatic performance improvements in search time, since N is usually quite large in an actual genome database, and the same database is repeatedly used for many queries. Our concept hinges on a technique that efficiently applies the compressed full-text index (FOCS 2000) for a secret-sharing scheme. We conducted an experiment using a human genomic sequence with the length of 10 million as the database and a query with the length of 100 and found that the query response time of our protocol was at least three orders of magnitude faster than a well-designed baseline protocol under the realistic computation/network environment.

2012 ACM Subject Classification Theory of computation → Pattern matching; Security and privacy → Privacy-preserving protocols; Theory of computation → Cryptographic protocols; Applied computing → Genomics

Keywords and phrases Private Genome Sequence Search, Secure Multiparty Computation, Secret Sharing, FM-index, Suffix Tree, Maximal Exact Match

Digital Object Identifier 10.4230/LIPIcs.WABI.2021.2

Supplementary Material *Software (Source Code)*: <https://waseda.box.com/v/wabi2021-suppl-pggs-src>

Funding This work is partially supported by JST CREST grant number JPMJCR19F6.

Kana Shimizu: Supported part in MEXT/JSPS KAKENHI Grant Number 19K12209 and 21H04871.

1 Introduction

The dramatic reduction in the cost of genome sequencing has prompted increased interest in personal genome sequencing over the last 15 years. Extensive collections of personal genome sequences have been accumulated both in academic and industrial organizations, and there is now a global demand for sharing the data to accelerate scientific research [13, 24]. As discussed in previous studies, disclosing personal genome information has a high privacy risk [10], so it is crucial to ensure that individuals' privacy is protected upon data sharing. At present, the most popular approach for this is to formulate and enforce a privacy policy, but it is a time-consuming process to reach an agreement, especially among stakeholders with different legal backgrounds, which slows down the pace of research. Therefore, there is a

¹ Corresponding author



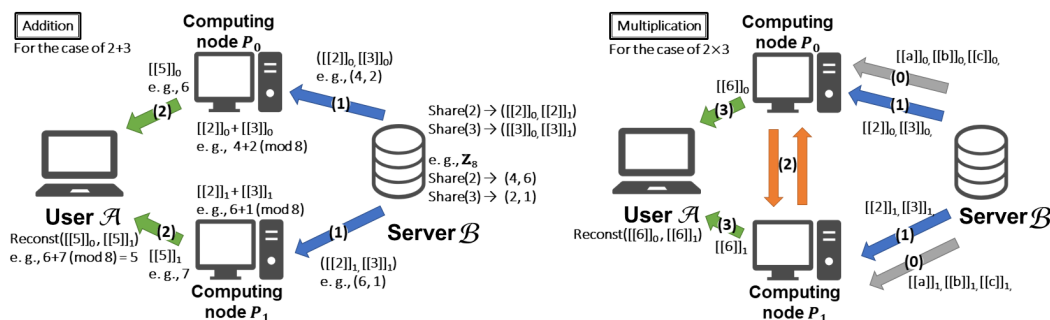
strong demand for privacy-preserving technologies that can potentially compensate for or even replace the traditional policy-based approach [3, 21]. One important application that needs a privacy-preserving technology is private genome sequence search, where different stakeholders respectively hold a query sequence and a database sequence and the goal is to let the query holder know the result while simultaneously keeping the query and the database private. Many studies have addressed the problem of how to compute exact or approximate edit distance or the longest common substring (LCS) through techniques based on homomorphic encryption [17, 8, 22] and secure multi-party computation (MPC) [15, 31, 33, 7, 2, 26, 23], or how to compute sequence similarity based on private set intersection [4]. While these studies can evaluate global sequence similarity for two sequences of similar length, other studies address the problem of finding a substring between a query and a long genome sequence or a set of long genome sequences, with the aim of evaluating local sequence similarity [28, 16, 30, 29, 18, 6, 25]. [28] proposed an approach to combine an additive homomorphic encryption and index structures such as FM-index [11] and the positional Burrows-Wheeler transform [9] to find the longest prefix of a query that matches a database (LPM) and a set-maximal match for a collection of haplotypes. [30] used a similar approach and improved the time and communication complexities for LPM on a protein sequence by using a wavelet matrix. [16] improved the round complexity of a set-maximal match, though the search time was more than one order of magnitude slower than [28] due to the heavy computational cost caused by the fully homomorphic encryption. [29] used the Goldreich-Micali-Wigderson protocol to build a suffix tree for a set-maximal match. According to experiments by [18], the search time of [29] is one order of magnitude slower than [28] and [18]. [18] used a garbled circuit to build a suffix tree for substring match and a set-maximal match under a different security assumption such that the tree-traversal pattern is leaked to the cloud server. [6] and [25] found fixed-length substring matches using a one-way hash function or homomorphic encryption on a public cloud under a security assumption such that the database is a public sequence and a query is leaked to a private cloud server.

In this study, we aim to improve privacy-preserving substring match under the security assumption such that both the query and the database sequence are strictly protected. We first propose a more efficient method for finding LPM, and then extend it to find the longest maximal exact match (LMEM), which is more practically important in bioinformatics. We designed the protocol for LMEM for ease of explanation, and the protocol can be applied to similar problems such as finding all maximal exact matches (MEMs) with a small modification. To our knowledge, this is the first study to address the problem of securely finding MEMs.

Our Contribution

The time complexity of the previous studies [28, 30] include the factor of N , and thus they do not scale well to a large database. For a similar reason, using secure matching protocols (e.g., [32]) for the shares (or tags in searchable encryption) of all substrings in a query and database is even worse in terms of time complexity. To achieve a real-time search on an actual genome database, we propose novel secret-sharing-based protocols that do not include the factor of N in the time, communication, and round complexities for the search time (i.e., the time after the input of a query until the end of the search).

The basic idea of the protocols is to represent the database string by a compressed index [11, 12] and store the index as a lookup table. LPM and MEMs are found by at most ℓ and 2ℓ table lookups respectively, where ℓ is the length of the query. More specifically, the table V is referenced in a recursive manner; i.e., one needs to obtain $V[j]$, where $j = V[i]$, given i . To ensure security, we need to compute $V[j]$ without seeing any element of V .



■ **Figure 1** Arithmetic addition and multiplication over secret sharing.

The key technical contribution of this study is an efficient protocol that achieves this type of recursive reference. We named the protocol secret-shared recursive oblivious transfer (ss-ROT). While the previous studies requires $O(N)$ time complexity to ensure security, the time, communication, and round complexities of ss-ROT are all $O(\ell)$ for ℓ recursive table lookups, except for the preparation of the table and generation of shares before the query input. Since the entire protocols mainly consist of ℓ table lookups for LPM, and 2ℓ table lookups and 2ℓ inner product computations for LMEM, the search times for LPM and LMEM do not depend on the database size.

We implemented the proposed protocol and tested it on substrings of a human genome sequence 10^3 to 10^7 in length and confirmed that the actual CPU time and data transfer overhead were in good agreement with the theoretical complexities. We also found that the search time of our protocol was three orders of magnitude faster than that of the previous method [28, 30]. For conducting further performance analysis, we designed and implemented baseline protocols using major techniques of secret-sharing-based protocols. The results showed that the search times of our protocols were at least two orders of magnitude faster than those of the baseline protocols.

2 Preliminaries

2.1 Secure Computation based on Secret Sharing

Here, we explain the 2-out-of-2 additive secret sharing $((2, 2)$ -SS) scheme and how to securely compute arithmetic/Boolean gates (Figure 1).

Secret Sharing and Secure Computation

In t -out-of- n secret sharing (e.g., [27]), we split the secret value x into n pieces, and can reconstruct x by combining more or an equal number of t pieces. We call the split pieces “share”. The basic security notion for secret sharing is that we cannot obtain any information about x even if we gather less than or equal to $(t-1)$ shares. In this paper, we consider a case with $(t, n) = (2, 2)$. A 2-out-of-2 secret sharing $((2, 2)$ -SS) scheme over \mathbb{Z}_{2^n} consists of two algorithms: **Share** and **Reconst**. **Share** takes as input $x \in \mathbb{Z}_{2^n}$ and outputs $([x]_0, [x]_1) \in \mathbb{Z}_{2^n}^2$, where the bracket notation $[x]_i$ denotes the arithmetic share of the i -th party (for $i \in \{0, 1\}$). We denote $[x] = ([x]_0, [x]_1)$ as their shorthand. **Reconst** takes as inputs $[x]_0$ and $[x]_1$ and outputs x . For arithmetic sharing $[x]_i$ and Boolean sharing $[x]_i^B$, we consider power-of-two integers n (e.g., $n = 16$) and $n = 1$, respectively.

Depending on the secret sharing scheme, we can compute arithmetic/Boolean gates over shares; that is, we can execute some kind of processing related to x without x . This means it is possible to perform some computation without violating the privacy of the secret data,

■ **Table 1** Secure subprotocols used in this paper.

	Input	Output
Equality	$\llbracket x \rrbracket, \llbracket y \rrbracket$	$\llbracket z \rrbracket^B$ s.t. $z = 1$ if $x = y$ otherwise $z = 0$
Comp	$\llbracket x \rrbracket, \llbracket y \rrbracket$	$\llbracket z \rrbracket^B$ s.t. $z = 1$ if $x < y$ otherwise $z = 0$
CastUp	$\llbracket x \rrbracket \in \mathbb{Z}_{2^n}, n'$	$\llbracket x \rrbracket \in \mathbb{Z}_{2^{n'}} (n < n')$
B2A	$\llbracket x \rrbracket^B$	$\llbracket x \rrbracket$
Choose	$\llbracket x \rrbracket, \llbracket y \rrbracket, \llbracket e \rrbracket \in \{0, 1\}$	$\llbracket z \rrbracket$ s.t. $z = x$ if $e = 1$, otherwise ($e = 0$) $z = y$

and is called secure (multi-party) computation. It is known that we can execute arbitrary computation by combining basic arithmetic/Boolean gates. In the following paragraphs, we show how to concretely compute these gates over shares.

Semi-Honest Secure Two-Party Computation Based on (2, 2)-Additive SS

We use a standard (2, 2)-additive SS scheme, defined by

- $\text{Share}(x)$: randomly choose $r \in \mathbb{Z}_{2^n}$ and let $\llbracket x \rrbracket_0 = r$ and $\llbracket x \rrbracket_1 = x - r$.
- $\text{Reconst}(\llbracket x \rrbracket_0, \llbracket x \rrbracket_1)$: output $\llbracket x \rrbracket_0 + \llbracket x \rrbracket_1$.

Note that one of the shares of x ($\llbracket x \rrbracket_0$ or $\llbracket x \rrbracket_1$) does not reveal any information about x . In Figure 1, the secret value $x = 2$ is split into $\llbracket x \rrbracket_0 = 4$ and $\llbracket x \rrbracket_1 = 6$. These are valid (2, 2)-additive shares because $4 + 6 \equiv 2 \pmod{8}$ holds. Even if we can see $\llbracket x \rrbracket_0 = 4$, we cannot decide the value of x since we execute a split of x uniformly at random. This means, in Figure 1, computing nodes P_0 and P_1 cannot obtain any information about x as long as these two nodes do not collude. On the other hand, we can compute arithmetic ADD/MULT gates over shares as follows:

- $\llbracket z \rrbracket \leftarrow \text{ADD}(\llbracket x \rrbracket, \llbracket y \rrbracket)$ can be done locally by just adding each party's share on x and on y . In Figure 1 (left), we show an example of secure addition. P_0/P_1 obtain shares 6/7 by adding their two shares. In this process, P_0/P_1 cannot find they are computing $2 + 3$.
- Multiplication is more complex than addition. There are various methods for multiplication over shares, most of which require communication between computing nodes. In this paper, we use the standard method for $\llbracket w \rrbracket \leftarrow \text{MULT}(\llbracket x \rrbracket, \llbracket y \rrbracket)$ based on Beaver triples (BT) [5]. Such a triple consists of $\text{bt}_0 = (a_0, b_0, c_0)$ and $\text{bt}_1 = (a_1, b_1, c_1)$ such that $(a_0 + a_1)(b_0 + b_1) = (c_0 + c_1)$. Hereafter, a , b , and c denote $a_0 + a_1$, $b_0 + b_1$, and $c_0 + c_1$, respectively. We use these BTs as auxiliary inputs for computing MULT. Note that we can compute them in advance (or in offline phase) since they are independent of inputs $\llbracket x \rrbracket$ and $\llbracket y \rrbracket$. We adopt a trusted initializer setting (e.g., [19, 20]); that is, BTs are generated by the party other than two computing nodes and then distributed. In the online phase of MULT, each i -th party P_i ($i \in \{0, 1\}$) can compute the multiplication share $\llbracket z \rrbracket = \llbracket xy \rrbracket$ as follows:
 1. P_i first computes $(\llbracket x \rrbracket_i - a_i)$ and $(\llbracket y \rrbracket_i - b_i)$, and sends them to P_{1-i} .
 2. P_i reconstructs $x' = x - a$ and $y' = y - b$.
 3. P_0 computes $\llbracket z \rrbracket_0 = x'y' + x'b_0 + y'a_0 + c_0$, and P_1 computes $\llbracket z \rrbracket_1 = x'b_1 + y'a_1 + c_1$. Here, $\llbracket z \rrbracket_0$ and $\llbracket z \rrbracket_1$ calculated with the above procedures are valid shares of xy ; that is, $\text{Reconst}(\llbracket z \rrbracket_0, \llbracket z \rrbracket_1) = xy$. We shorten the notations and write the ADD and MULT protocols simply as $\llbracket x \rrbracket + \llbracket y \rrbracket$ and $\llbracket x \rrbracket \cdot \llbracket y \rrbracket$, respectively.

