

曖昧検索機能を持つ動的な検索可能暗号方式

大塚 元気¹ 多田 充^{2,a)}

概要: 2012年, Kamaraら [2] は検索可能暗号方式 (SSE) が最低限実用的なものになるためには, (1) 検索にかかる計算量が準線形オーダーであること, (2) 適応的な選択キーワード攻撃に対し安全であること, (3) 大きすぎないインデックスサイズ及びインデックスに対するファイルの動的な追加・削除が可能であること, が必要であると述べ, これらの要件をすべて満たす SSE を提案した。しかし, この方式は検索の機能として, 完全一致によるものしか扱うことができない。SSE としてより実用的なものにするためには, 曖昧検索機能やレコメンド等の検索の機能を実現し, より柔軟な検索機能を実現すべきである。本論文では, 上記の要件を全て満たし, その上で, フォーマットの違いやスペリングミス等の誤字に対応する機能や, ユーザの入力から関連するキーワードを予測するレコメンド機能を実現する方法を述べる。

キーワード: 検索可能暗号方式 (SSE), 曖昧検索

Dynamic searchable symmetric encryption with fuzzy reference function

GENKI OTSUKA¹ MITSURU TADA^{2,a)}

Abstract: In 2012, Kamara et al.[2] told that searchable symmetric encryption (SSE) systems should provide the following functions for its being minimally practical: (1) The system requires, at most, quasi-linear computational complexity for searching; (2) The system is secure against adaptive chosen-keyword attacks; (3) The system enables dynamic file-addition and file-deletion for each index of not so large size. Then they presented an SSE with all the properties given above. However, that system can deal with perfect-matching search. It is better for SSE to have the functions for flexible search such as fuzzy-matching search and showing recommendations. In this paper, we show an SSE which realizes the functions for making up for misspelling and different formats of data, and that for recommendation derived from words given by the user.

Keywords: Searchable symmetric encryption (SSE), Fuzzy reference

1. はじめに

クラウドストレージは大変便利なサービスである。ただ, そのセキュリティは完璧とは言えない。例えば, 格納したデータを外部に漏らさないとも限らない。またそもそも,

クラウドストレージのサービス提供者がそのデータを見るかも知れない。そのため, プライベートなデータ, その他機密性が要求されるデータ, 人に見せたくないデータをクラウドストレージに格納する場合, 当該データを暗号化して格納する必要がある。しかし暗号化するという行為はそのデータに対する基本的な処理 — 例えば文書の検索 — ですら不可能にしてしまう。そのため, クラウドストレージの利便性を保ったまま, プライベートなデータに対する検索, つまり, 暗号化されたデータに対する検索を行いたいというのは自然な要求でもある。検索可能暗号とは, データを暗号化したまま検索を行う技術を指す。検索可能暗号

¹ 千葉大学大学院 理学研究科
Graduate School of Science, Chiba University, Inage, Chiba
263-8522, Japan

² 千葉大学 統合情報センター
Institute of Media and Information Technologies, Chiba Uni-
versity, Inage, Chiba 263-8522, Japan

a) m.tada@faculty.chiba-u.jp

の中でも、ユーザが所持し、クラウドストレージへ格納したいファイル群 $\mathbb{F} = (f_1, \dots, f_n)$ の他に、インデックスと呼ばれるデータ構造を \mathbb{F} から作成し、暗号の初期化時に、暗号化したファイル群 $\mathbb{C} = (c_1, \dots, c_n)$ と共にデータベースに格納し、検索時にはこのインデックスを利用し検索を行うタイプの検索可能暗号を「インデックススペースの検索可能暗号」と呼ぶ。本論文では、インデックススペースの検索可能暗号に対する考察をし、単に検索可能暗号と呼んだ場合はインデックススペースのものを指すことにする。著者らが知る限り、検索可能暗号に対する明示的な発表は2000年の Song ら [4] のものであり、以降様々な成果が報告されている。2012年、Kamara ら [2] は、検索可能暗号が実用的であるためには最低でも

- (C1) 検索に要する時間が準線形であること
- (C2) 適応的キーワード選択攻撃に対して安全であること
- (C3) インデックスのサイズが現実的であること
- (C4) インデックスを動的に更新、つまりファイルの追加と削除が効率的に可能であること

のすべてが満たされるべきであると主張し、同文献において、そのすべてを満たしたインデックススペースの検索可能暗号を提案した。しかし、以上のすべてを満たすことは検索可能暗号が実用的なものとなるための最低限な要件に過ぎず、実際、Kamara らの方法では完全一致による検索のみしか行うことができない。そのため誤字や表記の揺れ、フォーマットの違いに対応せず、また我々がよく Google 等の検索エンジンで検索をするときのような、レコメンド機能^{*1}も存在しない。さらに、インデックス作成時に使用するワードをすべて定義しておかなければならず、新規単語を追加したい場合はデータベースを再度、一から構築する必要がある。本論文では Kamara らの手法 [2] を拡張し、以下のようにして、これらの問題点を解決した。これは、クラウドストレージサービスを検索可能暗号を用いて展開するにあたり、Kamara らの手法では足りない機能面の拡張に焦点を当て、検索可能暗号としての必要最低限の項目をクリアしたうえで、以下の問題を解決したものである。

- (1) 誤字耐性：N-Gram[3]の技術を用い、検索時に、少しの誤字^{*2}が存在しても正しく検索できるようにした。
- (2) 新規単語の追加：N-Gramの技術を用いることにより、新規単語を出現させないようにインデックスを作成できるようにした。
- (3) レコメンド機能：レコメンド用のテーブルを新たに作成することによりレコメンド機能を実現した。特に、同じワードを表す複数の表現やフォーマットの異なる表現をレコメンドできるようにした。

本論文は以下のように構成される。第2章では提案手法や既存技術を説明するために必要な記号や概念の定義を行う。第3章では既存の方式を示すが、スペースの都合上、説明は割愛する。第4章では、課題解決のために必要な戦略を述べ、実際にその課題を解決するための構成を具体的に記述し、安全性に関する説明も行う。第5章では提案手法の問題点や、曖昧検索の具体例を述べ、第6章で本論文をまとめる。

