

# ギブスサンプラーに基づく アミノ酸配列モチーフの高精度抽出法

高橋 誉文<sup>†1</sup> 北上 始<sup>†1</sup> 福本 翔平<sup>†1</sup> 森 康真<sup>†1</sup> 田村 慶一<sup>†1</sup>

アミノ酸配列データベースから類似部分配列を抽出することとして知られている従来のギブスサンプリング法の抽出精度を向上するために、多重整列化に基づく新しい方法を提案する。従来のギブスサンプリング法の抽出精度は初期値に大きく左右される。この点に着目し、提案手法では、できるだけ良い初期値を計算するため、配列データセットに対して多重整列化を行い、ある幅のウインドウを多重整列上にスライドさせ、 $p$  値が最小となるウインドウ領域（類似部分配列）を初期値として選択する。多重整列化によって挿入されるギャップについては、ランダムに文字を当てはめる場合とすべての文字が等確率に現れる場合を比較する。また、ギブスサンプリングで利用される擬似度数に進化的な知識を導入し、抽出される類似部分配列としての配列モチーフ（進化的に保存される配列パターン）の抽出精度を向上している。

## Method for High-Precision Motif Extraction based on Gibbs Sampler in Amino Acid Sequences

YOSHIFUMI TAKAHASHI<sup>†1</sup> KITAKAMI HAJIME<sup>†1</sup> FUKUMOTO SYOUHEI<sup>†1</sup> MORI YASUMA<sup>†1</sup>  
TAMURA KEIICHI<sup>†1</sup>

In order to improve the extraction accuracy of the existing, well-known Gibbs sampling method for extracting similar subsequences from amino acid sequence databases, we propose a new extraction method based on multiple sequence alignment in the sequence dataset. The extraction accuracy of the existing Gibbs sampling method is highly dependent on the initial solution selected randomly. In focusing on this point, the proposed method performs multiple sequence alignment for the sequence dataset to calculate the best possible initial solution. After that, we slide the aligned sequences on a window of a certain width and select the window region including the set of subsequences, where  $p$ -value is minimized, as the initial solution. In order to confirm the effectiveness of the proposed method, we carried out comparative experiments with random distribution and equal distribution. Moreover, we improve the accuracy of the existing Gibbs sampling method by using an amino acid substitution matrix as the knowledge of molecular evolution for pseudocount.

### 1. はじめに

アミノ酸の配列モチーフは、生物進化の過程で保存された生物学的な機能やタンパク質の立体構造と深く関係する。ギブスサンプラーは、配列データベースから配列モチーフにできるだけ近い類似部分配列を抽出する方法であり、多くの研究[1][2][3]が行われている。また、ギブスサンプラーを用いて、配列モチーフに近い類似部分配列を抽出する Web サービス[4]も行われている。

機能が未知のアミノ酸配列から配列モチーフを特定することができれば、例えば、その構造の推定を初めとして、SURFNET[5][6]などにより、活性部位や結合ポケットなどのような機能部位の特定につなげることができるものと期待されている[7]。また、このような機能部位の解析は、医薬品の開発に重要な分子シミュレーションや分子設計などを支援するものと期待されている。

ギブスサンプラーは、Geman 兄弟[8]によって画像処理の分野で利用され、画像復元に対する有効な統計的手法として紹介された。パイオインフォマティクス分野では、ギブスサンプラーは、アミノ酸配列データセットの各配列か

ら配列モチーフに最も近い類似部分配列を1個のみ抽出する方法として、Lawrence ら[1]によって紹介された。本稿では、これをサイトサンプラーと呼び、各配列から0個以上の配列モチーフに最も近い類似部分配列を抽出するモチーフサンプラーと区別する[2][9][10]。いずれも、マルコフ連鎖モンテカルロ法[11][12]の1つとして分類されているギブスサンプラーに該当する。なお、モチーフサンプラーの計算では、サイトサンプラーの計算結果が利用されているが、モチーフサンプラーの計算精度を向上するには、サイトサンプラーの計算結果が生物進化の過程で保存された配列モチーフにできるだけ近いことが重要である。

パイオインフォマティクス分野でのギブスサンプラーは、アミノ酸配列データセットから配列モチーフの抽出問題に適用する研究[1][2][13][14][15]と、DNA 配列データセットから DNA 配列上のタンパク質結合部位（配列モチーフ）の識別問題に適用する研究[3][4][7][16][17][18]とに分類される。

アミノ酸配列データセットや DNA 配列データセットから配列モチーフに最も近い類似部分配列を抽出するために、Lawrence らのサイトサンプラーに対して、さまざまな改良法が提案されている。これらの手法は、機械学習や統計学などを導入し、数学的な最適解の高精度な探索に努力が払

<sup>†1</sup> 広島市立大学  
Hiroshima City University

われてきたが、生物進化の過程で発生するアミノ酸置換の相対頻度を表わす表（アミノ酸置換行列）をサイトサンプラーの計算モデルに十分に反映していないため、従来の計算モデルから数学的に厳密な類似部分配列が得られても、その類似部分配列はバイオインフォマティクス分野の専門家にとって興味ある配列モチーフから外れていることがある。

本論文では、アミノ酸配列データベースから配列モチーフにできるだけ近い類似部分配列を抽出するために、Lawrence らによって紹介されたサイトサンプラーに基づき、新しい計算モデルに基づく計算手法を提案する。

サイトサンプラー（ギブスサンプラー）は、ユーザにより与えられた配列モチーフの長さ  $K$  をもとに、 $N$  本の配列データを含む配列データセットから配列モチーフの候補である  $N \times K$  の整列行列（ $K$ -類似部分配列の集合）を出力する方法である。配列データセットに含まれる  $N$  本の配列データの長さは同一ではないが、簡単のため  $L$  とすると、配列データセットには、 $(L-K+1)^N$  個の  $N \times K$  の整列行列が存在するため、これらの中から配列モチーフに近い類似部分文字列を直接探索するのは現実的ではない。サイトサンプラーでは、配列データセットに含まれる  $N$  本の各配列データからランダムに選択された  $K$ -部分配列の集まりを候補配列集合として選択し、それらを初期値とすることにより、候補配列集合を更新しながら配列モチーフとして最も確からしい類似部分配列を見つけ出そうとする確率的アルゴリズムである。しかしながら、アミノ酸配列データセットを対象にした従来のサイトサンプラーには、以下のような問題がある。