We also write $\text{ADD}(\text{ADD}(\llbracket x_A \rrbracket, \llbracket x_B \rrbracket), \llbracket x_C \rrbracket)$ as $\Sigma_{c=\{A,B,C\}} \llbracket x_c \rrbracket$. Note that, similarly to the ADD protocol, we can also locally compute multiplication by constant c , denoted by $c \cdot \llbracket x \rrbracket$. We can easily extend the above protocols to Boolean gates. By converting $+$ and $-$ into

\oplus in the arithmetic ADD and MULT protocols, we can obtain the XOR and AND protocols, respectively. We can construct NOT and OR protocols from the properties of these gates. When we compute NOT($\llbracket x \rrbracket_0^B, \llbracket x \rrbracket_1^B$), P_0 and P_1 output $\neg \llbracket x \rrbracket_0^B$ and $\llbracket x \rrbracket_1^B$, respectively. When we compute OR($\llbracket x \rrbracket^B, \llbracket y \rrbracket^B$), we compute $\neg \text{AND}(\neg \llbracket x \rrbracket^B, \neg \llbracket y \rrbracket^B)$. We shorten the notations and write XOR, AND, NOT, and OR simply as $\llbracket x \rrbracket \oplus \llbracket y \rrbracket$, $\llbracket x \rrbracket \wedge \llbracket y \rrbracket$, $\neg \llbracket x \rrbracket$, and $\llbracket x \rrbracket \vee \llbracket y \rrbracket$, respectively. By combining the above gates, we can securely compute higher-level protocols. The functionality of the secure subprotocols [23] used in this paper are shown in Table 1. Due to space limits, we omit the details of their construction. Note that we can compute Choose by $\llbracket z \rrbracket = \llbracket y \rrbracket + \llbracket e \rrbracket \cdot (\llbracket x \rrbracket - \llbracket y \rrbracket)$. In this paper, we consider the standard simulation-based security notion in the presence of semi-honest adversaries (for 2PC), as in [14]. We show the definition in Appendix B. Roughly speaking, this security notion guarantees the privacy of the secret under the condition that computing nodes do not deviate from the protocol; that is, although computing nodes are allowed to execute arbitrary attacks in their local, they do not (maliciously) manipulate transmission data to other parties. The building blocks we adopt in this paper satisfy this security notion. Moreover, as described in [14], the composition theorem for the semi-honest model holds; that is, any protocol is privately computed as long as its subroutines are privately computed.

2.2 Index Structure for String Search

Notation and Definition

Σ denotes a set of ordered symbols. A string consists of symbols in Σ . We denote a lexicographical order of two strings S and S' by $S \leq S'$ (i.e., A < C < G < T and AAA < AAC). We denote the i -th letter of a string S by $S[i]$ and a substring starting from the i -th letter to the j -th letter by $S[i, j]$. The index starts with 0. The length of S is denoted by $|S|$. A reverse string of S (i.e., $S[|S| - 1], \dots, S[0]$) is denoted by \hat{S} . We consider a direction from the i -th position to the j -th position as rightward if $i < j$ and leftward otherwise.

Given a query w and a database S , we define the longest prefix that matches a database string (LPM) by $\max_{(0,j)} \{j | w[0, \dots, j] = S[k, \dots, l]\}$, where $0 \leq j < \ell$ and $0 \leq k \leq l < N$, and the longest maximal exact match (LMEM) by $\max_{(i,j)} \{j - i | w[i, \dots, j] = S[k, \dots, l]\}$, where $0 \leq i \leq j < \ell$ and $0 \leq k \leq l < N$.

FM-Index and related data structures

FM-Index [11] and related data structures [12] are widely used for genome sequence search. Given a query string w of length ℓ and a database string S of length N , [11] enables LPM to be found in $O(\ell)$ time regardless of N , and it also enables LMEM to be found in $O(\ell)$ if auxiliary data structures are used [12]. Given all the suffixes of a string S : $S[0, \dots, |S| - 1]$, $S[1, \dots, |S| - 1], \dots, S[|S| - 1]$, a suffix array is an array of positions $(p_0, \dots, p_{|S|-1})$ such that $S[p_0, \dots, |S| - 1] \leq S[p_1, \dots, |S| - 1] \leq S[p_2, \dots, |S| - 1], \dots, \leq S[p_{|S|-1}, \dots, |S| - 1]$. We denote the suffix array of S by SA and denote its i -th element by $SA[i]$. A Burrows-Wheeler transform (BWT) is a permutation of the sequence S such that its i -th letter becomes $S[SA[i] - 1]$. We denote a BWT of S by L and denote its i -th letter by $L[i]$. Let us define a rank of S for a letter $c \in \Sigma$ at position t by $\text{Rank}_c(t, S) = |\{j | S[j] = c, 0 \leq j < t\}|$ and a count of occurrences of letters that are lexicographically smaller than c in S by $\text{CF}_c(S) = \sum_{r < c} \text{Rank}_r(|S|, S)$, and the operation $\text{LF}_c(i, S) = \text{CF}_c(L) + \text{Rank}_c(i, L)$. The match between w and S is reported as a form of left-closed and right-open interval on SA , and the lower and upper bounds of the interval are respectively computed by LF. Given a

letter c and an interval $[f, g]$ that corresponds to suffixes that share the prefix x (i.e., $[f, g]$ reports the locations of the substring x in S), we can find a new interval that corresponds to all suffixes that share the prefix cx (i.e., locations of the substring cx) by

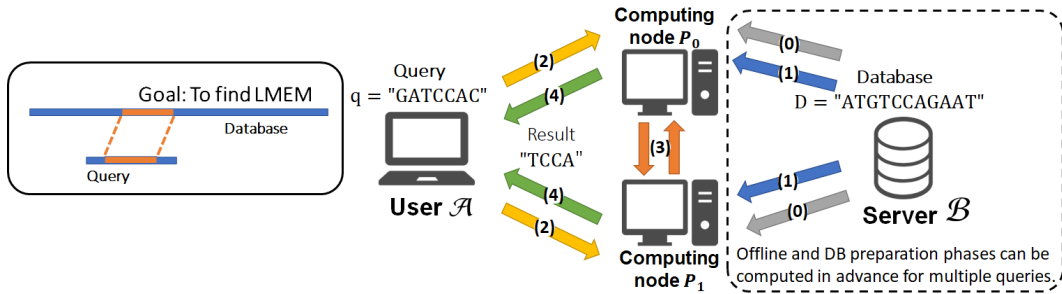
$$[f', g'] = [\text{LF}_c(f, S), \text{LF}_c(g, S)]. \quad (1)$$

The leftward extension of the match is called a backward search, which is the main functionality of FM-Index. By starting the search with the initial interval $[0, N)$ and conducting the backward searches for $w[\ell - 1], w[\ell - 2], \dots$, the longest suffix match is detected when $f = g$. Rank and CF are precomputed and stored in an efficient form that can be searched in constant time. Therefore, the longest suffix match can be computed in $O(\ell)$ time. LPM is found if the search is conducted on \hat{S} and match is extended by $w[0], w[1], \dots, w[\ell - 1]$.

Searching LMEM by repeating LPM for $w[0, \dots, \ell - 1], w[1, \dots, \ell - 1], w[2, \dots, \ell - 1], \dots, w[\ell - 1]$ takes $O(\ell^2)$ time. We can improve it to $O(\ell)$ time by using the longest common prefix (LCP) array and related data structures [12]. The LCP array, denoted by LCP, is an array that stores the length of the longest prefix of $S[\text{SA}[i - 1], |S| - 1]$ and $S[\text{SA}[i], |S| - 1]$ in $\text{LCP}[i]$ for $0 < i \leq N$. The lcp-interval $[i, j)$ of lcp-value d is an interval such that it satisfies $\text{LCP}[i] < d, \text{LCP}[j] < d, \text{LCP}[k] > d$ for all $k \in \{i + 1, \dots, j - 1\}$, and $\text{LCP}[k] = d$ for at least one $k \in \{i + 1, \dots, j - 1\}$, and is denoted by $d - [i, j)$. $d - [i, j)$ corresponds to all the suffixes that share the prefix $S[\text{SA}[i], \dots, \text{SA}[i] + d - 1]$. The parent interval of $d - [i, j)$ is the lcp-interval $h - [m, n)$ such that $h < d$ and $0 \leq m \leq i < j \leq n < N$, and there is no other lcp-interval $t - [r, s)$ such that $h < t < d$ and $0 \leq m \leq r \leq i < j \leq s \leq n < N$. The parent of the lcp-interval $[f, g)$ can be found by

$$[f', g'] = \begin{cases} [\text{PSV}[f_i], \text{NSV}[f_i]) & \text{LCP}[g_i] \leq \text{LCP}[f_i] \\ [\text{PSV}[g_i], \text{NSV}[g_i]) & (\text{otherwise}), \end{cases} \quad (2)$$

where $\text{PSV}[i] = \max\{j | 0 \leq j < i \wedge \text{LCP}[j] < \text{LCP}[i]\}$ and $\text{NSV}[i] = \min\{j | i \leq j < N \wedge \text{LCP}[j] < \text{LCP}[i]\}$. By finding a parent interval using PSV and NSV whenever it fails to extend the match, we can avoid useless backward searches, and thus LMEM is found at most 2ℓ backward searches. LCP, PSV and NSV are precomputed and stored in an efficient form that can be searched in constant time, so we can find LMEM in $O(\ell)$ time. See section 5.2 of [12] for more details of the data structures. Examples of the search by FM-Index, LCP, PSV, and NSV are provided in Appendix A.



■ **Figure 2** Schematic view of our goal and model. (0) Server (DB holder) distributes Beaver triples. (A reliable third party can serve as the trusted initializer instead.) (1) Server distributes shares of the database. (2) User (query holder) distributes shares of the query. (3) The computing nodes jointly calculate shares of the result. (4) The results are sent to User. The offline phase is (0), DB preparation phase is (1), and Search phase consists of (2)–(4).

■ **Table 2** Summary of complexities for our protocols and related protocols. BTime and Bsize are generation time and size of BTs. Dtime and Dsize are generation time for the shares of the database and size of the shares. Stime is the time for Search phase. Comm. is the size of data exchanged between computing nodes. Round is the number of data exchanges.

	Btime	Bsize	Dtime	Dsize	Stime	Comm.	Round
ss-ROT (ours)	0	0	ℓN	ℓN	ℓ	ℓ	ℓ
Secure LPM (ours)	ℓ	ℓ	ℓN	ℓN	ℓ	ℓ	ℓ
[30, 28] (LPM by AHE)	–	–	–	–	ℓN	$\ell\sqrt{N}$	ℓ
Baseline (LPM)	$\ell^2 N$	$\ell^2 N$	N	N	$\ell^2 N$	$\ell^2 N$	$\log \ell + \log N$
Secure LMEM (ours)	ℓ^2	ℓ^2	ℓN	ℓN	ℓ^2	ℓ^2	ℓ
Baseline (LMEM)	$\ell^3 N$	$\ell^3 N$	N	N	$\ell^3 N$	$\ell^3 N$	$\log \ell + \log N$

3 Proposed protocols

Problem Setting and Outline of Our Protocols

We assume that a query holder \mathcal{A} , a database holder \mathcal{B} , and two computing nodes P_0 and P_1 participate the protocol. \mathcal{A} holds a query string w of length ℓ and \mathcal{B} holds a database string T of length N . After the protocol is run, only \mathcal{A} knows LPM or LMEM between w and T . P_0 and P_1 do not obtain any information of w and T , except for ℓ and N .

