

Require: d : 転置インデックス
Require: X : 検索クエリの特徴集合
Require: τ : 必要なオーバーラップ数

```

1:  $X$  の要素を,  $|get(d, X[k])|$  の昇順に並び替える
2:  $M \leftarrow \{\}$ 
3: for  $k = 0$  to  $(|X| - \tau)$  do
4:   for all  $i \in get(d, X[k])$  do
5:      $M[i] \leftarrow M[i] + 1$ 
6:   end for
7: end for
8:  $R \leftarrow \{\}$ 
9: for  $k = (|X| - \tau + 1)$  to  $(|X| - 1)$  do
10:  for all  $i \in M$  do
11:    if binary_search(get( $d, X[k]$ ),  $i$ ) then
12:       $M[i] \leftarrow M[i] + 1$ 
13:    end if
14:    if  $\tau \leq M[i]$  then
15:       $i$  を  $R$  に追加
16:       $i$  を  $M$  から削除
17:    else if  $M[i] + (|X| - k - 1) < \tau$  then
18:       $i$  を  $M$  から削除
19:    end if
20:  end for
21: end for
22: return  $R$ 

```

図 2: τ オーバーラップ問題のアルゴリズム

性質 2 要素数が k の集合 Z と, 要素数が任意の集合 Y がある. 要素数が h となる任意の上位集合 $X \supseteq Z$ を考える. もし, $|Z \cap Y| < \theta$ ならば, $|X \cap Y| < \theta + h - k$ である.

図 2 に, これらの性質を利用した τ オーバーラップ問題の解法を示した. 性質 1 より, $|X \cap Y| \geq \tau$ となるためには, $(|X| - \tau + 1)$ 個の転置リストの中に文字列 y の SID が, 少なくとも 1 回は含まれていなければならない. そこで, 2 から 7 行目で $(|X| - \tau + 1)$ 個の転置リストを走査し, τ オーバーラップ問題の解となり得る文字列の SID を収集している. このとき, 転置リストに含まれる SID の数が少ないほど, 考慮すべき文字列の候補をコンパクトにできるので, 1 行目で転置リストの順番を要素数の少ない順に並び替えている. 9 行目から 21 行目では, それぞれの候補 SID (i) が残りの転置リスト中に含まれるかどうか, 二分探索で調べ, 頻度カウンタ $M[i]$ をインクリメントする. このとき, 性質 2 を用いて, 候補 i がこれから走査する転置リストの全てに含まれた場合の $|X \cap Y|$ の上限値を $M[i] + (|X| - k - 1)$ として求め, この値が τ よりも小さければ, 候補 i を枝刈りする.

3 評価実験

Google Web 1T コーパス (LDC2006T13) に含まれる全ての単語ユニグラム (13,588,391 文字列) をデータセットとし, 評価実験を行った. 1 つの文字列当たりに含まれる文字 tri-gram の数は約 10.3 で, データセット全体には 301,459 種類の文字 tri-gram が含まれる. 提案手法のシステムを C++ で実装し, 転置インデックスの構築には CDB++^{*} を用いた. 提案手法を実装したライブラリは, SimString[†] として公開している. 実験環境は, インテルの Xeon 5140 CPU (2.33 GHz) と 4GB の主記憶を搭載したサーバーで, 類似文字列検索に必要なデータベース (601 MB) を 557.4 秒で構築した.

従来手法として, DivideSkip と Locality Sensitive Hashing (LSH) を実装した. LSH では, データ中の全ての文字列を 64 ビットのハッシュ値に変換した. 検索文字列のハッシュ値 $h(x)$ とのハミング距離が δ 以内の文字列を高速に見つけ出すため, データ中の全てハッシュ値のビット列をランダムに入れ替え, 入れ替え後のハッシュ値を昇順に整列したリストを