#### (1) 初期値のランダム性

サイトサンプラーは、局所最適解を回避する工夫がなされていないため、計算精度が初期値に大きく依存し、計算結果に最適解が保証されていない。初期値のランダム性を低減するために、 $N$  本の配列からいくつかのサンプルをランダムに選択し、それらのサンプルのみから貪欲探索により最良の初期値を取り出す方法[2][10]があるが、初期値のランダム性に関する本質的な解決にはなっていない。

#### (2) 進化的知識としてのアミノ酸置換行列の欠如

配列データセットからランダムに選択された 1 本の配列データを  $Z$  とする。サイトサンプラーでは、候補配列集合から  $Z$  上に存在する部分配列  $X$  を探し、それまでに計算された文字の出現確率を表すプロファイル行列に適合する部分配列  $X'$  を  $Z$  上から確率的に選択し、候補配列集合内の  $X$  を  $X'$  に置き換え、候補配列集合を更新している。しかし、アミノ酸の置換のし易さを数値化したアミノ酸置換行列[20]についての知識が考慮されていないため、 $X'$  が適切に選択されない。

本論文では、これらの二つの問題点を解決するために、Thompson[21]や Larkin[22]の文献で提案されている案内木に基づく多重整列化(マルチプルアラインメント)に加え、Henikoff らが提案した擬似度数[23][24]をサイトサンプラーの計算モデルに組み込んだ新しい抽出法を提案する。なお、案内木は近隣結合法[25]などで作成された分子進化系統樹（生物進化の枝分かれの様子を描画した木構造）が多

く利用されている。提案手法の要点は以下のとおりである。

- (1) 初期値のランダム性の問題を解決するために、できるだけ良質な初期値を計算する。そのために、まず、配列データセットに対して、分子進化系統樹に基づく多重整列化[21][22]を利用し、 $N \times L$  の整列行列を導出する。この整列行列にスライディングウィンドウ法を適用し、ある幅をもつウィンドウを左から右に 1 文字ずつスライドすることにより、全てのウィンドウ領域を列挙する。これらのウィンドウ領域から配列モチーフに最も近いウィンドウ領域[26]を選択する。この選択の評価尺度として相対エントロピーを用いる。ただし、配列モチーフの両端にはどちらもギャップが含まれないため、ウィンドウ領域の両端の少なくとも一方に閾値以上のギャップ数が存在する場合は、そのウィンドウ領域を選択の対象から除外する。
- (2) 進化的知識としてのアミノ酸置換行列が欠如しているという問題に対処するために、プロファイル行列の計算式に現れる擬似度数に、アミノ酸置換行列[20][27]に基づく知識を組み込む[26]。アミノ酸置換行列は、進化的な知識であり、生物進化の過程で発生するアミノ酸置換の相対頻度を行列[20]で表現したものとして知られている。

以上のような提案方法について、実験により有効性を確認するために、5 種類のデータセットを用いて、提案手法により抽出された類似部分配列が既に登録されている配列モチーフにどれだけ近いかを評価した結果、1 件のデータセットを除き約 90% 以上も近いことがわかった。

本論文の構成は以下のとおりである。2 章では、関連研究について述べる。3 章では従来手法であるサイトサンプラー（ギブスサンプラー）について述べる。その中でプロファイル行列を用いた出現頻度の計算法や背景頻度の計算法を説明する。4 章では多重整列化に基づく提案手法について述べ、5 章では従来手法と提案手法による計算結果を比較して評価を行う。最後の 6 章では本論文のまとめと今後の課題について述べる。

## 2. 関連研究

アミノ酸配列や DNA 配列などの配列データセットから配列モチーフに最も近い類似部分配列を抽出する方法には、列挙法[28][29][30]、隠れマルコフモデル[31]、ギブスサンプラー[1][2][3][4][7][13][14][15][16][17][18][32][33]などがある。

列挙法とは、配列モチーフ内に存在するワイルドカード(任意の文字と一致する記号)を考慮して、PrefixSpan[28]のアプローチにより頻出パターンをすべて抽出する方法である。しかし、非常に多くの頻出パターンが列挙されるため、その中から配列モチーフに最も近い頻出パターンを探し出すことが困難である。進化的知識としてアミノ酸置換行列（生物進化の過程で発生するアミノ酸置換の相対頻度を表わす表）を計算モデルに組み込んでいないことも原因の一つと考えられる。

隠れマルコフモデルは、ユーザが予め設計したネットワーク構造に基づいて、状態遷移確率や出力確率などのモデルパラメータを計算する方法である。このネットワーク構造に対するモデルパラメータはプロファイル HMM と呼



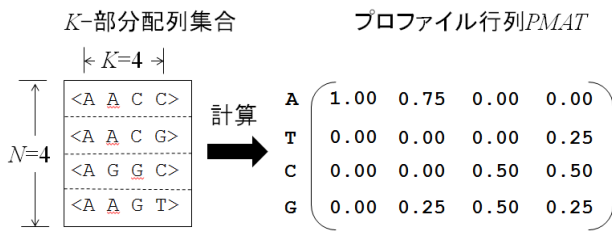


図 2 K-部分配列集合に対する出現頻度行列  $AMAT=(c_{ij})$

プロフィール行列として出現頻度行列  $AMAT=(c_{ij})$  を使用してみよう. アミノ酸配列データセット  $DS$  内に含まれる配列を  $Z$  とし, その長さを  $|Z|$  と表記すると, 配列  $Z$  には  $|Z|-K+1$  個の  $K$ -部分配列が存在する ( $|Z| \geq K$ ). それらの中に含まれる  $K$ -部分配列の 1 つを  $x = \langle \alpha_1 \alpha_2 \dots \alpha_K \rangle$  と表記すると, プロファイル行列を用いて, その出現頻度 (生起確率)  $P_x$  を次式のように計算することができる.

$$P_x = c_{11} \times c_{22} \times \dots \times c_{KK}$$

ただし,  $x$  の部位に存在する文字  $\alpha_i$  がプロフィール行列の  $i$  行目に対応し,  $c_{ij}$  は  $j$  列目における文字  $\alpha_i$  の出現する頻度を意味する. これにより計算された  $x$  の出現頻度  $P_x$  が高ければプロフィール行列の計算に用いられた  $K$ -部分配列集合に類似し, 低ければ類似しないと解釈される.