Our protocol consists of offline, DB preparation, and Search phases. In the offline phase, \mathcal{B} generates BTs (correlated randomness used for multiplication) and sends them to P_0 and P_1 . In the DB preparation phase, \mathcal{B} creates a lookup table and distributes its shares to P_0 and P_1 . In the Search phase, \mathcal{A} generates shares of the query and sends them to P_0 and P_1 , and P_0 and P_1 jointly compute the result without obtaining any information of the lookup table. Finally, \mathcal{A} obtains the results. Figure 2 shows the schematic view of our goal and model. Note that the Offline and DB preparation phases do not depend on a query string, so they can be computed in advance for multiple queries.

In Section 3.1, we propose the important building block ss-ROT that enables recursive reference to a lookup table. In Section 3.2, we describe how to design the lookup table based on FM-Index, and propose an efficient protocol for LPM by using the lookup table and ss-ROT. In Section 3.3, we describe the additional table design for auxiliary data structures, and propose the complete protocol for LMEM. Table 2 summarizes the theoretical complexities of the three protocols. For comparison, the complexities of the baseline protocols and a previous method for LPM based on an additive homomorphic encryption [28, 30] are shown. As we mentioned in Section 1, the baseline protocols are designed using major techniques of secret-sharing-based protocols. The detailed algorithms are described in Appendix C.

3.1 Secret-shared Recursive Oblivious Transfer

We define a problem called a secret-shared recursive oblivious transfer (ss-ROT) as follows.

► **Definition 1.** *We assume a database holder \mathcal{B} and two computing nodes P_0 and P_1 participate the protocol. \mathcal{B} holds a vector V of length N and $0 \leq V[i] < N$. Given the initial position p_0 and the depth of recursion ℓ ($2 \leq \ell$), the secret-shared recursive oblivious transfer protocol outputs shares of*

$$\underbrace{V[V[\dots V[p_0] \dots]]}_{\ell} \quad (3)$$

without leaking V to P_0 and P_1 .

Protocol 1 Secret-shared Recursive Oblivious Transfer (ss-ROT).

Input: Public input: p_0
Input: Private input of server: $R^j[i]$ ($i = 0, \dots, N-1, j = 0, \dots, \ell-1$)

- 1: (Preparation by \mathcal{B}) \mathcal{B} generates and distributes $\llbracket R^j[i] \rrbracket_0$ and $\llbracket R^j[i] \rrbracket_1$ to P_0 and P_1
- 2: **for** $0 \leq j \leq \ell-2$ **do** ▷ **Step 1**
- 3: P_0 and P_1 obtain a position p_{j+1} by $p_{j+1} = \text{Reconst}(\llbracket R^j[p_j] \rrbracket_0, \llbracket R^j[p_j] \rrbracket_1)$.
- 4: **end for**
- 5: P_0 and P_1 output $\llbracket R^{\ell-1}[p_{\ell-1}] \rrbracket_0$ and $\llbracket R^{\ell-1}[p_{\ell-1}] \rrbracket_1$. ▷ **Step 2**

For simplicity, we denote the recursion of Eq. 3 by $V^{(\ell)}[p_0]$ (e.g., $V[V[p_0]]$ is denoted by $V^{(2)}[p_0]$). In our protocol, all the random values are uniformly generated from \mathbb{Z}_{2^n} .

DB Preparation Phase. \mathcal{B} generates $\ell-1$ random values $r^0, \dots, r^{\ell-2}$ and computes the following vectors $R^0, \dots, R^{\ell-1}$. Each vector R^j has N elements.

$$R^j[i] = \begin{cases} (V[i] + r^j) \bmod N & (j = 0) \\ (V[(i - r^{j-1}) \bmod N] + r^j) \bmod N & (1 \leq j \leq \ell-2) \\ (V[(i - r^{j-1}) \bmod N]) \bmod N & (j = \ell-1) \end{cases} \quad (4)$$

\mathcal{B} computes $\text{Share}(R^j[i])$ and sends $\llbracket R^j[i] \rrbracket_0$ and $\llbracket R^j[i] \rrbracket_1$ to P_0 and P_1 , for $i = 0, \dots, N-1$ and $j = 0, \dots, \ell-1$.

Search Phase. The Search phase consists of two steps and is described in Lines 2–5 of Protocol 1. The input is the initial position p_0 and shares of R . The output is $\llbracket V^{(\ell)}[p_0] \rrbracket$. An example of a search is illustrated in Figure 3.

$$\begin{array}{llll} V = (2, 0, 3, 1) & (2, 0, \color{red}{3}, 1) \mathbf{v}[2] & R^0 = (2+r^0, 0+r^0, \color{red}{3+r^0}, 1+r^0) = (3, 1, \color{red}{0}, 2) \mathbf{R}^0[2] \\ p_0=2, \ell=4 & (2, 0, \color{red}{3}, \color{red}{1}) \mathbf{v}[3] & R^1 = (\color{red}{1+r^1}, 2+r^1, 0+r^1, 3+r^1) = (\color{red}{3}, \color{red}{0}, 2, 1) \mathbf{R}^1[0] \\ r^0=1, r^1=2, r^2=1 & (2, \color{red}{0}, \color{red}{3}, 1) \mathbf{v}[1] & R^2 = (3+r^2, 1+r^2, 2+r^2, \color{red}{0+r^2}) = (0, 2, \color{red}{3}, \color{red}{1}) \mathbf{R}^2[3] \\ & (\color{red}{2}, 0, 3, 1) \mathbf{v}[0] & R^3 & = (1, \color{red}{2}, \color{red}{0}, 3) \mathbf{R}^3[1] \end{array}$$

Figure 3 Example of a search when $V = (2, 0, 3, 1)$, $p_0 = 2$, and $\ell = 4$. The goal is to compute $\llbracket V^{(4)}[2] \rrbracket = \llbracket 2 \rrbracket$. Here we assume \mathcal{B} generates $r^0 = 1, r^1 = 2, r^2 = 1$. In Step 1 of Search phase, P_0 and P_1 jointly compute $\text{Reconst}(\llbracket R^0[2] \rrbracket_0, \llbracket R^0[2] \rrbracket_1)$ to obtain $R^0[2] = 0$. ($R^0[2]$ is randomized by r^0 , so any element of V is leaked.) In a similar way, P_0 and P_1 compute $R^1[0] = 3$ and $R^2[3] = 1$. In Step 2, P_0 and P_1 output $\llbracket R^3[1] \rrbracket_0$ and $\llbracket R^3[1] \rrbracket_1$ respectively. Since $R^0[2] = V[2] + r^0$, $R^1[V[2] + r^0] = V[V[2] + r^0 - r^0] + r^1$, $R^2[V[V[2] + r^0] + r^1] = V[V[V[2] + r^0] + r^1 - r^1] + r^2$, and $R^3[V[V[V[2] + r^0] + r^1] + r^2] = V[V[V[V[2] + r^0] + r^1] + r^2 - r^2]$, ss-ROT successfully computes $\llbracket V^{(4)}[2] \rrbracket$.

3.1.1 Security and Complexities

► **Theorem 2.** *ss-ROT is correct and secure in the semi-honest model.*

Due to space limits, the proof is shown in Appendix D.

In the DB preparation phase, \mathcal{B} generates shares of V of length N for ℓ times. Therefore, time and communication complexities are $O(\ell N)$. For the Search phase, Reconst is computed ℓ times in Step 1. Since the time, communication, and round complexities of Reconst are $O(1)$, those of the Search phase become $O(\ell)$.

3.2 Secure LPM

Construction of Lookup Table. The goal is to find LPM securely. To apply FM-Index for a prefix search, the reverse string of T (i.e., \hat{T}) is used. The backward search of FM-Index is formulated by Eq. 1. If we precompute $\text{LF}_c(i, \hat{T})$ for $i = 0, \dots, N$ and $c \in \{A, T, G, C\}$, and store them in a lookup table that consists of four vectors: V_A , V_C , V_G , and V_T such that $V_c[i] = \text{LF}_c(i, \hat{T})$, Eq. 1 is replaced by the following table lookup

$$f_{k+1} = V_{w[k]}[f_k], \quad g_{k+1} = V_{w[k]}[g_k]. \quad (5)$$

I.e., starting with the initial interval $[f_0 = 0, g_0 = N)$, we can compute the match by recursively referring to the lookup table while $f < g$.

Protocol Overview. The key idea of Secure LPM is to refer to V by ss-ROT, i.e., P_0 and P_1 jointly refer to V ℓ times in a recursive manner. To achieve backward search, P_0 and P_1 need to select $V_x[\cdot]$ for each reference, where x is a query letter to be searched with. This is achieved by expressing the query letter by unary code (Eq. 7) and computing the inner product of Eq. 7 and $(V_A[\cdot], V_C[\cdot], V_G[\cdot], V_T[\cdot])$. To find LPM, P_0 and P_1 need to check $f = g$ for each reference. We use the subprotocol Equality to check it securely. Since V is randomized with different numbers for searching f and g , the difference of the random numbers is precomputed and removed securely upon the equality check. \mathcal{A} receives only the result of each equality check to know LPM. For examples, LPM is the prefix of length $i - 1$ when $f = g$ for the i -th reference. If $f \neq g$ for all references, LPM is the entire query.

DB Preparation Phase. \mathcal{B} creates a lookup table and generates the following 4ℓ vectors in a similar manner to ss-ROT. For simplicity, we denote the length of V_c by $N' = N + 1$.

$$R_{c,f}^j[i] = \begin{cases} (V_c[i] + r_f^j) \bmod N' & (j = 0) \\ (V_c[(i - r_f^{j-1}) \bmod N'] + r_f^j) \bmod N' & (1 \leq j < \ell) \end{cases} \quad (6)$$

$R_{c,f}^j[i]$ is used for computing the lower bound f of the interval $[f, g)$. We also generate $R_{c,g}^j[i]$ for the upper bound g . R consists of 8ℓ vectors, each of length N' . Since the longest match is found when $f = g$, \mathcal{B} also generates a vector $r'[j] = r_f^j - r_g^j$ that is used for equality check of f and g . Then, \mathcal{B} sends shares of $R_{c,f}^j[i]$, $R_{c,g}^j[i]$, and $r'[j]$ to P_0 and P_1 .

Search Phase. Protocol 2 describes the algorithm in detail. \mathcal{A} generates four vectors q_A , q_C , q_G , q_T , each of length ℓ , as follows.

$$q_c[j] = \begin{cases} 1 & (c = w[j]) \\ 0 & (c \neq w[j]) \end{cases} \quad (7)$$

For each j , $(q_A[j], q_C[j], q_G[j], q_T[j])$ encodes $w[j]$ (e.g., $(q_A[j], q_C[j], q_G[j], q_T[j]) = (1, 0, 0, 0)$ if $w[j] = A$). The aim of the encode is to compute $\llbracket R_x[j] \rrbracket = \llbracket \sum_{c \in \Sigma} q_c[j] \cdot R_c[j] \rrbracket$ when $w[j] = x$. Figure 4 illustrates an example of the table lookup.

\mathcal{A} generates shares of q_A , q_C , q_G , q_T and distributes them to P_0 and P_1 . P_0 and P_1 compute $\text{LF}_{w[j]}(f', \hat{T}) + r_f^j$ and $\text{LF}_{w[j]}(g', \hat{T}) + r_g^j$ in Lines 5–8 without leaking f' and g' , where $[f', g']$ corresponds to the match of $w[0, j]$ and \hat{T} . In Lines 10–12, the equality of f' and g' is examined for all rounds. Note that f_j and g_j are randomized by different values r_f^{j-1} and r_g^{j-1} in order to conceal f' and g' , so our protocol computes the equality of $f_j - g_j - r'[j - 1]$ and 0. In Lines 15–16, \mathcal{A} receives all the results of equality checks (i.e., $\llbracket o[1] \rrbracket^B, \dots, \llbracket o[\ell] \rrbracket^B$) from P_0 and P_1 , and knows LPM by reconstructing them. For example, if $w = \text{GCT}$ and $o = (0, 0, 1)$, \mathcal{A} knows that LPM is GC.

■ **Protocol 2** Secure LPM.

