

高速に類似文字列ペアを発見するビット並列フィルター

山田太樹[†] 清水佳奈[†][†]早稲田大学基幹理工学部

1 序論

編集距離の近い文字列ペアの発見は、ゲノム情報解析や時系列情報の解析等、様々な分野において重要な課題である。例えばゲノム情報解析では、ゲノム配列のクラスタリングやアセンブリ（ゲノム配列の断片をつなげて一つの配列を復元するタスク）を行う際に、大量の編集距離を計算する必要がある。このような背景から本研究では、長さがそれぞれ n, m （ただし、 $n < m$ とする）の文字列が与えられた場合に、編集距離が閾値 θ 以内であるか否かを高速に判定する手法の開発を目的とする。編集距離そのものを高速に計算する従来研究では、Myers によるワード長 w のとき $O(nm/w)$ を達成するビット並列アルゴリズムが知られている [1]。一方で、文字列のセットが与えられた場合にそれらの中から類似するペアを発見する手法としては SlideSort 法 [2] や ED-Join 法 [3] が知られている。SlideSort や ED-Join では、実際に用いられるデータセットが疎であることに着目し、類似性の低いペアを計算量の低い方法によって予め除外し、編集距離の計算回数を減らすことで全体の計算量を大幅に向上させている。これらの手法が対象としている全ペア列挙問題の他にも、特定のパターンとの編集距離が近いデータの探索などでは、類似性の低いペアを評価する回数が多くなりがちであるため、類似文字列ペア発見の問題に対して直接的に編集距離を計算する Myers の手法が適しているとは言い難い。そこで本研究では、ED-Join 法や SlideSort 法で用いられている鳩ノ巣原理を応用して、編集距離が遠いペアを計算量 $O(\theta n/w)$ で判定することのできるビット並列フィルター F を構築し、 F が受理したペアに対してのみ編集距離を計算する手法を提案する。

2 提案手法

文字列 s の i 番目の文字を $s[i]$ と表す。 s と t の編集距離 $Editdist(s, t)$ について、以下の定理が成り立つことが知られている [2]。

定理. 等しい長さ l の文字列 s と t について、 $Editdist(s, t) \leq \theta$ であるならば、

$$s[p] = t[q], q \in \{p - \lfloor \theta/2 \rfloor, \dots, p + \lfloor \theta/2 \rfloor\}$$

を満たす q が存在するような p が少なくとも $l - \theta$ 個存在する。

上記定理で示される必要条件を満たさないペアの編集距離は必ず θ よりも大きくなるため、編集距離の計

図 1: $\theta = 2$ の時の s と t の一致する文字の組の例

算を行わなくとも類似性の低いペアと判断することができる。説明のため、 $\theta = 2$ における具体例を図 1 に示す。例では、 $\theta = 2$ より、少なくとも $5 - 2 = 3$ 箇所において s がオフセット幅 1 以内で t と一致すればよい。左の s と t の組では、4 つの一致が存在するため必要条件を満たす。一方、右の組では 2 つの一致しか存在しないため必要条件を満たさない。よって、右の組については編集計算の計算を行わずとも、非類似ペアと判断することができる。このように、上記定理は類似ペアを発見する際のフィルターとして利用することが可能である。本研究では、上述の必要条件を高速に計算するビット並列アルゴリズムを提案し、多数のペアを評価する際のならし計算量を削減することを目指す。

2.1 フィルタ F の基本構築方法

まず、オフセット $j \in \{-\lfloor \theta/2 \rfloor, \dots, \lfloor \theta/2 \rfloor\}$ について、 s と t の一致を示すビットベクトル

$$F_j[i] = \begin{cases} 1 & (s[i] \neq t[i+j]) \\ 0 & (s[i] = t[i+j]) \end{cases}$$

を求める。これは、以下により計算できる。

$$F_j = (s \oplus (t \gg j))$$

ただし、 \oplus は排他的論理和、 $t \gg j$ は配列 t を j の幅で右シフトする演算を示す。 F_i の論理積

$$M = F_{-\lfloor \theta/2 \rfloor} \wedge \dots \wedge F_{\lfloor \theta/2 \rfloor}$$

における各ビットは、 s と t がオフセット幅 $\lfloor \theta/2 \rfloor$ 以内で一致するか否かを示す。このため、 M に含まれる 1 のビットの数が θ 以下であれば、定理を満たす。具体例を図 2 に示す。例では s の 2 文字目と 4 文字目がそれぞれオフセット幅 1 以内で t に一致するので、 M の対応するビットは 0、それ以外は 1 となる。

2.2 ビット列のブロック化による性能の改良

2.1 の手法では、文字の種類が少ない場合に偶然の一致が多くなり、編集距離が遠いにもかかわらずフィルターに受理されてしまうペアが増加する問題がある。そこで、一文字の一致ではなく、長さ k の文字列ブロッ

Efficient bit-parallel filtering algorithm for finding similar string pair in edit distance
Taiki YAMADA[†] and Kana SHIMIZU[†]
[†]School of Fundamental Science and Engineering, Waseda University