2. 準備

ここでは、本論文で使用する表記法、検索可能暗号方式(SSE)、および、インデックスへの動的な変更を可能にした Dynamic SSE (DSSE) の定義を述べる。

2.1 表記法

列 v に対して、 v_i または $v[i]$ によりその i 番目の要素、 $\#v$ によりその要素数を表す。同様に A を配列とすると、 $A[i]$ でその i 番目の要素、 $\#A$ でセルの数を表す。 L がリストのとき、 $\#L$ でそのノード数を表す。 T を連想配列としたとき、キーと値の組 (s, v) が T に含まれていれば、 $T[s]$ は v を表し、 $\#T$ は格納されている組の数を表す。キーが s となる要素が存在するとき $s \in T$ と表す。

$W \subset \mathfrak{P}(\{0, 1\}^*)$ で可能な語全体の(有限)集合を表し、これを辞書とも呼ぶ。ただし、セキュリティパラメータ k に対して $\#W = \text{poly}(k)$ とする^{*3}。語 $w \in W$ に対して $|w|$ で w のビット長を表す。ただし、 $\forall w \in W, |w| = \text{poly}(k)$ とする。 $f = (w_1, \dots, w_m) \in W^m$ をファイルとしたとき、 $\#f$ は f に含まれる語数、すなわち m とし、 $|f|$ はファイルのビット長、つまり $\sum_{i=1}^m |w_i|$ とする。また、 \bar{f} は f の中から重複を除いたファイルとする。ただし、重複は後ろから除くとする。

ユーザのデータは $n (= \text{poly}(k))$ 個のファイルの集合 $\mathbb{F}(\subset 2^W) = (f_1, \dots, f_n)$ と見なせる。各ファイル $f_i = (w_{i,1}, \dots, w_{i,m_i})$ は $m_i = \text{poly}(k)$ 個の語を含み、一意的な識別子 (ID) をもち、これを $\text{id}(f_i)$ とする。キーワード $w \in W$ に対して、 \mathbb{F} 中の w を含むファイルのリストを、 \mathbb{F}_w で表す。各 f_i を定められた暗号方式 Enc および鍵 K で暗号化したものを c_i とし、 $\mathbb{C} = (c_1, \dots, c_n)$ としたとき、 w に対する \mathbb{F}_w の各要素を暗号化したものを \mathbb{C}_w と表記する。つまり、 $\mathbb{F}_w = (f_{i_1}, \dots, f_{i_\ell})$ とすると $\mathbb{C}_w = (c_{i_1}, \dots, c_{i_\ell})$ となる。

2.2 検索可能暗号方式 (SSE)

辞書 W 上のインデックススペース SSE 方式とは、5つの

^{*1} 文章入力時に、入力中の単語・文に対し、次に入力する可能性の高い単語・文をユーザに知らせる機能。

^{*2} どれだけの曖昧さを許すかについてはアプリケーションごとにパラメータとして定める。

^{*3} 本論文では、 x が k の多項式に満たないオーダー — 例えば $x = o(k)$ — であったとしても $x = \text{poly}(k)$ のように表記する。つまり $x = \text{poly}(k)$ は「 x は高々 (k の) 多項式オーダー」の意味であるとする。

多項式時間アルゴリズムの組：

$SSE = (\text{Gen}, \text{Enc}, \text{Trapdoor}, \text{Search}, \text{Dec})$

である。各アルゴリズムは以下のように定義される。

Gen : $1^k \mapsto K$: セキュリティパラメータ k に対する秘密鍵 K を出力する確率的アルゴリズムである。

Enc : $(K, \mathbb{F}) \mapsto (\gamma, \mathbb{C})$: 鍵 K および文書 (ファイル) のコレクション \mathbb{F} を入力として、暗号化した文書のコレクションおよび検索用のインデックス (γ) を出力するアルゴリズムである。

ただし、 $\text{Enc}(K, \mathbb{F})$ を $\text{Enc}_K(\mathbb{F})$ を表記することもある。

Trapdoor : $(K, w) \mapsto t$: 秘密鍵 K およびキーワード w を入力として、トラップドア t を出力する決定的アルゴリズムである。

前項と同様に $\text{Trapdoor}_K(w)$ と表記することもある。

Search : $(\gamma, t) \mapsto X$: インデックス γ およびトラップドア t を入力として、辞書順に並んだ文書の ID からなるリスト X を出力する決定的アルゴリズムである。

Dec : $(K, c_i) \mapsto f_i$: 秘密鍵 K および暗号化された文書 c_i を入力とし、復号された文書 f_i を出力する決定的アルゴリズムである。

$\text{Dec}(K, c_i)$ を $\text{Dec}_K(c_i)$ と表記する場合もある。

ユーザは自身の持つデータ (\mathbb{F}) をクラウドストレージサーバ (以下「サーバ」と称する) へ SSE を利用し格納したいと考えている。このとき、SSE の各アルゴリズムは以下のように実行される。

まず、ユーザは秘密鍵の生成 (Gen) を行う。この鍵 (K) はインデックスの作成時と検索時、そして、データの暗号化及び復号に用いる。ユーザは鍵を用い、サーバへ格納したいデータからインデックス (γ) と呼ばれる検索用の構造を生成 (Enc) する。同時に、鍵を用いファイルそのものも暗号化し、インデックスと一緒にサーバへ格納する。ユーザがファイルの検索を行う場合は、鍵および検索ワード w から、トラップドア t と呼ばれるクエリを作成 (Trapdoor) し、サーバへ送信する。クエリを受け取ったサーバはインデックスとクエリから検索処理 (Search) を行い、検索結果として、暗号化されたファイル群 (\mathbb{C}_w) をユーザへ返す。最後に、ユーザが検索結果のファイル群を復号 (Dec) することにより、検索は完了する。

2.3 Dynamic 検索可能暗号方式 (DSSE)

動的インデックスベース SSE 方式 (DSSE) とは、以下の9つの多項式時間アルゴリズムの組

$DSSE = (\text{Gen}, \text{Enc}, \text{SrchToken}, \text{AddToken}, \text{DelToken}, \text{Search}, \text{Add}, \text{Del}, \text{Dec})$