Input: Public input: $N, N' = N + 1, \ell, \Sigma = \{A, C, G, T\}, f_0 = 0, g_0 = N$
Input: Private input of user: query $q_c \in \{0, 1\}^\ell$ ($c \in \Sigma$)
Input: Private input of server: $R_{c,f}^j, R_{c,g}^j \in \mathbb{Z}_{N'}^{N'}$ ($c \in \Sigma, 0 \leq j < \ell$), $r' \in \mathbb{Z}_{N'}^\ell$

- 1: (Preparation by \mathcal{B}) \mathcal{B} distributes $\llbracket R_{c,f}^j \rrbracket, \llbracket R_{c,g}^j \rrbracket$ ($c \in \Sigma, 0 \leq j < \ell$), $\llbracket r' \rrbracket$ to P_0 and P_1
- 2: (Preparation by \mathcal{A}) \mathcal{A} distributes $\llbracket q_c \rrbracket$ ($c \in \Sigma$) to P_0 and P_1 .
- 3: (Computation by P_0 and P_1)
- 4: **for** $j = 0, \dots, \ell - 1$ **do**
- 5: $\llbracket f_j \rrbracket \leftarrow \sum_{c \in \Sigma} \text{MULT}(\llbracket R_{c,f}^j \rrbracket, \llbracket q_c \rrbracket)$ ▷ Select $\llbracket R_{w[j],f}^j \rrbracket$
- 6: $\llbracket g_j \rrbracket \leftarrow \sum_{c \in \Sigma} \text{MULT}(\llbracket R_{c,g}^j \rrbracket, \llbracket q_c \rrbracket)$ ▷ Select $\llbracket R_{w[j],g}^j \rrbracket$
- 7: $f_{j+1} \leftarrow \text{Reconst}(\llbracket f_j \rrbracket_0, \llbracket f_j \rrbracket_1)$ ▷ Update randomized lower bound
- 8: $g_{j+1} \leftarrow \text{Reconst}(\llbracket g_j \rrbracket_0, \llbracket g_j \rrbracket_1)$ ▷ Update randomized upper bound
- 9: **end for**
- 10: **for** $j = 1, \dots, \ell$ **parallelly do**
- 11: $\llbracket o[j] \rrbracket^B \leftarrow \text{Equality}(\llbracket f_j \rrbracket - \llbracket g_j \rrbracket - \llbracket r'[j-1] \rrbracket, \llbracket 0 \rrbracket)$ ▷ Equality check of upper and lower bounds.
- 12: **end for**
- 13: P_0 and P_1 send $\llbracket o \rrbracket_0^B, \llbracket o \rrbracket_1^B$ to \mathcal{A}
- 14: (Verification by \mathcal{A})
- 15: **for** $j = 1, \dots, \ell$ **do**
- 16: $o[j] \leftarrow \text{Reconst}(\llbracket o[j] \rrbracket_0^B, \llbracket o[j] \rrbracket_1^B)$
- 17: **end for**

Output: \mathcal{A} outputs $o[1], \dots, o[\ell]$ to determine LPM.

3.2.1 Security and Complexities

► **Theorem 3.** *Protocol 2 is correct and secure in the semi-honest model.*

Due to space limits, the proof is shown in Appendix E. \mathcal{A} may reveal T by making many queries. Such a problem is called output privacy. Although output privacy is outside of the scope of this paper, we should mention here that \mathcal{A} needs to make an unrealistically large number of queries for obtaining T by such a brute-force attack, considering that N is very long.

The DB preparation phase generates shares of $R_{c,f}^j$ and $R_{c,g}^j$ ($c \in \Sigma$ and $0 \leq j < \ell$); i.e., $8 \times \ell$ vectors of length N' . Therefore, the time and communication complexities are $O(\ell N)$. For the Search phase, MULT and Reconst are computed twice in Lines 4–9 for ℓ rounds and Equality is computed once in Lines 10–12 for ℓ rounds. Each time, the communication and round complexities of these subprotocols are $O(1)$, so those of the Search phase become $O(\ell)$.

3.3 Secure LMEM

Construction of Lookup Table. As described in Section 2.2, we can find a parent interval by a reference to LCP, PSV, and NSV. Therefore, in addition to V_c defined in Section 3.2, we prepare lookup tables that simply store all the outputs of them; i.e., $V_{cp}[i] = \text{LCP}[i]$, $V_{psv}[i] = \text{PSV}[i]$, and $V_{nsv}[i] = \text{NSV}[i]$.

$$\begin{array}{l}
V_A = (0, 0, 1, 1, 1) \quad q_A = (0, 0, 0) \quad R_A^0 = (0+r^0, 0+r^0, 1+r^0, 1+r^0, 1+r^0) \quad R_A^1 = (\dots, 1+r^1, \dots) \quad R_A^2 = (\dots, 0+r^2, \dots) \\
V_C = (1, 1, 1, 2, 2) \quad q_C = (0, 1, 0) \quad R_C^0 = (1+r^0, 1+r^0, 1+r^0, 2+r^0, 2+r^0) \quad R_C^1 = (\dots, 1+r^1, \dots) \quad R_C^2 = (\dots, 1+r^2, \dots) \\
V_G = (2, 2, 2, 2, 3) \quad q_G = (1, 0, 0) \quad R_G^0 = (2+r^0, 2+r^0, 2+r^0, 2+r^0, 3+r^0) \quad R_G^1 = (\dots, 2+r^1, \dots) \quad R_G^2 = (\dots, 2+r^2, \dots) \\
V_T = (3, 3, 3, 3, 4) \quad q_T = (0, 0, 1) \quad R_T^0 = (3+r^0, 3+r^0, 3+r^0, 3+r^0, 4+r^0) \quad R_T^1 = (\dots, 3+r^1, \dots) \quad R_T^2 = (\dots, 3+r^2, \dots)
\end{array}$$

$$\begin{array}{l}
w=GCT, \hat{T}=ACGT \quad (q_A[0], q_C[0], q_G[0], q_T[0]) \cdot (R_A^0[0], R_C^0[0], R_G^0[0], R_T^0[0]) = 2+r^0 \\
V_r[V_C[V_G[0]]] = 3 \quad (q_A[1], q_C[1], q_G[1], q_T[1]) \cdot (R_A^1[2+r^0], R_C^1[2+r^0], R_G^1[2+r^0], R_T^1[2+r^0]) = 1+r^1 \\
R_r^2[R^1[R^0[0]]] = 3+r^2 \quad (q_A[2], q_C[2], q_G[2], q_T[2]) \cdot (R_A^2[1+r^1], R_C^2[1+r^1], R_G^2[1+r^1], R_T^2[1+r^1]) = 3+r^2
\end{array}$$

■ **Figure 4** Example of a secure table lookup when $w = GCT$ and $\hat{T} = ACGT$. Only the lookup for a lower bound is shown. For simplicity, $R_{c,f}^j$ and r_f^j are denoted by R_c^j and r^j . $LF_{w[i]}(f_i, \hat{T})$ ($i = 0, 1, 2$) is computed by $V_G[0]$, $V_C[2]$, and $V_T[1]$. V is referenced securely by using R . $R_G^0[0]$ is computed by $\sum_{c \in \Sigma} q_c[0] \cdot R_c[0]$. $R_C^1[2+r^0]$ is computed by $\sum_{c \in \Sigma} q_c[1] \cdot R_c[2+r^0]$. $R_T^2[1+r^1]$ is computed by $\sum_{c \in \Sigma} q_c[2] \cdot R_c[1+r^1]$.

DB Preparation Phase. \mathcal{B} generates randomized vectors $R_{c,f}$, $R_{c,g}$ and $r^j[j] = r_f^j - r_g^j$ using the same algorithm in Section 3.2 for length 2ℓ . As shown in Eq. 2, V_{lcp} is referred by the upper and lower bounds of $[f, g]$. Therefore, \mathcal{B} generates following circular permutations of V_{lcp} such that $W_{l,f}$ and $R_{c,f}$, and $W_{l,g}$ and $R_{c,g}$, are permuted by the same random values, respectively. I.e.,

$$W_{l,x}^j[i] = \begin{cases} V_{lcp}[i] & (j = 0) \\ V_{lcp}[(i - r_x^{j-1}) \bmod N] & (1 \leq j < 2\ell) \end{cases},$$

where x is either f or g . V_{psv} is referred by both f and g , and is plugged in to f . Therefore, \mathcal{B} generates $W_{p,f}^j$ and $W_{p,g}^j$ such that both of them are randomized by r_f^j , and $W_{p,f}^j$ is permuted by r_f^{j-1} and $W_{p,g}^j$ is permuted by r_g^{j-1} as follows.

$$W_{p,f}^j[i] = \begin{cases} (V_{psv}[i] + r_f^j) \bmod N & (j = 0) \\ (V_{psv}[(i - r_f^{j-1}) \bmod N] + r_f^j) \bmod N & (1 \leq j < 2\ell) \end{cases}$$

$$W_{p,g}^j[i] = \begin{cases} (V_{psv}[i] + r_g^j) \bmod N & (j = 0) \\ (V_{psv}[(i - r_g^{j-1}) \bmod N] + r_g^j) \bmod N & (1 \leq j < 2\ell) \end{cases}$$

Similarly, V_{nsv} is referred by both f and g , and is plugged in to g . Therefore, \mathcal{B} generates $W_{n,f}^j[i]$ and $W_{n,g}^j[i]$ as follows.

$$W_{n,f}^j[i] = \begin{cases} (V_{nsv}[i] + r_f^j) \bmod N & (j = 0) \\ (V_{nsv}[(i - r_f^{j-1}) \bmod N] + r_f^j) \bmod N & (1 \leq j < 2\ell) \end{cases}$$

$$W_{n,g}^j[i] = \begin{cases} (V_{nsv}[i] + r_g^j) \bmod N & (j = 0) \\ (V_{nsv}[(i - r_g^{j-1}) \bmod N] + r_g^j) \bmod N & (1 \leq j < 2\ell) \end{cases}$$

\mathcal{B} distributes shares of $R_{c,f}$, $R_{c,g}$, r^j , $W_{l,f}$, $W_{l,g}$, $W_{p,f}$, $W_{p,g}$, $W_{n,f}$, and $W_{n,g}$ to P_0 and P_1 .

Search Phase. Protocol 3 describes the algorithm in detail. \mathcal{A} generates query vectors q_A , q_C , q_G , q_T by Eq. 7 and distributes shares of the vectors to P_0 and P_1 . In Line 6 of Protocol 3, $[\hat{f}, \hat{g}]$ is computed by the reference to R (i.e., a search based on a backward search) similarly to Lines 5–6 of Protocol 2. In Line 9, $[f_{ex}, g_{ex}]$ is computed by the reference to W (i.e., a search based on LCP, PSV and NSV). In Line 11, the interval is updated by either $[\hat{f}, \hat{g}]$ or $[f_{ex}, g_{ex}]$ based on the result of $f' = g'$ in Line 7, where $[f', g']$ corresponds to the interval that corresponds to a substring match.

In each round, we need to know a query letter to be searched with, so we need to maintain the right end position of the match in the query. The position moves toward the right while the match is extended, but remains the same when the interval is updated based on PSV and NSV. To memorize the position, we prepare shares of a unit bit vector u of length ℓ , in which the position t is memorized as $u[t] = 1$ and $u[i \neq t] = 0$. In Lines 14–18, u remains the same if the interval is updated based on PSV and NSV, and $u = (u, [\ell - 1], u[0], u[1], \dots, u[\ell - 2])$ otherwise. In Lines 19–21, the inner product of q_c ($c \in \Sigma$) and u becomes the encode of $w[t]$ that is used for the next round.

We also maintain the left end position of the match. While the match is extended, the position remains the same and it moves toward the right when the interval is updated by $[f_{ex}, g_{ex}]$. The new left end position can be computed by $p + m - c$ where p is the current position, m is the length of the current match, and c is the lcp-value of $[f_{ex}, g_{ex}]$ (i.e., the longest common prefix length of suffixes contained in $[f_{ex}, g_{ex}]$). The position is computed in Line 23. The match length is incremented by 1 for each extension. When the interval is updated by $[f_{ex}, g_{ex}]$, the match length is reduced to the lcp-value of $[f_{ex}, g_{ex}]$, which is computed by $\max(\text{LCP}[f], \text{LCP}[g])$. The match length is computed in Line 22. In Line 25, the longest match length and the corresponding left end position are updated. After all the positions in the query have been examined, LMEM and its left end position are sent to \mathcal{A} in Line 27.

3.3.1 Security and Complexities

► **Theorem 4.** *Protocol 3 is correct and secure in the semi-honest model.*

Due to space limits, the proof is shown in Appendix F.