さて, 配列データの本数が少ないことなどが原因で,  $AMAT=(c_{ij})$  の要素  $c_{ij}$  に 0 が出現すると, 他の要素の値がいくら大きくても出現頻度  $P_x$  が 0 になってしまうことがある. これを避けるため, 文献 [1][2][9][10] では, ベイズ統計解析を導入し, アミノ酸配列データセット  $DS$  から配列  $Z$  を取り除いた  $N-1$  本の  $K$ -部分配列に対する, プロファイル行列  $PMAT=(p_{ij})$  を以下のように定義している.

$$p_{ij} = \frac{(c_{ij} + b_i)}{((N-1) + B)} \quad (1)$$

ただし,  $N$  は  $DS$  に含まれる配列総数,  $B$  は  $\sqrt{N}$  と定める.  $N-1$  本の  $K$ -部分配列は,  $DS$  からある配列データ  $Z$  を除いた  $DS'$  から得られる. また,  $p_{ij}$  の  $i$  行目に該当する文字の全配列に対する相対出現頻度を  $f_i$  とすると,  $p_{ij}$  の  $i$  行目に該当する文字の疑似度数  $b_i$  は  $f_i \times B$  としており, これにより分子のゼロ除算を回避することができる.

### 3.2 背景頻度行列

図 3 の例で示されるように, 配列データセット  $DS$  から  $K$ -部分配列集合を除いた結果は背景配列集合と呼ばれる. 背景頻度行列  $BMAT$  とは, 背景配列集合に出現する文字  $\alpha$  の頻度を表現する  $M \times 1$  の行列  $BMAT=(q_i)$  である.  $q_i$  は,  $i$  行目に対応する文字  $\alpha_i$  の背景頻度を意味し, 文字  $\alpha_i$  の背景頻度は, 背景配列集合内に存在する文字の総数に対する文字  $\alpha_i$  の出現頻度 (生起確率) である.

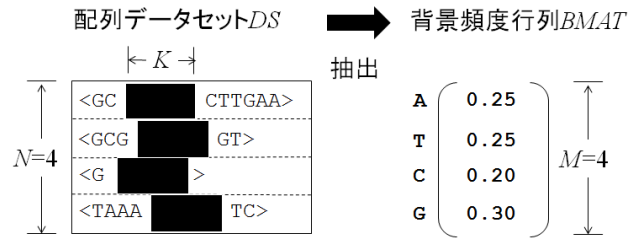


図 3 配列データセット  $DS$  と背景頻度行列  $(q_i)$

$DS$  内のある配列データ  $Z$  に存在する  $K$ -部分配列  $x = \langle \alpha_1 \alpha_2 \dots \alpha_K \rangle$  の背景頻度  $Q_x$  は以下のとおりである.

$$Q_x = q_1 \times q_2 \times \dots \times q_K$$

ただし,  $q_i$  は文字  $\alpha_i$  の背景出現を意味する. これにより計算された  $x$  の背景頻度が高ければ,  $K$ -部分配列集合に非類似となり, 低ければ類似していると解釈できる.

### 3.3 オッズ比

素朴なサイトサンプラーでは,  $DS$  内に含まれる配列  $Z$  からプロフィール行列に適合する  $K$ -部分配列  $x$  を計算するために出現頻度  $P_x$  だけを用いるが, 文献 [1][2][9][10] では, 次式で定義されるオッズ比  $A_x$  あるいは対数オッズ比  $\log_2 A_x$  が用いられている.

$$A_x = P_x \div Q_x$$

オッズ比  $A_x$  が高い  $x$  を配列データ  $Z$  から選択することは,  $K$ -類似部分配列集合に類似している (出現頻度  $P_x$  が高い) と同時に背景配列集合に似ていない (背景頻度  $Q_x$  が低い)  $x$  を配列データ  $Z$  から選択することを意味する. 逆にオッズ比  $A_x$  が低ければ,  $x$  は  $K$ -類似部分配列集合に似ていないと同時に背景配列集合に近い  $x$  を配列データ  $Z$  から選択することを意味する.

### 3.4 アルゴリズム

サイトサンプラーは  $N$  本の配列からなる  $DS$  からランダムに選択された配列  $Z$  を用いる事で, プロファイル行列と背景頻度行列を算出し, 出現頻度が高くかつ背景頻度の低い  $K$ -部分文字列集合を抽出する処理を行っている. そのアルゴリズムは以下のとおりである.

- ①  $DS$  の各配列に対して  $K$ -部分配列の開始点  $st_i$  をランダムに選び, それらを行列順に整列させた  $N$  本の  $K$ -部分配列集合  $\{st_1, st_2, \dots, st_N\}$  を初期値とする.
- ②  $DS$  からランダムに一つの配列  $Z$  を選択する.
- ③  $DS' = DS - \{Z\}$  から,  $N-1$  本の  $K$ -部分配列に対するプロフィール行列  $PMAT=(p_{ij})$  と背景配列集合に対する背景頻度行列  $BMAT=(q_i)$  を算出する.
- ④ 配列  $Z$  内に存在する  $|Z|-K+1$  個の  $K$ -部分配列  $x$  から, それぞれの出現頻度  $P_x$  および背景頻度  $Q_x$  を計算し, 類似度スコア  $A_x$  を計算する. すなわち,  $x$  のオッズ比  $A_x = P_x \div Q_x$  あるいは対数オッズ比  $\log_2 A_x$  を算出する.
- ⑤  $\{A_1, A_2, \dots, A_{|Z|-k+1}\}$  となった各値から, 比例した確率

で  $A_r$  を選択し,  $A_r$  に対応する  $K$ -部分配列集合の新たな開始点  $st_z'$  を  $st_z$  に置き換える.