であり、各アルゴリズムは以下のように定義される。ただし、Gen, Enc, Dec については前節と同様である。

SrchToken : $(K, w) \mapsto \tau_s$: 秘密鍵 K およびキーワード w を入力として、検索トークン τ_s を出力する確率的

アルゴリズムである。

AddToken : $(K, f) \mapsto (\tau_a, c_f)$: 秘密鍵 K および文書 (ファイル) f を入力として、追加トークン τ_a および追加するファイルを暗号化したもの c_f を出力する確率的アルゴリズムである。

DelToken : $(K, f) \mapsto \tau_d$: 秘密鍵 K および文書 (ファイル) f を入力として、削除トークン τ_d を出力する確率的アルゴリズムである。

Search : $(\gamma, \mathbb{C}, \tau_s) \mapsto I_w$: インデックス γ , 暗号文の列 \mathbb{C} および検索トークン τ_s を入力として、ファイル ID からなるリスト I_w を出力する決定的アルゴリズムである。

Add : $(\gamma, \mathbb{C}, \tau_a, c) \mapsto (\gamma', \mathbb{C}')$: インデックス γ , 暗号文の列 \mathbb{C} , 追加トークン τ_a および追加する暗号文 c を入力として、新しいインデックス γ' および新しい暗号文の列 \mathbb{C}' を出力するアルゴリズムである。

Del : $(\gamma, \mathbb{C}, \tau_d) \mapsto (\gamma', \mathbb{C}')$: インデックス γ , 暗号文の列 \mathbb{C} および削除トークン τ_d を入力として、新しいインデックス γ' および新しい暗号文の列 \mathbb{C}' を出力するアルゴリズムである。

ユーザは秘密鍵の生成 (Gen) を行う。生成された鍵 K はインデックスの作成時と検索時、データの暗号化及び復号、およびファイルの追加及び削除時に用いる。ユーザは鍵 K を用い、サーバへ格納したいデータからインデックスと呼ばれる検索用の構造を生成 (Enc) する。同時に、鍵を用いファイルそのものも暗号化し、インデックスと共にサーバへ格納する。ユーザがファイルの検索を行う際は、鍵および検索ワードから、検索トークンと称されるクエリを作成 (SrchToken) しサーバへ送信する。クエリを受け取ったサーバはインデックスとクエリから検索処理 (Search) を行い、検索結果として、暗号化されたファイル群をユーザへ返す。最後に、ユーザは検索結果のファイル群を復号 (Dec) することにより、検索を完了する。また、ユーザがファイルの追加、または削除を行う場合は、鍵および追加するファイルから追加トークン、または削除トークンを作成 (AddToken, DelToken) し、サーバへ送信する。トークンを受け取ったサーバはインデックスとトークンから追加 (Add) または削除 (Del) 処理を行う。

3. 既存方式 [2]

本論文では既存技術として Kamara ら [2] が考案した DSSE を取り上げ、本文中でしばしば参照し、表記についても踏襲しているが、その詳細については、スペースの都合上、割愛する。

4. 提案方式

本章では、Kamara らの方法 [2] に対する下記の課題：

- 新規単語の追加問題

- 誤字や表記の揺れの問題

の対応策として、第 4.1 節から第 4.4 節にかけて、

- N -Gram[3] を用いた新規単語を出現させない方法
- N -Gram を用いた誤字や表記の揺れを検索する方法
- レコメンドテーブルを用いて単純に単語をレコメンドする方法
- フォーマット変換をレコメンドテーブルを用いて実現する方法

を述べ、第 4.5 節において提案方式の詳細を述べ、第 4.6 節においてその安全性について述べる。

4.1 新規単語の追加問題に対する対応

既存方式では、テーブルに追加するワード全体の集合 W は予め決まっておき、使用状況に応じ新しい単語を追加する場合は、テーブルを一から作りなおす必要がある。しかし、 N -Gram の技術を用いることで、この問題を以下のように解決できる。ここで、ワード w の N -Gram とは、 w 中の連続する長さ N の部分文字列である。例えば、 w を「検索可能暗号」とすれば、その 2-Gram は「検索」「索可」「可能」「能暗」「暗号」の 5 つである。

N -Gram を利用し、新規単語の追加問題に対応する戦略は以下のようなものである。

- (1) 使用する文字全体の集合を C (ただし $|C| = \text{poly}(k)$) とする。
- (2) 登録する文字全体の集合 W を $W := C^N$ と定義する。
- (3) 文書 f をテーブルへ登録する場合、 f をひとつの文字列とみなし、その N -Gram を利用する。

つまり、テーブルへ登録されるワードは、文字列の N -Gram として現れる文字列すべて、即ち、長さが N の文字列全体である。これらをすべて登録する。インデックス作成時に、文書 f の情報をテーブルへ登録する方法は、 f を 1 つの文字列とみなし、現れた N -Gram に対応するリストへ f (の ID) を格納することになる。

4.2 誤字や表記の揺れへの対応

既存方式では、完全一致によるものしか扱えず、文書中や検索文字列に誤字があった場合、検索が困難になってしまう。そこで、 N -Gram を用い登録されていることを仮定し、以下のような変更を加えることによりこの問題を解決できる。なお、表記の揺れへの対応については、次節で述べるレコメンド機能を用い対応する。

インデックスへの追加処理：

テーブル A_s のノードに暗号化したポジションの情報の列を加える。これにより、ワードが文書のどの位置に出現したかが分かるようになる。なお、ワードの出現回数を漏洩させないために、ポジションの列の長さが固定長になるようダミーデータを加える。このとき、文書中に出てきた各 N -Gram のポジション情報を記憶できる個数の上限を

定める。

検索処理：

- (1) 検索語 w を N -Gram のワードに区切る。ここで $w \mapsto (g_1, \dots, g_m)$ とし、さらに w の長さを ℓ とする。
- (2) g_1, \dots, g_m を従来方式で検索する。このとき、暗号化されたポジション情報を得るので、ユーザは復号する。
- (3) g_1, \dots, g_m が現れた文書およびポジションから、完全一致や曖昧検索を行う。ここで、曖昧検索で合致したと見なす基準はアプリケーションごとに自由に定めることができる。例えば、 w の $\lceil m/2 \rceil$ 個以上の異なる N -Gram が 2ℓ の長さ以下の文字列に出現した文書を検索結果に含める、等である。