The DB preparation phase generates shares of $R_{c,f}^j$ and $R_{c,g}^j$ ($c \in \Sigma$, $0 \leq j < \ell$) and $W_{x,f}^j$ and $W_{x,g}^j$ ($x \in \{l, p, n\}$ and $0 \leq j < \ell$); $14 \times \ell$ vectors of length $N + 1$. Therefore, the time and communication complexities are $O(\ell N)$. For the Search phase, MULT is computed ℓ times in parallel in Lines 15–16. (These are not dependent on each other.) In Line 20, MULT is computed ℓ times in parallel, and Line 20 is computed in parallel four times in Lines 19–21. Lines 15–16 and Lines 19–21 are repeated for $2\ell - 1$ rounds. Other subprotocols are also computed for $2\ell - 1$ rounds. The time, communication, and round complexities are $O(1)$ for MULT, and independent computation of MULT for ℓ times does not increase the round complexity. The time, communication and round complexities are $O(1)$ for the other subprotocols used in Protocol 3. Therefore, the complexities of the Search phase are $O(\ell^2)$ for time and communication, and $O(\ell)$ for the number of rounds.

4 Experiment

We implemented Protocol 2 (Secure LPM) and Protocol 3 (Secure LMEM). For comparison, we also implemented baseline protocols (Baseline (LPM) and Baseline (LMEM)). Details of the baseline protocols are provided in Appendix C. All protocols were implemented by Python 3.5.2. The dataset was created from Chromosome 1 of the human genome. We extracted substrings of length $N = 10^3, 10^4, 10^5, 10^6$, and 10^7 for databases, and $\ell = 10, 25, 50, 75$, and 100 for queries. Share was run with $n = 16$ and $n = 32$ for $N < 10^5$ and $10^5 \leq N$ in the proposed protocols, and $n = 1$ for a Boolean share and $n = 8$ for an arithmetic share in the baseline protocols. We did not implement a data transfer module, and each protocol is implemented as a single program. Therefore, the search time of the protocols was measured by the time consumed by either one of P_0 and P_1 . To assess the influence of communication

Protocol 3 Secure LMEM.

Input: Public input: $N, N' = N + 1, \ell, \Sigma = \{A, T, G, C\}, f_0 = 0, g_0 = N$
Input: Private input of user: query $q_c \in \{0, 1\}^\ell$ ($c \in \Sigma$)
Input: Private input of server: $R_{c,f}^j, R_{c,g}^j, W_{l,f}^j, W_{l,g}^j, W_{p,f}^j, W_{p,g}^j, W_{n,f}^j, W_{n,g}^j \in \mathbb{Z}_{N'}^N$ ($c \in \Sigma, 0 \leq j < \ell$), $r' \in \mathbb{Z}_{N'}^\ell$

- 1: (Preparation by \mathcal{B}) \mathcal{B} generates shares of input vectors. \mathcal{B} also generates shares of variables: $\llbracket u \rrbracket = \llbracket (1, 0, \dots, 0) \rrbracket$, $\llbracket p_{max} \rrbracket = \llbracket 0 \rrbracket$, $\llbracket p \rrbracket = \llbracket 0 \rrbracket$, $\llbracket m_{max} \rrbracket = \llbracket 0 \rrbracket$, $\llbracket m \rrbracket = \llbracket 0 \rrbracket$. All shares are sent to P_0 and P_1 .
- 2: (Preparation by \mathcal{A}) \mathcal{A} generates and distributes $\llbracket q_c \rrbracket$ ($c \in \Sigma$) to P_0 and P_1 .
- 3: (Computation by P_0 and P_1)
- 4: Set shares of the initial letter $\llbracket z_c \rrbracket = \llbracket q_c[0] \rrbracket$ ($c \in \Sigma$).
- 5: **for** $j = 0, \dots, 2\ell - 1$ **do**
- 6: $\llbracket \hat{f}_j \rrbracket \leftarrow \sum_{c \in \Sigma} \text{MULT}(\llbracket R_{c,f}^j[f_j] \rrbracket, \llbracket z[c] \rrbracket)$, $\llbracket \hat{g}_j \rrbracket \leftarrow \sum_{c \in \Sigma} \text{MULT}(\llbracket R_{c,g}^j[g_j] \rrbracket, \llbracket z[c] \rrbracket)$
- 7: $\llbracket e1 \rrbracket^B \leftarrow \text{Equality}(\llbracket \hat{f}_j \rrbracket - \llbracket \hat{g}_j \rrbracket - \llbracket r'[j-1] \rrbracket, \llbracket 0 \rrbracket)$, $\llbracket e1 \rrbracket \leftarrow \text{B2A}(\llbracket e1 \rrbracket^B)$ \triangleright If $f = g$.
- 8: $\llbracket e2 \rrbracket^B \leftarrow \text{Comp}(\llbracket W_{l,f}^j[f_j] \rrbracket, \llbracket W_{l,g}^j[g_j] \rrbracket)$, $\llbracket e2 \rrbracket \leftarrow \text{B2A}(\llbracket e2 \rrbracket^B)$ \triangleright If $\text{LCP}[f] < \text{LCP}[g]$
- 9: $\llbracket f_{ex} \rrbracket \leftarrow \text{Choose}(\llbracket W_{p,g}^j[g_j] \rrbracket, \llbracket W_{p,f}^j[f_j] \rrbracket, \llbracket e2 \rrbracket)$, $\llbracket g_{ex} \rrbracket \leftarrow \text{Choose}(\llbracket W_{n,g}^j[g_j] \rrbracket, \llbracket W_{n,f}^j[f_j] \rrbracket, \llbracket e2 \rrbracket)$
- 10: $\llbracket l_{ex} \rrbracket \leftarrow \text{Choose}(\llbracket W_{l,g}^j[g_j] \rrbracket, \llbracket W_{l,f}^j[f_j] \rrbracket, \llbracket e2 \rrbracket)$
- 11: $\llbracket f_{j+1} \rrbracket \leftarrow \text{Choose}(\llbracket f_{ex} \rrbracket, \llbracket \hat{f}_j \rrbracket, \llbracket e1 \rrbracket)$, $\llbracket g_{j+1} \rrbracket \leftarrow \text{Choose}(\llbracket g_{ex} \rrbracket, \llbracket \hat{g}_j \rrbracket, \llbracket e1 \rrbracket)$
- 12: $f_{j+1} \leftarrow \text{Reconst}(\llbracket f_{j+1} \rrbracket_0, \llbracket f_{j+1} \rrbracket_1)$, $g_{j+1} \leftarrow \text{Reconst}(\llbracket g_{j+1} \rrbracket_0, \llbracket g_{j+1} \rrbracket_1)$ \triangleright Update f, g
- 13: $\llbracket e' \rrbracket \leftarrow \text{B2A}(\text{ADD}(\llbracket e1 \rrbracket^B, \llbracket 1 \rrbracket^B))$
- 14: **for** $i = 0, \dots, \ell - 1$ **parallely do** \triangleright Maintain right end of the match.
- 15: $\llbracket u1[i] \rrbracket \leftarrow \text{MULT}(\llbracket u[i] \rrbracket, \llbracket e1 \rrbracket)$ $\triangleright u1 = u$ and $u2 = (0, \dots, 0)$ if $\hat{f} - \hat{g} - r' = 0$,
- 16: $\llbracket u2[i] \rrbracket \leftarrow \text{MULT}(\llbracket u[i] \rrbracket, \llbracket e' \rrbracket)$ $u1 = (0, \dots, 0)$ and $u2 = u$ otherwise.
- 17: $\llbracket u[i] \rrbracket \leftarrow \text{ADD}(\llbracket u2[(i-1) \bmod \ell] \rrbracket, \llbracket u1[i] \rrbracket)$ $\triangleright u$ is incremented iff. $\hat{f} - \hat{g} - r' \neq 0$.
- 18: **end for**
- 19: **for** $c \in \Sigma$ **parallely do**
- 20: $\llbracket z_c \rrbracket \leftarrow \sum_{i \in \{0, \dots, \ell-1\}} \text{MULT}(\llbracket q_c[i] \rrbracket, \llbracket u[i] \rrbracket)$ \triangleright Select next letter to be searched with.
- 21: **end for**
- 22: $\llbracket m' \rrbracket \leftarrow \text{Choose}(\llbracket l_{ex} \rrbracket, \text{ADD}(\llbracket 1 \rrbracket, \llbracket m \rrbracket), \llbracket e1 \rrbracket)$ \triangleright Calculate match length
- 23: $\llbracket p \rrbracket \leftarrow \text{Choose}(\text{ADD}(\text{ADD}(\llbracket p \rrbracket, \llbracket m \rrbracket), \llbracket -l_{ex} \rrbracket), \llbracket p \rrbracket, \llbracket e1 \rrbracket)$ \triangleright Update left end position of match
- 24: $\llbracket e3 \rrbracket \leftarrow \text{B2A}(\text{Comp}(\llbracket m' \rrbracket, \llbracket m_{max} \rrbracket))$, $\llbracket m \rrbracket \leftarrow \llbracket m' \rrbracket$
- 25: $\llbracket m_{max} \rrbracket \leftarrow \text{Choose}(\llbracket m_{max} \rrbracket, \llbracket m \rrbracket, \llbracket e3 \rrbracket)$, $\llbracket p_{max} \rrbracket \leftarrow \text{Choose}(\llbracket p_{max} \rrbracket, \llbracket p \rrbracket, \llbracket e3 \rrbracket)$ \triangleright Update max
- 26: **end for**
- 27: P_0 and P_1 send $\llbracket m_{max} \rrbracket_0, \llbracket m_{max} \rrbracket_1, \llbracket p_{max} \rrbracket_0$, and $\llbracket p_{max} \rrbracket_1$ to \mathcal{A}
- 28: (Verification by \mathcal{A}) $\text{max} \leftarrow \text{Reconst}(\llbracket m_{max} \rrbracket_0, \llbracket m_{max} \rrbracket_1)$ and $p_{max} \leftarrow \text{Reconst}(\llbracket p_{max} \rrbracket_0, \llbracket p_{max} \rrbracket_1)$.

Output: \mathcal{A} outputs m_{max} and p_{max} to report LMEM.

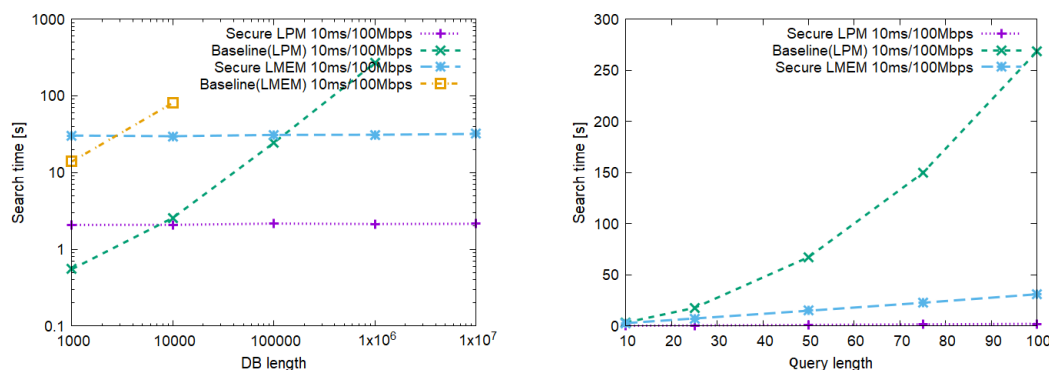
on a realistic environment, we theoretically estimated delays caused by network bandwidth and latency. We assume three environments: LAN (0.2 ms/10 Gbps), WAN₁ (10 ms/100 Mbps), and WAN₂ (50 ms/10 Mbps). During the run of Search phase, we stored all the data that were transferred from P_0 to P_1 in a file and measured the file size as an actual communication size. Note that the communication is symmetric and data transfer size from P_0 to P_1 is equal to that from P_1 to P_0 . Based on the data transfer size D byte, we estimate the communication delay by $D/k + 1000 \times eT$, where k is bandwidth, e is latency and T is a round of communication. All the protocols were run with a single thread on the same machine equipped with Intel Xeon 2.2 GHz CPU and 256 GB memory. We also tested the

■ **Table 3** Offline time (Time), offline size (Size), DB preparation time (Time), DB preparation size (Size), Search time on a local machine (Time), Search communication size (Size), estimated Search time for three environments: LAN (0.2 ms/10 Gbps), WAN₁ (10 ms/100 Mbps), and WAN₂ (50 ms/10 Mbps), for $N = 10^4$ (only for Baseline (LMEM)), $10^5, 10^6, 10^7$, and $\ell = 100$. The size unit is MB and the time unit is sec except for the cell describing “20 h<”.