- ⑥ ②~⑤をユーザが定めた回数分繰り返す. 繰り返し回数は多いほど良い結果が出力されるが, その分計算時間が大幅に増加する.

### 3.5 相対エントロピー

抽出した  $K$ -部分配列集合が配列モチーフとしてどのぐらい近いかを評価するために, 相対エントロピーが利用されている[1][2][9][10]. 相対エントロピーは, 次のように定義される.

$$F = \sum_{i=1}^K \sum_{j=1}^M c_{ij} \log_2 \left( \frac{p_{ij}}{q_i} \right) \quad (2)$$

ただし, この定義では,  $M \times K$  の出現頻度行列  $AMAT=(c_{ij})$ ,  $M \times K$  のプロファイル行列  $PMAT=(p_{ij})$ ,  $M \times 1$  の背景頻度行列  $BMAT=(q_i)$  が用いられている.  $K$ -部分配列集合の相対エントロピーが大きければ配列モチーフに近く, 小さければ配列モチーフから離れているものと判断している.

### 3.6 サイトサンプラーの問題点

サイトサンプラーのアルゴリズムは, 計算精度が初期値 (3.4 節の①で与えられる開始点) に大きく依存する. 初期値による影響を少なくするため, SA 法[41]や遺伝的アルゴリズム[38]の利用が考えられる. しかし, 従来の相対エントロピーを SA 法の目的関数や遺伝的アルゴリズムの適応度に利用されるため, 配列モチーフとは異なる類似部分配列が得られることがある. すなわち, 従来の相対エントロピーの計算式は,  $K$ -部分配列集合内の配列どうしの類似度を表すが,  $K$ -部分配列集合内の各配列と配列モチーフとの近さを表すものではない.

サイトサンプラーで計算される  $K$ -類似部分配列を配列モチーフにできるだけ近づけるためには, できるだけ品質の高い初期値を計算することや相対エントロピーの計算式を改善することが重要となる.

## 4. 提案手法

提案手法では, ランダムに初期値を与えることを避けるため, 多重整列を用いて初期値 ( $K$ -類似部分配列集合) の探索を行う. この多重整列は, 予め, 分子進化系統樹を用いた多重整列化 (マルチプルアラインメント) により求めたものである. 従来のサンプラーでは, プロファイル行列は, 処理手順④における出現度数を初めとして, 処理手順⑥の相対エントロピーの計算で利用されており, プロファイル行列の計算式には, 擬似度数が組み込まれている. 提案手法では, 従来の擬似度数の計算式に, アミノ酸置換のし易さを数値化したアミノ酸置換行列 (進化的な知識) を導入している. 以上の提案手法によるサイトサンプラーにより, 配列モチーフにできるだけ近い  $K$ -類似部分配列を抽出しようとしている.

本章では, 先ず, 多重整列化の方法について述べた後, 多重整列から安定した初期値を探索する方法について述べ

る. 次に, プロファイル行列の擬似度数の計算式にアミノ酸置換行列を導入する方法について述べる. 最後に, 配列モチーフとしての  $K$ -類似部分配列集合を抽出するアルゴリズムについて述べる.

### 4.1 多重整列化の方法

$N$ 本の系列データ (時系列データやテキストデータなど) に対する多重整列化では, 編集距離を尺度とする動的計画法[42]が利用されている ( $N \geq 3$ ). この多重整列化は, 非類似度スコア (編集距離) の累積が最小となるように, 系列データの適当な場所にギャップを挿入し, 系列データ間の文字の対応付けを行うことにより, すべての系列データの長さを揃えている.

$N$ 本の配列データを含むアミノ酸配列データセットに対する多重整列化では, 非類似度スコアの累積が最小となる方法を利用せず, 類似度スコアを尺度とする動的計画法[42]が利用されている. これにより累積類似度スコアが最大になるように, 系列データの適当な場所にギャップを挿入し, 配列データの長さを揃えている [43][44]. なお, 類似度スコアの計算には, 生物進化の過程で発生する文字どうしの置換のしやすさを表す置換行列を利用している. また, 進化的に近縁の配列データどうしを整列化した方が精度向上につながるため, 分子進化系統樹を案内木とする多重整列化を行っている [21][22].

$N$ 本の配列データを多重整列化することにより, 配列長が  $L$  になったとしよう. 以下では, 多重整列化した結果を  $N \times L$  行列とみなし, これを  $N \times L$  の整列行列と呼ぶ. 図4に多重整列化により得られる整列行列の例を示す. ギャップと呼ばれる記号 (-) は, 多重整列化において, 文字の挿入や削除の操作によって追加される記号である. これにより, 大域的に一番類似するように配列データの長さを一致させている. 図では, 多重整列化によりアミノ酸配列データセット  $DS$  の長さを一致させることにより得られる整列行列を  $DS'$  としている. 分子進化学の専門家は, この整列行列  $DS'$  を用いて, 配列モチーフ領域を経験的に探索している[31].