4.3 レコメンド機能の追加

既存方式にはレコメンド機能は存在しない。ここで、より検索をスムーズに行えるように単語間のレコメンド機能を以下のように追加する。

下準備：

- (1) 暗号化されていないレコメンドテーブルを用意する。
暗号化されていないレコメンドテーブルとは、ワードとワードのリストの対応表のことである。
- (2) レコメンドテーブルの暗号化
 - T_s, A_s と同様の構造をもったテーブル T_r, A_r を用意する。
 - ワード w でレコメンドされるワード群を w_1, \dots, w_n としたとき、 A_r にはワード w_1, \dots, w_n を暗号化したものをランダムな場所にリストとして格納し、 T_r にはワード w のユーザ鍵でのハッシュ値に対応する場所に上記リストの先頭アドレスを格納する。

レコメンドテーブル検索時：

- (1) 検索キーワードを w とする。
- (2) w のユーザ鍵でのハッシュ値をサーバへ送信する。
- (3) サーバはリストを復元、暗号化されたワードのリストをユーザへ返す。
- (4) ユーザはリストの各ノードを復号し、レコメンド情報を得る。

4.4 フォーマット変換レコメンド機能の追加

単語間の関連だけではフォーマットの違いに対応できないため、より柔軟なレコメンド機能を実現できるよう、簡易的なフォーマット変換も行えるように拡張する。

下準備：

- (1) フォーマット変換で用いるメタ文字の集合を M とする。ただし、 M 内の各メタ文字は
 - 単一の文字を表す
 - メタ文字以外からメタ文字への対応が一意になるという条件を満たすとする。
- (2) 暗号化されていないレコメンドテーブルを用意する。

このとき、検索ワードとレコメンドリストに格納するワードにメタ文字を用いてもよい。

(3) レコメンドテーブルの暗号化

- T_s, A_s と同様の構造をもったテーブル T_r, A_r を用意する。
- ワード w でレコメンドされるワード群を w_1, \dots, w_n としたとき、 A_r にはワード w_1, \dots, w_n を暗号化したものをランダムな場所にリストとして格納し、 T_r にはワード w のユーザ鍵でのハッシュ値に対応する場所に上記リストの先頭アドレスを格納する。

レコメンドテーブル検索時：

- (1) 検索キーワードを w とする。
- (2) w を可能な限りメタ文字で置き換えた後、 w のユーザ鍵でのハッシュ値をサーバへ送信する。
- (3) サーバはリストを復元、暗号化されたワードのリストをユーザへ返す。
- (4) ユーザはリストを復号し、レコメンド情報を得る。
- (5) レコメンド情報にメタ文字が含まれていた場合は置き換えたうえでユーザにレコメンド情報として伝える。

4.5 提案方式

ここでは提案方式の構成を述べる。なお、第 2.3 節で示した関数 (アルゴリズム) の他、曖昧検索処理を行う関数 $FuzSearch$ 、レコメンドテーブルの検索クエリを作成する関数 $RecSrchToken$ 、および、レコメンドテーブルの検索処理関数 $RecSearch$ を追加している。提案方式において、 $SKE1 = (Gen1, Enc1, Dec1)$ および $SKE2 = (Gen2, Enc2, Dec2)$ を適切な対称鍵暗号方式とする。また、 $F : \{0,1\}^k \times \{0,1\}^* \rightarrow \{0,1\}^k$ 、 $G : \{0,1\}^k \times \{0,1\}^* \rightarrow \{0,1\}^*$ 、 $P : \{0,1\}^k \times \{0,1\}^* \rightarrow \{0,1\}^k$ 、および $Q : \{0,1\}^k \times \{0,1\}^* \rightarrow \{0,1\}^*$ を擬似ランダム関数とし、 $H_1 : \{0,1\}^* \rightarrow \{0,1\}^*$ および $H_2 : \{0,1\}^* \rightarrow \{0,1\}^*$ をランダムオラクルとする。 $z \in \mathbb{N}$ をフリーリストの初期サイズとする。使用できる通常文字の有限集合を Σ とする (このとき $W \subset \Sigma^*$ である)。 s を 1 つの文書に現れることが可能な同じ N -Gram 語の個数の最大値とする。メタ文字の集合を Ω とする。ただし、 Σ の任意の文字に対応する Ω のメタ文字は多くとも 1 つとする。

我々が提案するインデックススペース DSSE 方式：

$$DSSE = (Gen, Enc, SrchToken, AddToken, DelToken, Search, Add, Del, Dec, FuzSearch, RecSrchToken, RecSearch)$$

は以下のように構成される。

Gen(1^k) :

- (i) $K_1 \xleftarrow{\$} \{0,1\}^k$, $K_2 \xleftarrow{\$} \{0,1\}^k$, $K_3 \xleftarrow{\$} \{0,1\}^k$, $K_4 \xleftarrow{\$} SKE1.Gen1(1^k)$, $K_5 \xleftarrow{\$} \{0,1\}^k$, $K_6 \xleftarrow{\$} SKE2.Gen2(1^k)$
- (ii) $K = (K_1, K_2, K_3, K_4, K_5, K_6)$ を出力する。

Enc(K, \mathbb{F}) :

(i) 準備：

- (a) A_s, A_d をサイズ $\#C/8 + z$ の配列とする。
- (b) T_s, T_d を、それぞれ、サイズ $\#\Sigma^N + 1$, $\#\mathbb{F}$ の連想配列とする。
- (c) T_r をサイズ t の連想配列とする。ここで、 t はレコメンドを機能させたい有限個の語 $w \in (\Sigma \cup \Omega)^*$ の数である。
- (d) A_r をサイズ $10t$ の配列とする。
- (e) $\mathbf{0}$ で $(\log \#A_s)$ 個の 0 の列を表すとする。
- (f) $free$ を W に現れていない文字列とする。