	N	Offline		DB preparation		Search		Estimated time on network		
		Time	Size	Time	Size	Time	Size	LAN	WAN ₁	WAN ₂
Secure	10^5	0.166	0.013	123	305	0.141	0.010	0.181	2.162	10.249
LPM	10^6	0.141	0.013	1248	3051	0.113	0.010	0.153	2.134	10.221
(ours)	10^7	0.150	0.013	12628	30517	0.126	0.010	0.167	2.147	10.234
[30]	10^5	-	-	-	-	691	163	691	707	838
	10^6	-	-	-	-	7817	517	7818	7863	8261
	10^7	-	-	-	-	20 h<	-	-	-	-
Baseline (LPM)	10^5	3995	184	0.146	0.095	13	122	13	24	118
	10^6	38767	1841	1.522	0.954	164	1227	165	268	1196
	10^7	20 h<	-	-	-	-	-	-	-	-
Secure	10^5	7.619	1.704	435	1068	4.817	0.999	5.577	42.900	195.654
LMEM	10^6	7.882	1.704	4467	10681	4.926	0.999	5.686	43.009	195.763
(ours)	10^7	8.457	1.704	46384	106811	5.740	0.999	6.501	43.824	196.578
Baseline (LMEM)	10^4	12747	611	0.015	0.010	46	407	46	80	389
	10^5	20 h<	-	-	-	-	-	-	-	-

C++ implementation of [30], which is based on AHE. The algorithm for LPM in [28] for the string with $|\Sigma| \leq 4$ (e.g., genome sequence) is the same as [30]. [30] is implemented as a server-client software, and the client and the server were run with individual single threads on the same machine. Therefore, the results of [30] do not include delays caused by bandwidth limitation and latency, so we also estimated delays based on the data transfer size and round of communication in the same manner. Each run of the program was terminated if the total runtime of all phases exceeded 20 hours.

Comparison to Baseline Protocols. Table 3 shows the offline time and size, DB preparation time and size, and Search time and communication size for $N = 10^5, 10^6, 10^7$, and $\ell = 100$. It also shows the result of Baseline (LMEM) for $N = 10^4$, as the runs for $N > 10^4$ did not finish within 20 hours. The Search times and communication sizes of Secure LPM and Secure LMEM are several orders of magnitudes faster and smaller than those of Baseline (LPM) and Baseline (LMEM). Since the round and communication complexities of the proposed protocols do not depend on N , their estimated Search time remains small even on WAN environments. The left panel of Figure 5 shows the estimated Search time on WAN₁ for $N = 10^3, 10^4, \dots, 10^7$ and $\ell = 100$. The times of Secure LPM and Secure LMEM do not increase, while those of the baseline protocols increase linearly to N . The right panel of Figure 5 shows the estimated Search time on WAN₁ for $\ell = 10, 25, \dots, 100$ for $N = 10^6$. We can not show the results of Baseline (LMEM) because none of its runs were finished within the time limit. As shown in the graph, the time of Secure LPM increases linearly to ℓ and that of Baseline (LPM) increases proportionally to ℓ^2 , which are in good agreement with the theoretical complexities in Table 2. According to the graph, the time of Secure LMEM also increases linearly to ℓ though its time and communication complexities are $O(\ell^2)$. This is because the CPU times are much smaller than the delays caused by network latency that are influenced by the round complexity $O(\ell)$.



■ **Figure 5** Estimated time (actual search time on a local machine + estimated data transfer time) for various N and ℓ .

We have preliminary results for testing Secure LPM and Baseline (LPM) on the actual network (10 ms/100 Mbps). The results are 40 sec for Secure LPM and 1739 sec for Baseline (LPM) when $N = 10^6$. Though both of the preliminary implementations have room for improvement in the performance of data transfer, the results also indicate that our protocol outperforms the baseline protocol and the previous study.

The time and size of Secure LPM and Secure LMEM are several orders of magnitude better than those of the baseline protocols for the offline phase, and vice versa for the DB preparation phase. The total time of the offline and DB preparation phases of our protocols are more than one order magnitude faster than that of baseline protocols. For the total size of the offline and DB preparation phases, Secure LMEM was better than Baseline (LMEM), but Baseline (LPM) was better than Secure LPM though the complexity is better for Secure LPM. This is because the majority of the shares were Boolean in the baseline protocols, while all of the shares were arithmetic in the proposed protocols.

Comparison to [30]. [30] is a two-party MPC based on AHE. Each homomorphic operation is time consuming and has no offline and DB preparation phases. As shown in Table 3, the Search time of Secure LPM is four orders of magnitude faster than [30] for $N = 10^6$. Since time complexity of [30] includes a factor of N , the difference in Search time becomes greater as N becomes large. Moreover, our protocols have a further advantage in communication for a query response when the network environment is poor, as the round complexity of [30] and our protocols are the same while [30] requires $O(\sqrt{N})$ communication size. The entire runtimes including all the phases are still six times faster for $N = 10^5$ and $N = 10^6$. We can compute LMEM by examining [30] for all the positions in a query string, but this approach consumed 3406 sec and 2.6 GByte of communication for $N = 10^4$.

5 Discussion

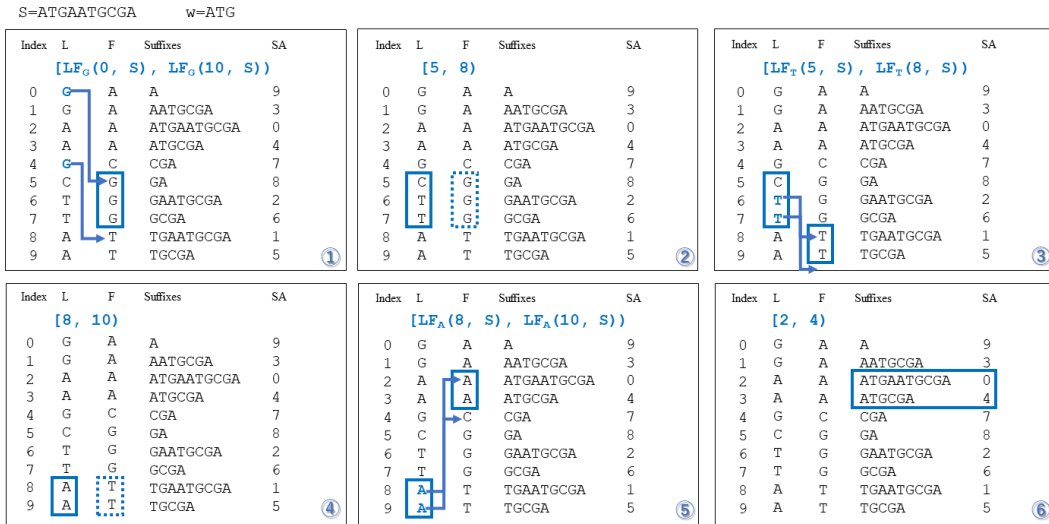
As clearly shown by the results, Search time of the proposed protocols are significantly efficient. Considering the importance of query response time for real applications, it is realistic to reduce Search time at the cost of DB preparation time. Since the total times for offline and DB preparation phases of the proposed protocols were significantly better than those of the well-designed baseline protocols, we consider the trade-off between Search and DB preparation times in our approach to be efficient. For further reduction of DB preparation time, parallelizing the share generation is a feasible approach. Regarding the DB preparation

phase, the data transfer between the server and the computing nodes is problematic when the number of queries and the length of the database are large. One potential solution to mitigate the problem is to use an AES-based random number generation that is similar to the technique used in [1]. To explain it briefly, when the server needs to distribute a share of x , (1) the server and P_0 generate the same randomness r using a pre-shared key and a pseudorandom function, and (2) the server computes $x - r$ and sends it to P_1 . Although P_0 's computation cost increases, we can remove the data transfer from the server to P_0 . In our protocols, the generation of shares in the DB preparation phase cannot be outsourced because they are dependent on the database. Designing an efficient algorithm to outsource the share generation is an important open question.

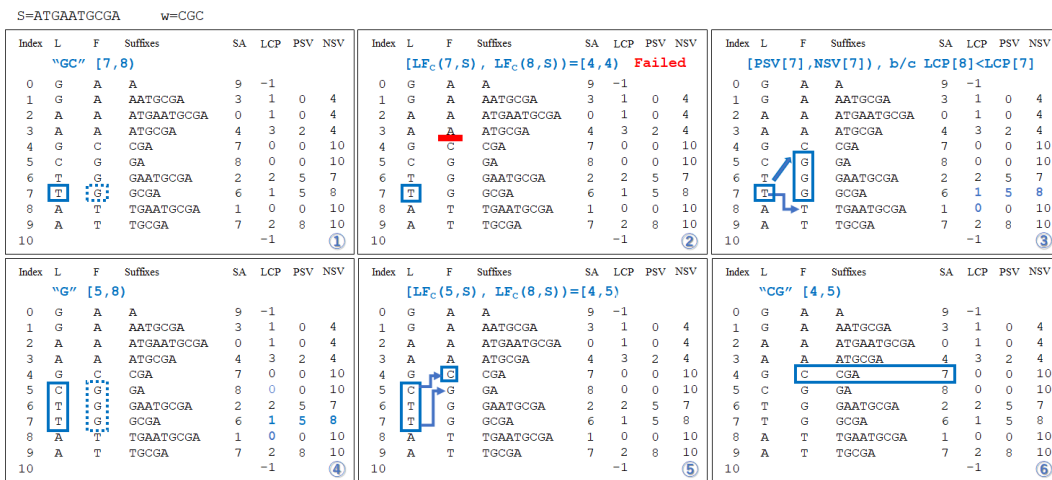
References

- 1 Toshinori Araki, Jun Furukawa, Yehuda Lindell, Ariel Nof, and Kazuma Ohara. High-throughput semi-honest secure three-party computation with an honest majority. In *Proc. of CCS 2016*, pages 805–817, 2016. doi:10.1145/2976749.2978331.
- 2 Gilad Asharov, Shai Halevi, Yehuda Lindell, and Tal Rabin. Privacy-preserving search of similar patients in genomic data. *PoPETs*, 2018(4):104–124, 2018. doi:10.1515/popets-2018-0034.
- 3 Md Momin Al Aziz, Md. Nazmus Sadat, Dima Alhadidi, Shuang Wang, Xiaoqian Jiang, Cheryl L. Brown, and Noman Mohammed. Privacy-preserving techniques of genomic data - a survey. *Briefings Bioinform.*, 20(3):887–895, 2019.
- 4 Pierre Baldi, Roberta Baronio, Emiliano De Cristofaro, Paolo Gasti, and Gene Tsudik. Countering GATTACA: efficient and secure testing of fully-sequenced human genomes. In *Proc. of CCS 2011*, pages 691–702, 2011. doi:10.1145/2046707.2046785.
- 5 Donald Beaver. Efficient multiparty protocols using circuit randomization. In *Proc. of CRYPTO 1991*, pages 420–432, 1991. doi:10.1007/3-540-46766-1_34.
- 6 Yangyi Chen, Bo Peng, XiaoFeng Wang, and Haixu Tang. Large-scale privacy-preserving mapping of human genomic sequences on hybrid clouds. In *Proc. of NDSS 2012*, 2012. URL: <https://www.ndss-symposium.org/ndss2012/large-scale-privacy-preserving-mapping-human-genomic-sequences-hybrid-clouds>.
- 7 Ke Cheng, Yantian Hou, and Liangmin Wang. Secure similar sequence query on outsourced genomic data. In *Proc. of AsiaCCS 2018*, pages 237–251, 2018. doi:10.1145/3196494.3196535.
- 8 Jung Hee Cheon, Miran Kim, and Kristin E. Lauter. Homomorphic computation of edit distance. In *Proc. of FC 2015*, pages 194–212, 2015. doi:10.1007/978-3-662-48051-9_15.
- 9 Richard Durbin. Efficient haplotype matching and storage using the positional burrows–wheeler transform (pbwt). *Bioinformatics*, 30(9):1266–1272, 2014.
- 10 Yaniv Erlich and Arvind Narayanan. Routes for breaching and protecting genetic privacy. *Nature Reviews Genetics*, 15(6):409–421, 2014.
- 11 Paolo Ferragina and Giovanni Manzini. Opportunistic data structures with applications. In *Proc. of FOCS 2000*, pages 390–398, 2000. doi:10.1109/SFCS.2000.892127.
- 12 Johannes Fischer, Veli Mäkinen, and Gonzalo Navarro. An(other) entropy-bounded compressed suffix tree. In *Proc. of CPM 2008*, pages 152–165, 2008. doi:10.1007/978-3-540-69068-9_16.
- 13 Marc Fiume, Miroslav Cupak, Stephen Keenan, Jordi Rambla, Sabela de la Torre, Stephanie OM Dyke, Anthony J Brookes, Knox Carey, David Lloyd, Peter Goodhand, et al. Federated discovery and sharing of genomic data using beacons. *Nature biotechnology*, 37(3):220–224, 2019.
- 14 Oded Goldreich. *The Foundations of Cryptography - Volume 2, Basic Applications*. Cambridge University Press, 2004. doi:10.1017/CB09780511721656.
- 15 Yan Huang, David Evans, Jonathan Katz, and Lior Malka. Faster secure two-party computation using garbled circuits. In *Proc. of USENIX 2011*, 2011. URL: http://static.usenix.org/events/sec11/tech/full_papers/Huang.pdf.