^{*}<http://www.chokkan.org/software/cdbpp/>

[†]<http://www.chokkan.org/software/simstring/>

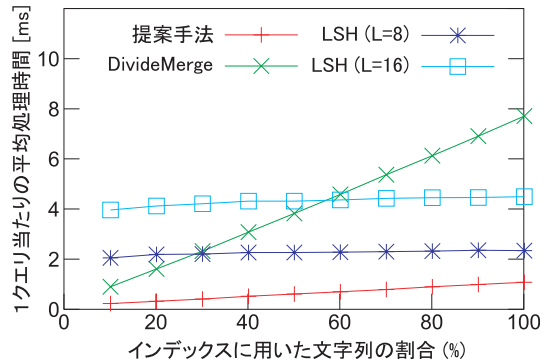


図 3: Google Web 1T コーパスにおけるクエリ処理時間

L 個用意した. このように作られた整列済みハッシュ値リストに二分探索を適用し, 探したいハッシュ値 (ビット列入れ替え済み) に最も近い文字列, 及びその周辺の $B - 1$ 個の文字列を候補とし, ハッシュ値 $h(x)$ とのハミング距離が δ 以内の文字列の中で, 類似度が α 以上のものを出力した. ここで, δ, L, B は LSH の類似文字列検索の速度と精度のトレードオフを制御するパラメータであり, $\delta = 16, B = 16$ に固定し, $L = 8$ もしくは 16 とした.

図 3 に, データセットの 10%–100% の文字列をインデックスし, 同データセットからランダムに選んだ 1,000 件をクエリ文字列として, 閾値 $\alpha = 0.8$ で類似文字列検索を行う際の, 1 クエリ当たりの平均レスポンス時間を示した. 提案手法が最も高速で, 100% のデータサイズの時に, 1 クエリを 0.99 [ミリ秒/クエリ] で処理していた. 図 3 の描画範囲外であるが, クエリ文字列とデータセット中の全ての文字列との類似度を計算するナイーブな実装は, 166.1 [秒/クエリ] の処理速度で, 提案手法は約 17 万倍高速に類似文字列検索が行える. また, これも図 3 に載せていないが, τ オーバーラップ問題を解く際, 性質 1 と性質 2 を利用せず, 転置リスト中の SID を全て走査する実装は, 403.7 [ミリ秒/クエリ] の処理速度であり, 提案手法は 407 倍高速であった. このことから, τ オーバーラップ問題を解く際に, 性質 1 や性質 2 で文字列の候補を絞り込んだり, 探索途中で枝刈りをするものの効果が伺える.

提案手法は, 既存手法の DivideSkip 法よりも約 8 倍高速であった. LSH の処理時間は, 約 2.2 [ミリ秒/クエリ] ($L = 8$ のとき), 及び約 4.3 [ミリ秒/クエリ] ($L = 16$ のとき) であった. LSH は δ, L, B の値を小さくすることでクエリ処理を高速化できるが, 検索の精度 (再現率) が低下する. 今回の実験条件において, LSH の類似文字列検索の再現率は, 36.5% ($L = 8$), 40.6% ($L = 16$) であり, 常に 100% の再現率が保障される提案手法とは大きな差がある. また, LSH はメモリ使用量が多くなりやすく, 今回の実験では, 2.6GB ($L = 8$), 4.3GB ($L = 16$) の主記憶を消費していた. これに対し, 提案手法のメモリ使用量は 601MB であった. 以上の実験結果から, 提案手法は大規模な類似文字列検索を, 実用的な計算量と記憶領域で実現していることが分かった.

参考文献

- [1] Chen Li, Jiaheng Lu, and Yiming Lu. Efficient merging and filtering algorithms for approximate string searches. In *ICDE '08: Proceedings of the 2008 IEEE 24th International Conference on Data Engineering*, pp. 257–266, 2008.
- [2] Alexandr Andoni and Piotr Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Communications of the ACM*, Vol. 51, No. 1, pp. 117–122, 2008.
- [3] Arvind Arasu, Venkatesh Ganti, and Raghav Kaushik. Efficient exact set-similarity joins. In *VLDB '06: Proceedings of the 32nd International Conference on Very Large Data Bases*, pp. 918–929, 2006.