(ii) 検索用のリストを作成：

各 N -Gram $g \in \Sigma^N$ に対し、以下を実行する。

- (a) 配列 A_s にランダムに格納された $\#\mathbb{F}_g$ 個のノード $(N_1, \dots, N_{\#\mathbb{F}_g})$ からなるリスト L_g を構成する：

$$N_i := (\langle id_i, addr_s(N_{i+1}),$$

$$Enc2_{K_6}(\langle Q_{K_5}(pos_{i,g,1}), \dots, Q_{K_5}(pos_{i,g,s}) \rangle) \oplus H_1(K_w, r_i), r_i)$$

と定義する。ここで、 id_i は、 \mathbb{F}_g の i 番目のファイルの ID、 $pos_{i,g,j}$ は、 \mathbb{F}_g の i 番目のファイル中での g の j 回目の出現位置、 $r_i \xleftarrow{\$} \{0,1\}^k$ 、 $K_w \leftarrow P_{K_3}(w)$ 、 $addr_s(N_{\#\mathbb{F}_g+1}) = \mathbf{0}$ 、 $pos_{i,g,j}$ が存在しない場合はランダムなデータで埋める、とする。

- (b) 検索テーブル $T_s \in L_g$ の先頭ノードへのポインタを格納する：

$$T_s[F_{K_1}(g)] \leftarrow \langle addr_s(N_1), addr_d(N_1^*) \rangle \oplus G_{K_2}(g)$$

ただし、 N_1^* は N_1 の対、つまり、 A_d のノードであり、その 4 番目のエントリーポイントが A_s のノード N_1 を指しているもの、とする。

(iii) 削除用のリストを作成：

各ファイル $f \in \mathbb{F}$ に対して以下を実行する。

- (a) 削除配列 A_d にランダムに格納された $\#f$ 個の対ノード $(D_1, \dots, D_{\#f})$ からなるリスト L_f を構成する：

- 各エントリ D_i は N -Gram g と関連付けられているため、 L_g のあるノード N と関連付けられている。
- N_{+1} は、 L_w 内における N の次のノードを表す。
- N_{-1} は、 L_w 内における N の前のノードを表す。

このとき D_i を

$$D_i := (\langle addr_d(D_{i+1}), addr_d(N_{-1}^*), addr_d(N_{+1}^*), addr_s(N), addr_s(N_{-1}), addr_s(N_{+1}), F_{K_1}(g) \rangle$$

$$\oplus H_2(K_f, r'_i), r'_i)$$

と定義する。ここで、 r'_i は k ビットのランダムな列、 $K_f \leftarrow P_{K_3}(f)$ 、 $addr_d(N_{\#\mathbb{F}+1}) = \mathbf{0}$ とする。

- (b) 削除テーブル $T_d \in L_f$ の先頭ノードへのポインタを格納する：

$$T_d[F_{K_1}(f)] \leftarrow addr_d(D_1) \oplus G_{K_2}(f)$$

(iv) レコメンド用のリストを作成：
メタ文字を用いても良い、レコメンドを機能させたい
 s 個の各語 $w \in (\Sigma \cup \Omega)^*$ に対して以下の (a)~(c) を実行する。

(a) 配列 A_r にランダムに格納された w_r 個のノード
(R_1, \dots, R_{w_r}) からなるリスト L_w を構成する：

$$R_i \leftarrow (\text{Enc1}_{K_4}(w_i), \text{addr}_r(R_{i+1})) \\ \oplus H_1(K_w, r_i, r_i)$$

と定義する。ここで $r_i \stackrel{\$}{\leftarrow} \{0, 1\}^k$, $K_w \leftarrow P_{K_3}(w)$,
 $\text{addr}_r(R_{w_r+1}) = \mathbf{0}$ とする。

(b) レコメンドテーブル $T_r \leftarrow L_w$ の先頭ノードへのポ
インタを格納する：

$$T_r[F_{K_1}(w)] \leftarrow \text{addr}_r(R_1) \oplus G_{K_2}(w)$$

(c) 配列 A_r の余った要素をランダムな文字列で埋
める。

(v) 暗号化されていないフリーリスト L_{free} を作る：

(F_1, \dots, F_z) および (F'_1, \dots, F'_z) を、それぞれ、 A_s およ
び A_d からランダムに選んだ、値が格納されていない z
個のセルとする。

(a) 検索テーブル T_s にフリーリストの先頭の情報を
格納する：

$$T_s[free] \leftarrow (\text{addr}_s(F_z), 0^{\log \#A})$$

(b) 検索配列 A_s へノードの情報を格納する：

$$A_s[\text{addr}_s(F_i)] \leftarrow (0^{\log \#F}, \text{addr}_s(F_{i-1}), \text{addr}_d(F'_i)) \\ \text{とする。ただし、} \text{addr}_s(F_0) = 0^{\log \#A} \text{ とする。}$$

(vi) A_s, A_d の値が格納されていないセルをランダムな文字
列で埋める。

(vii) ファイルを暗号化する。

$1 \leq i \leq \#F$ に対して $c_i \leftarrow \text{SKE1.Enc1}_{K_4}(f_i)$ を実行
する。

(viii) (γ, C) を出力する。ここで、 $\gamma \leftarrow (A_s, T_s, A_d, T_d)$ およ
び $C \leftarrow (c_1, \dots, c_{\#F})$ とする。

SrchToken(K, w) :

(i) w を N -Gram へ分割し、それを (g_1, \dots, g_m) とする。

(ii) $1 \leq i \leq m$ に対して $\tau_{s_i} \leftarrow (F_{K_1}(g_i), G_{K_2}(g_i), P_{K_3}(g_i))$
とする。

(iii) $(\tau_{s_1}, \dots, \tau_{s_m})$ を出力する。

Search($\gamma, C, \{\tau_{s_1}, \dots, \tau_{s_m}\}$) :

サーバ側の処理は以下の通りである。

(i) 各 τ_{s_i} に対して以下を実行する。

(a) 引数をチェックする。

(i) τ_{s_i} を $(\tau_{i,1}, \tau_{i,2}, \tau_{i,3})$ のように分解する。

(ii) T_s に $\tau_{i,1}$ が存在しない場合は空リストを返す。

(b) 検索テーブルのエントリポイントを探す：

$$(\alpha_1, \alpha'_1) \leftarrow T_s[\tau_{i,1}] \oplus \tau_{i,2}$$

(c) $N_1 \leftarrow A_s[\alpha_1]$ とし、 $\tau_{i,3}$ で復号する。つまり、 N_1 を
(ν_1, r_1) として、 $\nu_1 \oplus H_1(\tau_{i,3}, r_1)$ を計算し、これ