- 16 Y. Ishimaki, H. Imabayashi, K. Shimizu, and H. Yamana. Privacy-preserving string search for genome sequences with the bootstrapping optimization. In *Proc. of IEEE Big Data 2016*, pages 3989–3991, 2016.
- 17 Somesh Jha, Louis Kruger, and Vitaly Shmatikov. Towards practical privacy for genomic computation. In *Proc. of IEEE S&P 2000*, pages 216–230, 2008. doi:10.1109/SP.2008.34.
- 18 Md Safiur Rahman Mahdi, Md Momin Al Aziz, Noman Mohammed, and Xiaoqian Jiang. Privacy-preserving string search on encrypted genomic data using a generalized suffix tree. *Informatics in Medicine Unlocked*, 23:100525, 2021.
- 19 Payman Mohassel, Ostap Orobets, and Ben Riva. Efficient server-aided 2pc for mobile phones. *PoPETs*, 2016(2):82–99, 2016. doi:10.1515/popets-2016-0006.
- 20 Payman Mohassel and Yupeng Zhang. Secureml: A system for scalable privacy-preserving machine learning. In *Proc. of IEEE S&P 2017*, pages 19–38, 2017. doi:10.1109/SP.2017.12.
- 21 Muhammad Naveed, Erman Ayday, Ellen W Clayton, Jacques Fellay, Carl A Gunter, Jean-Pierre Hubaux, Bradley A Malin, and XiaoFeng Wang. Privacy in the genomic era. *ACM Computing Surveys (CSUR)*, 48(1):1–44, 2015.
- 22 Koji Nuida, Satsuya Ohata, Shigeo Mitsunari, and Nuttapong Attrapadung. Arbitrary univariate function evaluation and re-encryption protocols over lifted-elgamal type ciphertexts. *IACR Cryptology ePrint Archive*, 2019:1233, 2019. URL: <https://eprint.iacr.org/2019/1233>.
- 23 Satsuya Ohata and Koji Nuida. Communication-efficient (client-aided) secure two-party protocols and its application. In *proc. of FC 2020*, pages 369–385, 2020.
- 24 Anthony A Philippakis, Danielle R Azzariti, Sergi Beltran, Anthony J Brookes, Catherine A Brownstein, Michael Brudno, Han G Brunner, Orion J Buske, Knox Carey, Cassie Doll, et al. The matchmaker exchange: a platform for rare disease gene discovery. *Human mutation*, 36(10):915–921, 2015.
- 25 Victoria Popic and Serafim Batzoglou. A hybrid cloud read aligner based on minhash and kmer voting that preserves privacy. *Nature communications*, 8(1):1–7, 2017.
- 26 Thomas Schneider and Oleksandr Tkachenko. EPISODE: efficient privacy-preserving similar sequence queries on outsourced genomic databases. In *Proc. of AsiaCCS 2019*, pages 315–327, 2019. doi:10.1145/3321705.3329800.
- 27 Adi Shamir. How to share a secret. *Communications of the ACM*, 22(11):612–613, 1979.
- 28 Kana Shimizu, Koji Nuida, and Gunnar Rätsch. Efficient privacy-preserving string search and an application in genomics. *Bioinformatics*, 32(11):1652–1661, 2016. doi:10.1093/bioinformatics/btw050.
- 29 Katerina Sotiraki, Esha Ghosh, and Hao Chen. Privately computing set-maximal matches in genomic data. *BMC Medical Genomics*, 13(7):1–8, 2020.
- 30 H. Sudo, M. Jimbo, K. Nuida, and K. Shimizu. Secure wavelet matrix: Alphabet-friendly privacy-preserving string search for bioinformatics. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 16(5):1675–1684, 2019.
- 31 Xiao Shaun Wang, Yan Huang, Yongan Zhao, Haixu Tang, XiaoFeng Wang, and Diyue Bu. Efficient genome-wide, privacy-preserving similar patient query based on private edit distance. In *Proc. of CCS 2015*, pages 492–503, 2015. doi:10.1145/2810103.2813725.
- 32 Masaya Yasuda, Takeshi Shimoyama, Jun Kogure, Kazuhiro Yokoyama, and Takeshi Koshihara. Secure pattern matching using somewhat homomorphic encryption. In Ari Juels and Bryan Parno, editors, *Proc. of CCSW’13*, pages 65–76, 2013. doi:10.1145/2517488.2517497.
- 33 R. Zhu and Y. Huang. Efficient and precise secure generalized edit distance and beyond. *IEEE Transactions on Dependable and Secure Computing*, pages 1–1, 2020.



■ **Figure A1** An example of search by FM-Index.



■ **Figure A2** An example of search by FM-ndex, LCP array, PSV and NSV.

A Examples of a Search with FM-Index and Auxiliary Data Structures

Let us show examples of a search with FM-Index, LCP array, PSV and NSV. In addition to the data structures defined in Section 2.2, we also define a string F such that $F[i] = S[SA[i]]$. For the case of $S = ATGAATGCGA$, the indices become $SA = (9, 3, 0, 4, 7, 8, 2, 6, 1, 5)$, $L = GGAAGCTTAA$, and $F = AAAACGGT$. Figure A1 illustrates the example of a backward search to find the longest suffix of the query (ATG) that matches the database, and Figure A2 illustrates the search for MEMS with the query (CGC) by using LCP array, PSV, and NSV. As shown in the upper center panel of Figure A2, the search failed when the backward search with ‘C’ after finding the interval [7, 8) that corresponds to GC. Since $LCP[8] < LCP[7]$, the parent lcp-interval becomes $[PSV[7] = 5, NSV[7] = 8)$, which corresponds to ‘G’. The match CG is then searched with the backward search with ‘C’ from the parent lcp-interval.

B Semi-Honest Security

Here, we recall the simulation-based security notion in the presence of semi-honest adversaries (for two-party computation), as in [14].

► **Definition 5.** Let $f : (\{0, 1\}^*)^2 \rightarrow (\{0, 1\}^*)^2$ be a probabilistic 2-ary functionality and $f_i(\vec{x})$ denote the i -th element of $f(\vec{x})$ for $\vec{x} = (x_0, x_1) \in (\{0, 1\}^*)^2$ and $i \in \{0, 1\}$; $f(\vec{x}) = (f_0(\vec{x}), f_1(\vec{x}))$. Let Π be a 2-party protocol to compute the functionality f . The view of party P_i for $i \in \{0, 1\}$ during an execution of Π on input $\vec{x} = (x_0, x_1) \in (\{0, 1\}^*)^2$ where $|x_0| = |x_1|$, denoted by $\text{VIEW}_i^\Pi(\vec{x})$, consists of $(x_i, r_i, m_{i,1}, \dots, m_{i,t})$, where x_i represents P_i 's input, r_i represents its internal random coins, and $m_{i,j}$ represents the j -th message that P_i has received. The output of all parties after an execution of Π on input \vec{x} is denoted as $\text{OUTPUT}^\Pi(\vec{x})$. Then, for each party P_i , we say that Π privately computes f in the presence of semi-honest corrupted party P_i if there exists a probabilistic polynomial-time algorithm \mathcal{S} such that

$$\{(\mathcal{S}(i, x_i, f_i(\vec{x})), f(\vec{x}))\} \equiv \{(\text{VIEW}_i^\Pi(\vec{x}), \text{OUTPUT}^\Pi(\vec{x}))\},$$

where the symbol \equiv means that the two probability distributions are statistically indistinguishable.

As described in [14], the composition theorem for the semi-honest model holds; that is, any protocol is privately computed as long as its subroutines are privately computed.

C Our Secure Baseline LPM and LMEM

In this section, we show our secure baseline LCP and LMEM based on secret sharing. We explain how to construct LCP, since we can obtain LMEM by (parallelly) executing LCP for all positions in the query. Note that $\vec{x} = (x_1, x_2, \dots)$, \vec{x}_i denotes an i -th element of \vec{x} , $[\vec{t}] = ([\vec{t}]_0, [\vec{t}]_1)$, and $(|\vec{x}|, |\vec{y}|) = (L, N)$. Here, we assume $N > L$. When $[\vec{x}] = ([x_1], [x_2], \dots, [x_p])$, $[\vec{x}] \gg 1$ means $([0], [x_1], \dots, [x]_{p-1})$. In our protocol, we use two subprotocols as follows:

- All-AND takes a list $[\vec{t}]$ (with p Boolean shares) as input and outputs $[t_1 \wedge \dots \wedge t_p]^B$. We can compute this function with $[p]$ communication rounds (by appropriate parallelization) and $O(NL)$ -bit data transfer. The total data transfer size for this All-AND is $\{-\frac{1}{3}L^3 + \frac{1}{2}(N+1)L^2 - \frac{1}{6}(3N+1)L\}$ -bit since we execute All-AND at most $(N-1)$ times.
- All-OR takes a list $[\vec{u}]$ (with p Boolean shares) as input and outputs $[u_1 \vee \dots \vee u_p]^B$. We can compute this function with $[p]$ communication rounds (by appropriate parallelization) and $O(N)$ -bit data transfer. The total data transfer size for this All-OR is $\{-\frac{1}{2}L^2 + \frac{1}{2}(2N-1)L\}$ -bit since we execute All-OR for L times.

Our protocol is as in Protocol A1.

In the following, we explain the details of our baseline longest common prefix search protocol using an example that strings $\vec{x} = \text{“TGA”}$ and $\vec{y} = \text{“ATTGC”}$. In this example, $w = 2$ since there exists “TG” in \vec{y} , but “TGA” does not. First, we check whether $w = 1$ or not. To achieve this functionality, we check whether the first character of \vec{x} (i.e., “T” in the example) exists in \vec{y} using Equality for N times. If there exists at least one True in this calculation result, it means $w = 1$ (or larger), and we can achieve this functionality using All-OR. Then, we check whether $w = 2$ or not; that is, we check whether there exists “TG” in the subsequence of \vec{y} . Here, we only need to consider (“AT”, “TT”, “TG”, “GC”) as subsequences in our example. We search for the perfect matching using Equality and

■ **Protocol A1** Baseline Secure LPM.

Functionality: Compute the length of the longest common prefix w

Input: Strings $\llbracket \vec{x} \rrbracket$ and $\llbracket \vec{y} \rrbracket$, where $(|\vec{x}|, |\vec{y}|) = (L, N)$

Output: $\llbracket w \rrbracket$

```

1: for  $i = 1, \dots, L$  do
2:   for  $j = 1, \dots, N$  do
3:      $\llbracket \vec{s}_{i,j} \rrbracket^B = \text{Equality}(\llbracket \vec{x}_i \rrbracket, \llbracket \vec{y}_j \rrbracket)$ 
4:   end for
5: end for
6: for  $i = 1, \dots, L$  do
7:    $P_I$  ( $I \in \{0, 1\}$ ) locally generates an empty list  $\llbracket \vec{u} \rrbracket_I$ .
8:   for  $j = 1, \dots, N - i + 1$  do
9:     if  $i = 1$  then
10:       $P_I$  locally adds  $\llbracket \vec{s}_{1,j} \rrbracket_I^B$  to  $\llbracket \vec{u} \rrbracket_I$ .
11:     else
12:       $P_I$  locally generates an empty list  $\llbracket \vec{t} \rrbracket_I$ .
13:      for  $k = 1, \dots, i$  do
14:         $P_I$  locally add  $\llbracket \vec{s}_{k,k+j-1} \rrbracket_I^B$  to  $\llbracket \vec{t} \rrbracket_I$ .
15:      end for
16:       $P_I$  adds All-AND( $\llbracket \vec{t} \rrbracket$ ) to  $\llbracket \vec{u} \rrbracket_I$ .
17:     end if
18:   end for
19:    $\llbracket \vec{v}_{L-i+1} \rrbracket^B = \text{All-OR}(\llbracket \vec{u} \rrbracket)$ 
20: end for
21:  $\llbracket \vec{v} \rrbracket^B = \llbracket \vec{v} \rrbracket^B \oplus (\llbracket \vec{v} \rrbracket^B \ggg 1)$ 
22:  $\llbracket \vec{v} \rrbracket = \text{B2A}(\llbracket \vec{v} \rrbracket^B)$ 
23:  $\llbracket w \rrbracket = \sum_{\ell=1}^L \llbracket \vec{v}_\ell \rrbracket \cdot \ell$ 
24: return  $\llbracket w \rrbracket$ 

```

All-AND in our protocol. The condition $w = 2$ (or larger) holds if there is at least one perfect matching, and we can achieve this functionality using All-OR. We can compute the cases of $w \leq 3$ using almost the same strategy. After that, we extract the number of the longest common prefix. In the above procedure, we can obtain (False, True, True) for $(w = 3, w = 2, w = 1)$, respectively. We can extract the leftmost True using 1-bit right-shift and XOR, which is a common technique for constructing secure protocols. Finally, we can obtain a final output $\llbracket w \rrbracket$ using B2A and the inner product (with constant numbers). Note that we can optimize Equality by replacing simple OR. This is because all characters in \vec{x} and \vec{y} are {"A", "T", "G", "C"}, and we can represent them using 2-bit arithmetic sharing. With an appropriate parallelization, we can execute Protocol A1 with $O(\log \lceil L \rceil + \log \lceil N \rceil)$ communication rounds.