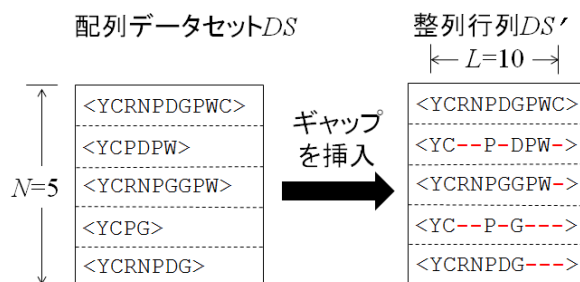


図4 多重整列化の例

多重整列化により配列データ上に挿入されたギャップは長さ  $K$  の配列モチーフ領域上にも挿入される. このため, 整列行列上に存在する配列モチーフ領域の長さ  $K'$  は一般

に  $K$  以上になってしまう。すなわち、長さ  $K$  の配列モチーフ領域において、アミノ酸の数よりもギャップの数が多く存在する列の数を  $K_g$  とすると、 $K' = K + K_g$  の関係が成立する。

#### 4.2 整列行列を用いた初期値の選択法

整列行列上のある特定領域（連続する  $K$  個の列）に配列モチーフが多く含まれることが経験的に知られている。これを踏まえ、多重整列化により得られる整列行列  $DS'$  から良好な初期値を探索するための新しい方法を提案する。以下では、ギャップが含まれる整列行列からプロファイル行列を計算する方法、初期値を見つけ出すためのスライディングウィンドウ法、スライディングウィンドウ法の精度向上に重要となる非配列モチーフ領域の判定法について述べる。

##### (1) プロファイル行列の計算法

多重整列化によって挿入されたギャップは、配列の長さを揃える為に挿入された空白である。このため、整列行列において、ギャップが混在する列に存在する文字だけを考慮し、文字の出現頻度を計算すると、他のギャップの少ない列に比べて、その値が大きくなり、あたかもモチーフ領域の一部と見做されてしまう。これを避ける為に、 $N \times L$  整列行列（多重整列）上に存在するある  $N \times K'$  配列領域の  $M \times K'$  プロファイル行列を算出するには、予めギャップを考慮した計算方法を定める必要がある ( $K' \leq L$ )。

著者らは、ギャップを考慮した計算方法として2つの方法を提案する。まず、整列行列  $DS'$  内に存在する各ギャップに対して、20種類のアミノ酸文字をランダムに割り当てる方法である。図5は各ギャップに文字をランダムに割り当てた例である。左側はギャップを含む  $5 \times 10$  整列行列  $DS'$  であり、右側は各ギャップをランダムに割り当てた結果として得られる行列  $DS''$  である。このような行列  $DS''$  を用いると  $M \times K'$  プロファイル行列  $PMAT=(p_{ij})$  を容易に計算することができる。以上により、ギャップが多く含まれている列にランダム性を持たせて、プロファイル行列を計算している。次に、すべての文字にギャップ文字の出現頻度を均等に分ける方法の均等割り付けである。図5の左図の6文字目を例に説明する。出現頻度は、文字Dが2、文字Gが1、そしてギャップ文字が2である。ギャップ文字の出現頻度2を20種類のアミノ酸に割り振ると、文字Dが2.1、文字Gが1.1、その他の文字が0.1となる。このように、出現頻度を均等に分けることで、ランダムに割り当てた時に偏りがないようにしている。

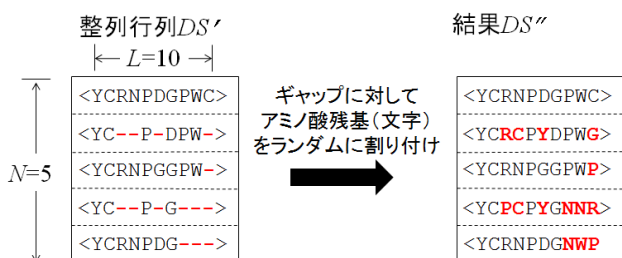


図5 ギャップに対してランダムに文字を割り当てた例

##### (2) スライディングウィンドウ法

スライディングウィンドウ法とは、整列行列の左端列から右端列へ向かってウィンドウをスライドさせることにより、配列モチーフに近いウィンドウ領域を見つけ出し、その領域をサイトサンプラーの初期値とする方法である。

ウィンドウの長さ  $K'$  はその位置によって異なり、可変である。すなわち、ウィンドウ内に存在する非ギャップ列の数  $K$  はウィンドウの位置に依存せず固定だが、ギャップ列の個数  $K_g$  はウィンドウの位置に依存する ( $K' = K + K_g$ )。

ウィンドウをスライドさせ、配列モチーフに近い領域を見つけ出す尺度として、相対エントロピー  $F$  を利用する。相対エントロピー  $F$  の計算には、クラスタごとに、プロファイル行列  $PMAT=(p_{ij})$  と背景頻度行列  $BMAT$  を利用している。

##### (3) 非配列モチーフ領域の判定法

PROSITE などのモチーフデータベースに登録されている配列モチーフの最左端または最右端に、ギャップは存在しない。このため、ウィンドウの最左端列または最右端列にギャップが存在している場合は、そのウィンドウに配列モチーフが多く含まれる可能性があるかどうか吟味する必要がある。

ウィンドウの最左端列または最右端列が、以下を満たすとき、そのウィンドウを非配列モチーフ領域と呼ぶ。

$$N_g \geq N \times \eta \quad (3)$$

ただし、 $N_g$  は列に存在するギャップ数、 $\eta$  は閾値を意味する。 $\eta$  は値をいくつか変化させてみたところ、0.1 が一番良い結果となったので、実験ではその値を採用する。

スライディングウィンドウ法では、非配列モチーフと判断されるウィンドウは、初期値の探索から外される。

#### 4.3 進化的な知識を導入した擬似度数の計算法

提案手法における出現頻度や相対エントロピーの計算で用いられるプロファイル行列では、従来のサイトサンプラーにはない進化的な知識を導入する。具体的には、プロファイル行列の定義に利用されている擬似度数にアミノ酸置換のし易さを数値化したアミノ酸置換行列（進化的な知識）を導入する。このために、式(1)のプロファイル行列  $PMAT=(p_{ij})$  を次のように定義する[45]。

$$p_{ij} = \frac{c_{ij} + b_{ij}}{N_j + B_j} \quad (4)$$

ただし、 $b_{ij}$  は以下で定義される擬似度数、 $N_j = \sum c_{ij} [1 \leq i \leq 20]$  を意味する。

$$b_{ij} = B_j \times \sum \left( \frac{c_{kj}}{N} \times \frac{g_{ki}}{g_i} \right) \quad [1 \leq k \leq 20] \quad (5)$$