を $\text{info}_{i,1}$ とする。実際、これは

$$(\text{id}_1, \text{addr}_s(N_2), \text{Enc2}_{K_6}(\langle Q_{K_5}(\text{pos}_{1,g,1}), \\ \dots, Q_{K_5}(\text{pos}_{i,g,s}) \rangle))$$

に等しい。

(d) $j \geq 2$ に対して、 $\alpha_{j+1} = \mathbf{0}$ となるまで、前項と同
様に N_i を復号する。

(ii) すべての $\tau_{s_i} (1 \leq i \leq m)$ の検索結果に存在する ID と
そのその出現情報を、どの τ_{s_i} の検索結果であったかを
を区別できる形式で出力する。

ユーザ側の処理は以下の通りである。

(i) ポジションの情報を復号する。

(ii) $\text{info}_{i,j}$ を用いて、すべての $\text{info}_{i,*} (1 \leq i \leq m)$ に現れ、
かつ、 $\text{info}_{i,*}$ から $\text{info}_{m,*}$ まで、 pos の情報が1つずつ
増加しているものが存在する ID をすべて出力する。

FuzSearch($\gamma, C, \{\tau_{s_1}, \dots, \tau_{s_m}\}$) :

Search($\gamma, C, \{\tau_{s_1}, \dots, \tau_{s_m}\}$) とほぼ同様であり、違いは、
ユーザ側の処理の (ii) である。 $\text{info}_{i,j}$ を用いて曖昧検索を
実施するが、どのような条件を定めるかは柔軟に設定でき
る。第5章においてその構成例を述べる。

RecSrchToken(K, w) :

(i) w の各文字を可能な限りメタ文字で置き換える。

(ii) $\tau_r \leftarrow (F_{K_1}(w), G_{K_2}(w), P_{K_3}(w))$ を出力する。

RecSearch(γ, C, τ_r) :

(i) 引数のチェックを行う： τ_r を (τ_1, τ_2, τ_3) と分解し、 T_r
に τ_1 が存在しない場合は空リストを返す。

(ii) レコメンドテーブルのエントリポイントを探す：

$$(\alpha_1, \alpha'_1) \leftarrow T_r[\tau_1] \oplus \tau_2$$

(iii) $R_1 \leftarrow A_r[\alpha_1]$ として τ_3 で復号する。つまり、 R_1 を
(ν_1, r_1) と分解し、 $\nu_1 \oplus H_1(\tau_3, r_1)$ を計算する。実際、
これは $(\text{Enc1}_{K_4}(w_1), \text{addr}_r(R_2))$ に等しい。

(iv) $i \geq 2$ に対して、 $\alpha_{i+1} = \mathbf{0}$ となるまで、前項と同様に
 R_i を復号する。

(v) $\{\text{Enc1}_{K_4}(w_1), \dots, \text{Enc1}_{K_4}(w_r)\}$ を出力する。

AddToken(K, f) :

(i) 準備：

($g_1, \dots, g_{\#f}$) を f 中の一意的な N -Gram の列とし、初
出現順は f と変わらないものとする。また、

$$\tau_a \leftarrow (F_{K_1}(f), G_{K_2}(f), \lambda_1, \dots, \lambda_{\#f})$$

とする。ただし、 $1 \leq i \leq \#f$ に対して、 $r_i \stackrel{\$}{\leftarrow} \{0, 1\}^k$,
 $r'_i \stackrel{\$}{\leftarrow} \{0, 1\}^k$ であり、さらに、

$$\lambda_i \leftarrow (F_{K_1}(g_i), G_{K_2}(g_i),$$

$$\langle \text{id}(f), \mathbf{0}, \text{Enc2}_{K_6}(\langle Q_{K_5}(\text{pos}_{g_i,1}),$$

$$\dots, Q_{K_5}(\text{pos}_{g_i,s}) \rangle \rangle \oplus H_1(P_{K_3}(g_i), r_i), r_i,$$

$$\langle \mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0}, F_{K_1}(g_i) \rangle \oplus H_2(P_{K_3}(f), r'_i),$$

$$r'_i)$$

であるとする。

(ii) $c_f \leftarrow \text{SKE1.Enc1}_{K_4}(f)$

(iii) (τ_a, c_f) を出力する。

Add(γ, c, τ_a) :

- (i) τ_a を $(\tau_1, \tau_2, \lambda_1, \dots, \lambda_{\#f})$ と分解し, T_d に τ_1 が存在しない場合は \perp を返す。(引数のチェック)
- (ii) インデックスの更新を行う。具体的には, $1 \leq i \leq \#f$ に対して以下を実行する。
 - (a) フリーリストの最後尾 φ および対応する削除配列のエントリ φ^* を探す。つまり, $(\varphi, \mathbf{0}) \leftarrow T_s[\text{free}]$, $(\varphi_{-1}, \varphi^*) \leftarrow A_s[\varphi]$ とする。
 - (b) 検索テーブルのフリーエントリを最後から1つ手前にセットする: $T_s[\text{free}] \leftarrow (\varphi_{-1}, \mathbf{0})$
 - (c) 1つめのノード N_1 へのポインタを復元する: $(\alpha_1, \alpha_1^*) \leftarrow T_s[\lambda_i[1]] \oplus \lambda_i[2]$
 - (d) 新しいノードを φ へ格納し, ポインタを N_1 を指すようにする: $A_s[\varphi] \leftarrow (\lambda_i[3] \oplus \langle \mathbf{0}, \alpha_1, 0, \dots, 0 \rangle, \lambda_i[4])$
 - (e) 検索テーブルを更新する: $T_s[\lambda_i[1]] \leftarrow (\varphi, \varphi^*) \oplus \lambda_i[2]$
 - (f) N_1 の対を更新する: $(D_1, r) \leftarrow A_d[\alpha_1^*]$, および, $A_d[\alpha_1^*] \leftarrow (D_1 \oplus \langle \mathbf{0}, \varphi^*, \mathbf{0}, \mathbf{0}, \varphi, \mathbf{0}, \mathbf{0} \rangle, r)$
 - (g) $A_s[\varphi]$ の対を更新する: $A_d[\varphi^*] \leftarrow (\lambda_i[5] \oplus \langle \varphi_{-1}^*, \mathbf{0}, \alpha_1^*, \varphi, \mathbf{0}, \alpha_1, \lambda_i[1] \rangle, \lambda_i[6])$
 - (h) $i = 1$ ならば, 削除テーブルを更新する: $T_d[\tau_1] \leftarrow \langle \varphi^*, \mathbf{0} \rangle \oplus \tau_2$
- (iii) ファイルを追加する: $C \leftarrow C \cup \{c\}$