D Proof of Theorem 2

Proof. Correctness and security of ss-ROT protocol are proved as follows.

Proof of correctness. We assume the following equation.

$$p_i = (V^{(i)}[p_0] + r^{i-1}) \bmod N \quad (8)$$

In Step 1, for $j = 0$, the protocol computes p_1 by reconstructing $R^0[p_0]$. From the definition of $R^j[i]$ in Eq. 4,

$$p_1 = R^0[p_0] = (V^{(1)}[p_0] + r^0) \bmod N. \quad (9)$$

For $j = k$, the protocol computes p_{k+1} by reconstructing $R^k[p_k]$. From the definition of $R^j[i]$ in Eq. 4 and the assumption of Eq. 8,

$$\begin{aligned} p_{k+1} = R^k[p_k] &= (V[(p_k - r^{k-1}) \bmod N] + r^k) \bmod N \\ &= (V[V^{(k)}[p_0]] + r^k) \bmod N \\ &= (V^{(k+1)}[p_0] + r^k) \bmod N. \end{aligned} \quad (10)$$

Eq. 8 holds for $i = 1$ by Eq. 9. It also holds for $i = k + 1$ under the assumption that Eq. 8 holds for $i = k$. Therefore by induction, Eq. 8 holds for $i = 1, \dots, \ell - 1$.

In Step 2, P_0 and P_1 output $\llbracket R^{\ell-1}[p_{\ell-1}] \rrbracket$. Since Eq. 8 holds for $i = \ell - 1$,

$$R^{\ell-1}[p_{\ell-1}] = (V[(p_{\ell-1} - r^{\ell-2}) \bmod N]) \bmod N$$

is transformed into $(V^{(\ell)}[p_0]) \bmod N$ by plugging in $p_{\ell-1} = V^{(\ell-1)}[p_0] + r^{\ell-2}$. Therefore the final output of ss-ROT becomes $\llbracket (V^{(\ell)}[p_0]) \bmod N \rrbracket$. The above argument completes the proof of correctness of Theorem 2.

Proof of security. Due to space limits, we only show a sketch of the proof. In the DB preparation phase of ss-ROT, \mathcal{B} does not disclose any private values, and P_0 and P_1 receive the shares. In the Search phase, all the messages exchanged between P_0 and P_1 are shares except for the result of Reconst in Step 1. In the j -th step of the loop in Step 1, $p_{j+1} = R^j[p_j] = (V^{(j+1)}[p_0] + r^j) \bmod N$ is reconstructed. Since the reconstructed value is randomized by r^j , no information is leaked. Note that for each vector R^j , all the elements $R^j[0], \dots, R^j[N-1]$ are randomized by the same value r^j , but only one of them is reconstructed, and different random numbers $r^0, \dots, r^{\ell-1}$ are used for $R^0, \dots, R^{\ell-1}$. In Step 2, P_0 and P_1 output a result, and no information other than the result is leaked. ◀

E Proof of Theorem 3

Proof. Correctness and security of Protocol 2 are proved as follows.

Proof of correctness. The lookup table V simply stores all possible outputs of LF. Therefore, backward search (Eq. 1) is equivalent to Eq. 5. For the case of querying w , $V_{w[k-1]}[\dots V_{w[0]}[p_0] \dots]$ becomes lower bound f (for $p_0 = 0$) or upper bound g (for $p_0 = N$) of the interval that corresponds to the prefix match of length k . In Line 5 of Protocol 2, $\llbracket R_{A,f}^k[f_k] \times q_A[k] + R_{C,f}^k[f_k] \times q_C[k] + R_{G,f}^k[f_k] \times q_G[k] + R_{T,f}^k[f_k] \times q_T[k] \rrbracket$ is computed. Since $q_{w[j]}[j] = 1$ and $q_c[j] = 0$ ($c \neq w[j]$), it is equivalent to $\llbracket R_{w[k],f}^k[f_k] \rrbracket$. Line 6 computes $\llbracket R_{w[k],g}^j[g_k] \rrbracket$ in the same manner. Each vector $R_{c,f}^j$ in Eq. 6 is generated in the same manner as R^j in Eq. 4. Since Eq. 6 uses the common random values r_f^j and r_f^{j-1} for $R_{A,f}^j, R_{C,f}^j, R_{G,f}^j, R_{T,f}^j$, we can recursively reference V_c ($c \in \{A, C, G, T\}$), which is obvious from the correctness of ss-ROT. Therefore, the recursion by Line 5 and Line 7 can compute $(V_{w[k-1]}[\dots V_{w[0]}[f_0] \dots] + r_f^{k-1}) \bmod N'$, and the recursion by Line 6 and Line 8 can also compute $(V_{w[k-1]}[\dots V_{w[0]}[g_0] \dots] + r_g^{k-1}) \bmod N'$.

The longest match is found when the interval width becomes 0. Since $f_k = (V_{w[k-1]}[\dots V_{w[0]}[f_0] \dots] + r_f^{k-1}) \bmod N'$ and $g_k = (V_{w[k-1]}[\dots V_{w[0]}[g_0] \dots] + r_g^{k-1}) \bmod N'$ are randomized, Line 11 computes $f_k - g_k - (r'[k-1] = r_f^{k-1} - r_g^{k-1})$ to obtain the correct interval width. Line 11 also computes the equality of 0 and the interval width for each round. By reconstructing all the results in Lines 15–17, \mathcal{A} knows the round, in which the interval width becomes 0; i.e., he/she knows LPM. The above argument completes the proof of correctness of Theorem 3.

Proof of security. Due to space limitation, we only show a sketch of the proof. For Lines 1–2 of Protocol 2, \mathcal{A} and \mathcal{B} do not disclose any private values, and P_0 and P_1 receive the shares. For Lines 3–13, it is guaranteed by the subprotocols ADD, MULT, and Equality that all the messages exchanged between P_0 and P_1 are shares except for the output of Reconst in Lines 7–8. (See Section 2.1 for details of the subprotocols.) In Lines 7–8, reconstructed values are $R_{w[j],f}^k[f_j]$ and $R_{w[j],g}^k[g_j]$. Since the values are $(V_{w[j]}[f_j] + r_f^j) \bmod N'$ and $(V_{w[j]}[g_j] + r_g^j) \bmod N'$ according to Eq. 6, it is obvious that V is randomized for all rounds $j = 0, \dots, \ell - 1$, and no information is leaked. For Lines 14–17, only the output of Equality at Line 11 is reconstructed. The reconstructed values are either 1 or 0 according to Equality, and no information other than the result is leaked. ◀

F Proof of theorem 4

Proof. Correctness and security of Protocol 3 are proved as follows.

Proof of correctness. V , R , r' and q are generated by the same algorithm used in Protocol 2. Therefore, Line 6 is equivalent to a backward search, and $e1$ is the result of the equality check of 0 and the width of the obtained interval in Line 7. The lookup tables V_{lcp} , V_{psv} , and V_{nsv} store all the outputs of LCP, PSV and NSV, and W_l , W_p , and W_n are generated based on V_{lcp} , V_{psv} , and V_{nsv} , respectively. Since $W_{l,f}^j$ and $W_{l,g}^j$ are circular permutations of V_{lcp} by the same random values r_f^{j-1} and r_g^{j-1} that are used for generating $R_{c,f}$ and $R_{c,g}$ ($c \in \Sigma$) respectively, Line 8 can compute $\text{LCP}[g_j] \leq \text{LCP}[f_j]$ and $e2$ holds the result. By using Choose and $e2$, either $[W_{p,f}^j[f_j], W_{n,f}^j[f_j]]$ or $[W_{p,g}^j[g_j], W_{n,g}^j[g_j]]$ is selected. $W_{p,f}^j$ and $W_{p,g}^j$ are permuted by r_f^{j-1} and r_g^{j-1} , but are randomized by the identical random value r_f^j . Similarly, $W_{n,f}^j$ and $W_{n,g}^j$ are permuted by r_f^{j-1} and r_g^{j-1} , but are randomized by r_g^j . Since $W_{p,f}[f_j]$ and $W_{n,g}[g_j]$ are generated in the same manner as $R_{c,f}$ and $R_{c,g}$, it is obvious that the reference by them is correct. The reference by $W_{n,f}[f_j]$ is transformed into

$$\begin{aligned} X_g^{j+1}[W_{n,f}^j[f_j]] &= V_x[W_{n,f}^j[f_j] - r_g^j] + r_g^{j+1} \\ &= V_x[V_{nsv}[f_j - r_f^{j-1}] + r_f^j - r_g^j] + r_g^{j+1} \\ &= V_x[V_{nsv}[f_j - r_f^{j-1}]] + r_g^{j+1} \end{aligned} \quad (11)$$

and the reference by $W_{p,f}[g_j]$ is transformed into

$$\begin{aligned} X_f^{j+1}[W_{p,f}^j[g_j]] &= V_x[W_{p,f}^j[g_j] - r_f^j] + r_f^{j+1} \\ &= V_x[V_{psv}[g_j - r_g^{j-1}] + r_g^j - r_f^j] + r_f^{j+1} \\ &= V_x[V_{psv}[g_j - r_g^{j-1}]] + r_f^{j+1} \end{aligned} \quad (12)$$

where X^{j+1} is any one of R_c^{j+1} , W_p^{j+1} and W_n^{j+1} , and V_x is the corresponding lookup table; i.e., either one of V_c , V_{psv} and V_{nsv} . Note that V_x could be a different table for each $j + 1$, but we abuse the same notation for simplicity of notation. Since f_j and g_j are described in the form of $V_x^{(j)}[p_0] + r_f^{j-1}$ and $V_x^{(j)}[p'_0] + r_g^{j-1}$ based on Eq. 8, Eq. 11 and Eq. 12 are transformed into $V_x^{(j+2)}[p_0] + r_g^{j+1}$ and $V_x^{(j+2)}[p'_0] + r_f^{j+1}$, which also satisfy the recursion form of Eq. 8. Thus, the intervals $[W_{p,f}^j[f_j], W_{n,f}^j[f_j]]$ and $[W_{p,g}^j[g_j], W_{n,g}^j[g_j]]$ are correct intervals and Line 9 is equivalent to computing Eq. 2.

Lines 14–18, u remains the same if $e1 = 0$ and is permuted such that $u[i] = u[(i-1) \bmod \ell]$ otherwise. Therefore Lines 19–21 can choose the letter to be searched with. The match length and the start position are obtained based on $e1$ in Lines 22–23, and the longest value and the corresponding position are selected in Lines 24–25. The shares of the length and start position of LMEM are sent to \mathcal{A} , and \mathcal{A} reconstructs them. Then, Protocol 3 outputs them. The above argument completes the proof of correctness of Theorem 4.

Proof of security. Due to space limits, we only show a sketch of the proof. For Lines 1–2 of Protocol 3, \mathcal{A} and \mathcal{B} do not disclose any private values, and P_0 and P_1 receive the shares. For Lines 3–27, it is guaranteed by the subprotocols ADD, MULT, Equality, and Choose that all the messages exchanged between P_0 and P_1 are shares except for the output of Reconst in Line 12. (See Section 2.1 for details of the subprotocols.) In Line 12, the reconstructed values are $f_{i+1} = V_x^{(j+1)}[p_0] + r_f^j$ and $g_{j+1} = V_x^{(j+1)}[p_0] + r_g^j$, according to Eq. 8, Eq. 11, and Eq. 12. Since f_{j+1} and g_{j+1} are randomized by r_f^j and r_g^j , respectively, for all rounds $j = 0, \dots, 2\ell - 1$, no information is leaked. In Line 28, \mathcal{A} reconstructs only the search result (the length and start position of LMEM). ◀