$B_j$  は  $B_j = m \times R_j$  と定義されており、 $R_j$  を1クラスタにおける  $j$  列目の文字の種類数とする。 $m$  は試験的な検索実験により決定される正の数であり、 $m=5\sim6$  が最も有効であると報告されている。また、 $g_{ki}$  はアミノ酸  $k$  からアミノ酸  $i$  への置換頻度で、次式で表される。

$$g_{ki} = g_k \times g_i \times 2^{s(k,i)} \quad (6)$$

$s(k,i)$ はアミノ酸  $k$  からアミノ酸  $i$  への置換のし易さを数値化した類似スコアであり、この類似スコアは、BLOSUM62 と呼ばれるアミノ酸置換行列から取得している。 $g_i$  は  $DS'$  内の  $i$  の出現確率で、 $G_i$  は  $g_{ki}$  の文字ごとの総和であり、 $G_i = \sum g_{ki} [1 \leq k \leq 20]$  と定められている。

進化的な知識を考慮した式 (4) は、式 (2) で紹介した相対エントロピーの計算に利用している。また、相対エントロピーでは、モチーフ領域が本当に類似しているかが分かりにくいいため、評価には相対エントロピーをもとに計算する  $p$  値を用いる。 $p$  値の式は次のように定義される。

$$p - value = (n + 1)2^{-F} \quad (7)$$

#### 4.4 アルゴリズム

提案手法では配列データセットを多重整列化し、それにより得られる整列行列  $DS'$  に対してスライディングウィンドウ法を適用する。これにより探索されたクラスタをサイトサンプラーの初期値とする。次に、サイトサンプラーを実行することにより、配列データセットの各配列から配列モチーフにできるだけ近い  $K$ -類似部分配列を抽出する。提案手法のアルゴリズムは以下のとおりである。

【入力】 配列データセット  $DS$ 、配列モチーフの長さ  $K$   
 アミノ酸配列を表現する文字集合  $\Omega$ 、閾値  $\eta$ 、  
 アミノ酸置換行列  $s(k,i)$ 、 $k \in \Omega$ 、 $i \in \Omega$

【出力】 配列データセットの各配列から抽出される  $K$ -類似部分配列

- ① 分子進化系統樹を案内木として利用する多重整列化を配列データセットに適用し、整列行列  $DS'$  を求める。
- ② 4.2 節の(1)で述べたように、 $DS$  の表現に利用されている文字の集合から文字をランダムに選び、それを整列行列  $DS'$  内のギャップに割り当て  $DS'$  を得る。
- ③ 4.2 節の(2)で述べたように、 $DS'$  に対してスライディングウィンドウ法を適用する。すなわち、 $N \times K'$  のクラスタ (長さ  $K'$  の部分配列集合) を 1 文字ずつスライドさせ、最後のクラスタの処理が終わるまで、以下の処理を繰り返す。
  - ・4.2 節の(3)で述べた非配列モチーフ領域の判定法を用いて、当該クラスタが非配列モチーフ領域と判定されたときは、そのクラスタを候補配列集合から外す。
  - ・4.3 節の式(4)で述べた擬似度数を計算し、その結果を用いて、当該クラスタに対するプロファイル行列  $PMAT$  を算出する。
  - ・算出されたプロファイル行列  $PMAT$  を用いて、当該クラスタの  $p$  値を 4.3 節で述べた方法で計算する。
- ④  $p$  値が算出されたクラスタの集合から  $p$  値の値が最も小さいクラスタ ( $K$ -部分配列集合) をサイトサンプラーの初期値として選択する。
- ⑤ 初期値を用いて、3.4 節のサイトサンプラーを実行する。ただし、このサイトサンプラーで必要となるプロファイル行列  $PMAT$  の計算では、4.3 節の式(4)で述べた擬似度数を用いる。

#### 5. 評価実験

提案手法の有効性を確認するために、PROSITE と呼ばれるモチーフライブラリー (モチーフデータベース) [46] を使用した。モチーフライブラリーに含まれるアミノ酸配列データセットのそれぞれには、現在までに見つかっている配列モチーフとそれを含む配列データが数多く収集されている。評価実験では、モチーフライブラリーの中から 5 種類のアミノ酸配列データセット [47][48][49][50][51] を選んだ。これらのデータセットの選定に際しては、データ件数の多いものや少ないものを初めとして、配列長がほぼ同じものから大きく異なるものが含まれるように配慮した。

表 1 に 5 種類のアミノ酸配列データセットの概要を示す。表中のクラスタの長さ  $K'$  は、本来の配列モチーフ長  $K$  に、多重整列化によって挿入された配列モチーフ内のギャップ長  $Kg$  を加えたものを意味する。なお、予備実験として、閾値  $\eta$  をいくつか変化させてみたところ、0.1 が一番良い結果が出たので、評価実験では、その値を採用している。

#### 5.1 実験方法

多重整列化を行うプログラムとして、分子進化系統樹を案内木として利用する CLUSTAL X [23][24] を使用している。CLUSTAL X は、計算途中で配列間の位置関係が凍結されてしまうため、完全に多重整列化された結果が出力されないと言われている。これを改善するため、CLUSTAL X などにより計算出力された多重整列に対して、適当な場所にギャップをさらに追加する反復改善法 [52] が開発されている。しかし、著者らの予備実験では、CLUSTAL X の方が反復改善法よりも多重整列中にギャップ数が少なく、良好な初期値が得られている。このため、多重整列化のプログラムとして CLUSTAL X を利用している。CLUSTALX のパラメータは、初期パラメータを使用している。

以下では、提案手法に有用性を確かめるために導入した尺度として、計算精度について説明する。

表 1 実験に使用したデータセット

番号	モチーフ名	登録番号	クラスタの長さ ( $K' = K + Kg$ )	件数
1	Kringle	PS00021	14=14 + 0	95
2	Homeobox	PS00027	115=24 + 91	1321
3	PTS_EIIA	PS00372	22=17 + 5	51
4	HTH_ASNC	PS00519	37=27 + 10	43
5	HTH_DEOR	PS00894	35=35 + 0	82