DelToken(K, f) :

- (i) $\tau_d \leftarrow (F_{K_1}(f), G_{K_2}(f), P_{K_3}(f), \text{id}(f))$ を出力する。

Del(γ, C, τ_d) :

- (i) τ_d を $(\tau_1, \tau_2, \tau_3, \text{id})$ と分解し, T_d に τ_1 のエントリがなければ \perp を返す。(引数のチェック)
- (ii) L_f の先頭ノードを求める: $\alpha'_1 \leftarrow T_d[\tau_1] \oplus \tau_2$
- (iii) $1 \leq i \leq \#f$ に対して以下を実行する。
 - (a) D_i を復号する: $(D_i, r) \leftarrow A_d[\alpha'_i]$, $(\alpha_1, \dots, \alpha_6, \mu) \leftarrow D_i \oplus H_2(\tau_3, r)$
 - (b) D_i を削除する: $A_d[\alpha'_i] \xleftarrow{\$} \{0, 1\}^{6 \log \#A + k}$
 - (c) フリーリストの最後のノードを探す: $(\varphi, 0^{\log \#A}) \leftarrow T_s[\text{free}]$
 - (d) 検索テーブルのフリーエントリを D_i の対を指すようにする: $T_s[\text{free}] \leftarrow \langle \alpha_4, 0^{\log \#A} \rangle$
 - (e) D_i の対の場所を解放する: $A_s[\alpha_4] \leftarrow (\varphi, \alpha'_i, 0, \dots, 0)$
 - (f) 前方ノードの処理を行う: N_{-1} を D_i の対の前のノードとする。このとき, $(\beta_1, \beta_2, \gamma_1, \dots, \gamma_s, r_{-1}) \leftarrow A_s[\alpha_5]$; $A_s[\alpha_5] \leftarrow (\beta_1, \beta_2 \oplus \alpha_4 \oplus \alpha_6, \gamma_1, \dots, \gamma_s, r_{-1})$ として, N_1 の次を指すポインタを更新し, $(\beta_1, \dots, \beta_6, \mu^*, r_{-1}^*) \leftarrow A_d[\alpha_2]$; $A_d[\alpha_2] \leftarrow (\beta_1, \beta_2, \beta_3 \oplus \alpha'_i \oplus \alpha_3, \beta_4, \beta_5, \beta_6 \oplus \alpha_4 \oplus \alpha_6, \mu^*, r_{-1}^*)$ として, N_{-1} の対 (D_i の前のノード) を更新する。
 - (g) 後方ノードの処理を行う: N_{+1} を D_i の対の次の

ノードとする。このとき,

$$(\beta_1, \dots, \beta_6, \mu^*, r_{+1}^*) \leftarrow A_d[\alpha_3];$$

$$A_d[\alpha_3] \leftarrow (\beta_1, \beta_2 \oplus \alpha'_i \oplus \alpha_2, \beta_3, \beta_4,$$

$$\beta_5 \oplus \alpha_4 \oplus \alpha_5, \beta_6, \mu^*, r_{+1}^*)$$

として, N_{+1} の対のポインタを更新する。

(h) $\alpha'_{i+1} \leftarrow \alpha_1$ とする。

(iv) id に一致する暗号文を C から削除する。

(v) T_d から τ_1 を削除する。

Dec(K, c) :

- (i) K を (K_1, \dots, K_6) のように分解する。
- (ii) $m \leftarrow \text{SKE1.Dec1}_{K_4}(c)$ を出力する。

4.6 提案方式の安全性

スペースの都合上, 概略のみ述べる。安全性に関する表記方法は [2] と同様である。提案方式の安全性に関しては, ノードにポジション情報を含めた点とレコメンドテーブルを利用した点に関して, 別々に述べる。まず, ノードにポジション情報を含めた点に関しては以下が示される。

定理 1 SKE1 および SKE2 が CPA-安全であり, かつ, F, G, P, Q が擬似ランダム関数であるならば, ランダムオラクルモデルの下で, 提案 DSSE 方式は, 適応的選択キーワード攻撃にたいして $(\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3, \mathcal{L}_4)$ -安全である。□

この定理は背理法, 即ち, 提案方式が適応的選択キーワード攻撃に対し, ランダムオラクルモデルの下で $(\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3, \mathcal{L}_4)$ -安全でなく, かつ, SKE1 が CPA-安全であり, さらに, F, G, P, Q が擬似ランダムであるならば, SKE2 が CPA-安全でなくなることを示すことができる。実際, その証明はシミュレータ S を構成することによる。 S の構成については Kamara ら [2] で示されているものと同様であるが, SKE2 を利用する箇所, つまりノード N_i に関するシミュレートを行う部分に関しては, ノード情報に SKE2 の暗号文空間からランダムに取った暗号文 c を追加する。ここで仮定より, 提案方式が安全でないこと, および, Kamara ら [2] の方式が安全であることが示されていることより,

$$\text{Enc2}_{K_6}(\langle Q_{K_5}(\text{pos}_{i,g,1}), \dots, Q_{K_5}(\text{pos}_{i,g,s}) \rangle)$$

とランダムな暗号文 c を無視できない確率で識別できる攻撃者 A が存在する (A はポジション情報を任意に指定可能)。従って, A は $\text{Enc2}_{K_6}(m)$ とランダムな暗号文 c を無視できない確率で識別できる (m は A が任意に指定した平文)。つまり SKE2 は CPA-安全ではなくなる。

次に, レコメンドに関する安全性を考察する。レコメンド機能に用いた構造は Curtmola ら [1] の SSE-1 と同様であり, 従って, 安全性もそれと同等であることが示される。

5. 考察

本章では, 提案方式に対する以下の考察を行う。

変更に伴う問題点: 提案方式で Kamara ら [2] の方式から変更した点による不都合な点を述べる。

曖昧検索： N -Gram を用いた具体的な曖昧検索の方式と具体例を示す。

上記の他、レコメンドテーブルの活用法などを述べる必要があるが、紙面?の都合上割愛する。

5.1 変更に伴う問題点

ここでは、変更点に伴う問題点をいくつか挙げる。

長さ N 未満の文字列を検索できなくなる：

長さ N 未満の文字列は N -Gram に分けることが不可能なため、検索ができない。そのため、 N は大きくし過ぎないほうが良い*4。

検索結果に間違いが生ずる：

検索結果に非常にまれにだがミスが出現することがある。なぜなら、ポジション情報は暗号化されたものとランダムなダミーデータに分かれるが、ダミーデータが偶然復号でき、他のポジション情報に対して検索結果に含めるような条件の値になっているかもしれないからである。しかし、この現象が起こることは非常に稀であり無視できる。実際、 $Q_{K_5}(\text{pos}_{i,g,1}), \dots, Q_{K_5}(\text{pos}_{i,g,j})$ ($j \leq s$) が偶然復号でき、かつ、それが他のポジション情報に対し検索結果へ含めるような条件の値になる確率は、まず、検索語 w の長さに関しては、その長さが N の場合(つまり、ポジション情報が復号できた時点で検索結果へ含めることになる場合)が最高である。また、ポジション情報に関しては、 $Q_{K_5}(\text{pos}_{i,g,2}), \dots, Q_{K_5}(\text{pos}_{i,g,s})$ がダミーデータの場合には、検索結果に間違いが生じるのは $Q_{K_5}(\text{pos}_{i,g,2}), \dots, Q_{K_5}(\text{pos}_{i,g,j})$ ($j \leq s$) がすべて復号できる場合である。このとき、その確率は $j = 2$ の場合が最高である。以上より、ダミーデータにより検索結果へ間違いが生じる確率は $Q_{K_5}(\text{pos}_{i,g,2})$ が偶然、ランダムに選んだデータと一致する確率を超えない。 Q が擬似ランダムであるなら、その確率は無視できる。

その他：

上記以外にも、実行効率が低下する、レコメンドテーブルのサイズを予め決める必要がある、レコメンドテーブルが重複する、ユーザ負担が増加する、などの問題が挙げられるが、その詳細については割愛する。

5.2 曖昧検索

曖昧検索の構成として、具体的な例を示す。

具体的な曖昧検索の構成例：

具体的な条件として

- 長さ l の検索語 w を m 個の N -Gram (g_1, \dots, g_m) に分けたとき、 w の $\lceil m/2 \rceil$ 個以上の異なる N -Gram が $2l$ の長さ以下の文字列に出現した場合、 f を検索結果に含める

*4 N の増加に伴って N -Gram の数も指数関数的に増加するので、その理由からも N は大きくしない方がよいであろう。

などが考えられる。具体的にこの条件において、 $N = 2$ とした場合、 $w =$ “検索可能暗号” という検索語に対して、その N -Gram は $g_1 =$ “検索”， $g_2 =$ “索可”， $g_3 =$ “可能”， $g_4 =$ “能暗”， $g_5 =$ “暗号” である。このとき、文字が欠落している「検索可暗号」という文字が入った文書は検索結果に含まれる。なぜなら、「検索可暗号」という文字列中には、 g_1, g_2, g_5 の $3 \geq \lceil 5/2 \rceil$ 個の 2-Gram が出現し、その文字列長は $5 \leq 2 \times 6$ であるからである。その他、“検索可能な暗号”，“可能検索暗号” や “検索不能暗号” のような文字列を含む文書も、検索結果に含まれる。一方、“検索できる暗号” の場合、出現する 2-Gram は g_1, g_5 の 2 つであり、 $2 < \lceil 5/2 \rceil$ であることから、検索結果には含まれない。

6. まとめ

現在、クラウドストレージの発達などに伴い、SSE の技術に対する研究は急速に進んでいる。しかし、検索にかかるオーダが準線形であること、適応的な選択キーワード攻撃に対し安全であること、大きすぎないインデックスサイズ、およびインデックスに対するファイルの動的な追加・削除が可能であること、といった条件を満たすものは多くない。Kamara ら [2] の方式はこれらを満たした方式であったが、より実用に近づけるために、曖昧検索の技術やレコメンド、オートコンプリート等の技術を実現する必要がある。本論文では、 N -Gram を用いた曖昧検索およびレコメンド機能の実現方法を述べた。今後の課題として、サーバの負担軽減、具体的には、インデックスサイズの縮小が挙げられる。提案方式では、レコメンドテーブルをユーザの数だけ用意する必要があったが、レコメンドテーブルは 1 つで十分であり、サーバごとに 1 つのレコメンドテーブルを用いたレコメンド機能を実装すべきと思われる。また、レコメンドそのものの機能に関しても、単語やフォーマット間だけでなく、ユーザが入力する文章から残りの検索内容をレコメンドできるようにすることが望まれる。

参考文献

- [1] R. Curtmola, J. Garay, S. Kamara and R. Ostrovsky: “Searchable symmetric encryption: Improved definitions and efficient constructions”, Proceedings of the 13th ACM conference on Computer and Communications Security (CCS), pp.79-88, 2006.
- [2] S. Kamara, C. Papamanthou and T. Roeder: “Dynamic Searchable Symmetric Encryption”, Proceedings of The 2012 ACM conference on Computer and Communications Security (CCS), pp.965-976, 2012.
- [3] C. E. Shannon: “A mathematical theory of communication”, Bell System Technical Journal, vol.27, pp.379-423, 623-656, 1948.
- [4] D. Song, D. Wagner and A. Perrig: “Practical techniques for searches on encrypted data”, In IEEE Symposium on Research in Security and Privacy, pp.44-55. IEEE Computer Society, 2000.