従来手法と提案手法の実行において、処理の繰り返し回数は、千回とした。更に両手法の試行回数はそれぞれ 10 回とし、それぞれ抽出した結果の平均を精度としている。この精度の定義は、以下のとおりである。

$$\text{精度}(\%) = \frac{B}{B + C} \times 100 \quad (8)$$

ただし、 $B$  は従来手法あるいは提案手法を適用することにより抽出された配列モチーフ領域、 $C$  は非モチーフ領域を意味する。図 6 を例として挙げると、抽出した  $K$ -部分配列集合の全文字数を分母とし、その中でマッチした配列モチ

ーフ領域の全文字数を分子とする事で、 $K$ -類似部分配列集合内でマッチした配列モチーフ領域が全体の何割であるかを表すことが出来る。よって式(8)の数値が高いほど、抽出された  $K$ -部分配列集合は配列モチーフ領域と一致している部分が多く、低いほど、配列モチーフ領域から外れていると解釈できる。

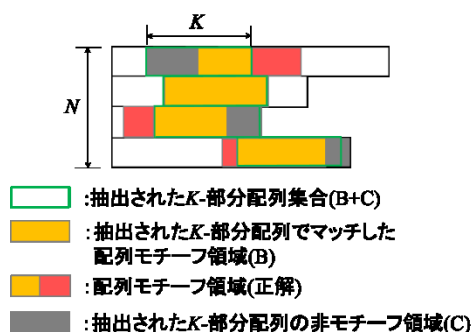


図 6  $K$ -部分配列集合内で判断する配列モチーフ領域

### 5.2 実験結果

表 2 に従来手法と提案手法であるランダム割り付けと均等割り付けの精度を示す。また、この表に、5つのデータセットのそれぞれに対する多重整列化の局在度も示す。

表 2 提案手法と従来手法との精度結果比較

番号	従来手法(%)	ランダム(%)	均等(%)
1	64.44	86.54	86.54
2	87.75	99.80	99.80
3	49.18	98.50	99.04
4	41.76	55.12	62.46
5	17.06	89.88	89.88

表 2 を見る限りでは提案方式が従来手法よりも優れている事、ランダム割り付けより均等割り付けの精度がよい事が分かる。しかし、4番のデータセットにおいては、他のデータより配列モチーフの精度が低い。

### 5.3 考察

提案手法が効果をあげた理由としては、従来手法には考慮されていない**進化的な知識**を、多重整列化や相対エントロピーの計算に用いているためである。また、初期値を探索した後に適用するサイトサンプラーに関しても**進化的な知識**を導入したプロファイルを計算しているので、抽出精度がより向上したと考えられる。

また、ランダム割り付けより均等割り付けの精度がよい理由として、均等に出現頻度を分けることでギャップ以外の文字の出現頻度の順位が変化しないからだと考えられる。

### 6. まとめ

本研究では、アミノ酸配列データセットから配列モチーフに相当する類似部分配列を高精度に抽出する方法を提

案した。解の品質を大きく左右する初期値をできるだけ良質なものにするために、ランダムに初期値を選択することを止め、**進化的な知識**が導入された多重整列化をアミノ酸配列データベースに適用し、これにより得られた整列行列から高い品質をもつ初期値を選択した。初期値の選択では、整列行列に対して、ウインドウをスライドさせるスライディングウインドウ法を適用している。スライディングウインドウ法で初期値を選択するとき、ウインドウ領域の両端の一方に、閾値  $\eta$  を超えるギャップ数がある場合、そのウインドウ領域を初期値の選択から除外した。なお、初期値としてのウインドウ領域を選択するために、ギャップ文字をランダム割り付けと均等割り付けの2つの方法を用い、 $p$  値の尺度を利用している。また、このような初期値の選択の他に、従来の GS アルゴリズムに含まれている(1)部分配列の出現頻度や(2)候補配列集合の相対エントロピーの計算に**進化的な知識**を導入した。

提案方法の評価実験を5つのデータセットを用いて実施した結果、従来方法に比べて高精度な配列モチーフとしての類似部分配列集合の抽出が可能になった。

今後の課題として、偽陽性を取り除くための方法の検討が挙げられる。例えば、それぞれの文字に対する背景頻度の計算に  $n$  次のマルコフ過程を導入する方法などがある。この他の課題としては、サイトサンプラーが前提としている配列モチーフの長さ  $K$  を遺伝的アルゴリズムや島モデル等により自動決定する方法の検討などがある。また、配列モチーフを抽出するだけでなく、テキストデータにおける規則的な共通部分の探索、Web 文章の履歴や顧客の購買履歴の分析の利用可能性についての検証は残されている。

### 参考文献

- [1] Charles E. Lawrence, Stephen F. Altschul, Mark S. Boguski, Jun S. Liu, Andrew F. Neuwald, and John C. Wootton: Detecting Subtle Sequence Signals: A Gibbs Sampling Strategy for Multiple Alignment, *Science*, Vol. 262, No. 513, pp.208-214, October 1993.
- [2] Liu Li-fang, Jiao Li-cheng: A Greedy Two-stage Gibbs Sampling Method for Motif Discovery in Biological Sequences, *Journal of Information Science and Engineering*, Vol.26, pp.2309-2318, 2010.
- [3] William Thompson, Eric C. Rouchka, and Charles E. Lawrence: "Gibbs Recursive Sampler: finding transcription factor binding sites," *Nucleic Acids Research*, Vol. 31, Issue 13, pp. 3580-3585, 2003.
- [4] The Gibbs Motif Sampler Homepage: <http://ccmbweb.ccv.brown.edu/gibbs/gibbs.html>
- [5] Roman A. Laskowski: "SURFNET: A program for visualizing molecular surfaces, cavities and intermolecular interactions," *Journal of Molecular Graphics*, Volume 13, Issue 5, pp.323-330, October 1995.
- [6] 欧州バイオインフォマティクス研究所: SURFNET, <http://www.ebi.ac.uk/thornton-srv/software/SURFNET/>
- [7] David Lee, Oliver Redfern, and Christine Orengo: Predicting protein function from sequence and structure, *Nature*



- Reviews *Molecular Cell Biology*, Vo.8, pp.995-1005, December 2007.
- [8] Stuart Geman and Donald Geman: "Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.6, No.6, pp.721-741, 1984.
- [9] Eric C. Rouchka: A Brief Overview of Gibbs Sampling, *Bioinformatics Technical Report Series*, No. TR-ULBL-2008-02, University of Louisville, 9 pages, March 24, 2008
- [10] Liu Li-fang, Jiao Li-cheng, and Huo Hong-wei: A Greedy Two-stage Gibbs Sampling Method for Motif Discovery in Biological Sequences, 2008 International Conference on BioMedical Engineering and Informatics, pp.13-17, 2008.
- [11] Hastings, W.K: Monte Carlo Sampling Methods Using Markov Chains and their Applications, *Biometrika*, Vol.57, pp.97-109, 1970.
- [12] 国友直人・山本拓[監修], 北川源四郎・竹村彰通[編]: 21世紀の統計科学, 第III巻 数理・計算の統計科学, 第10章 マルコフ連鎖モンテカルロ法入門, 東京大学出版会, 2008年.
- [13] Andrew F. Neuwald, Jun S. Liu, and Charles E. Lawrence: Gibbs motif sampling: Detection of bacterial outer membrane protein repeats, *Protein Science*, Cambridge University Press, Vol.4, pp.1618-1632, 1995.
- [14] 河野 修久, 加藤 智之, 田村 慶一, 北上 始: 配列データベースから類似部分配列を抽出するためのGS最適化手法に関する考察, 電子情報通信学会 第19回データ工学ワークショップ(DEWS2008), Online Proceedings, 8 pages, 2008年3月9日~3月11日.
- [15] Nobuhisa Kono, Hajime Kitakami, Keiichi Tamura, and Yasuma Mori: Extracting Similar Subsequences by Gibbs Sampling with Distributed MGG, *Proceedings of the 2009 International Conference on Parallel & Distributed Processing Techniques & Applications (PDPTA'09)*, pp.669-675, Las Vegas in USA, July 13-16 in 2009.
- [16] Jun S. Liu: "The Collapsed Gibbs Sampler in Bayesian Computations with Applications to a Gene Regulation Problem," *Journal of the American Statistical Association*, Vol.89, No.427, pp.958-966, September 1994.
- [17] Liu Li-fang, Jiao Li-cheng: "Motif GibbsGA: Sampling Transcription Factor Binding Sites Coupled with PSFM Optimization by Genetic Algorithm, *Journal of Convergence Information Technology*," Vol.5, No10, pp.141-148, 2010.
- [18] William A. Thompson, Lee A. Newberg, Sean Conlan, Lee Ann McCue, and Charles E. Lawrence: The Gibbs Centroid Sampler, *Nucleic Acids Res*, Vol.35, Issue suppl 2, pp.W232-W237, 2007.
- [19] Neil C. Jones and Pavel A. Pevzner: An Introduction to Bioinformatics Algorithms, The MIT Press, 2004.
- [20] Mark P. Styczynski, Kyle L. Jensen, Isidore Rigoutsos, and Gregory Stephanopoulos: "BLOSUM62 miscalculations improve search performance," *Nature Biotechnology*, Vol.26, No.3, pp.274-275, 2008.
- [21] Julie D. Thompson, Toby J. Gibson, Frédéric Plewniak, François Jeanmougin, and Desmond G. Higgins: "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice," *Nucleic Acids Research*, Oxford University Press, Vol.22, No.22, pp.4673-80, November 11 1994.
- [22] M. A. Larkin, G. Blackshields, N. P. Brown, R. Chenna, P. A. McGettigan, H. McWilliam, F. Valentin, I. M. Wallace, A. Wilm, R. Lopez, J. D. Thompson, T. J. Gibson and D. G. Higgins: Clustal W and Clustal X version 2.0., *Bioinformatics*, Vol. 23, Issue 21, pp.2947-2948, 2007.
- [23] Steven Henikoff and Jorja G. Henikoff: Amino Acid Substitution Matrices from Protein Blocks, *Proceedings of Natural Academy of Science of the United States of America*, Vol.89, pp.10915-10919, November 1992.
- [24] Jorja G. Henikoff and Steven Henikoff: Using substitution probabilities to improve position-specific scoring matrices, *Computer Applications in the Biosciences*, Vol.12, No.2, pp.135-143, April 1996.
- [25] Naruya Saitou and Masatoshi Ne: "The Neighbor-joining Method: A New Method for Reconstructing Phylogenetic Trees," *Molecular Biology and Evolution*, Volume 4, Issue 4, pp.406-425, 1987.
- [26] 福本 翔平, 北上 始, 森 康真: 多重整列に基づくモチーフの統計的抽出法 電子情報通信学会 第13回情報科学技術フォーラム (FIT2014) 論文集, D-008, Online Proceeding, 2014.
- [27] Steven Henikoff and Jorja G. Henikoff: "Amino Acid Substitution Matrices from Protein Blocks," *Proceedings of the National Academy of Sciences of the United States of America (PNAS)*, Vol.89, No.22, pp.10915-10919, 1992.
- [28] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, Jianyong Wang, Helen Pinto, Qiming Chen, Umeshwar Dayal, and Mei-Chun Hsu: "Mining Sequential Patterns by Pattern-Growth: The Prefix Span Approach," *IEEE Transaction on Knowledge and Data Engineering*, Vol. 16, No. 11, pp.1424-1440, November 2004.
- [29] 加藤 智之, 北上 始, 森 康真, 田村 慶一, 黒木 進: 極小かつ非冗長な可変長ワイルドカード領域を持つ頻出パターンの抽出, 電子情報通信学会和文論文誌D「データ工学特集号」, Vol.J90-D, No.2, pp.281-291, 2007年2月.
- [30] 加藤 智之, 森 康真, 黒木 進, 北上 始: 可変長配列パターン抽出法におけるギブスサンプリングを用いた不要パターンの除去方式, 日本データベース学会論文誌(DBSJ Letters), Vol.6, No.1, pp.65-68, 2007年6月.
- [31] Richard Durbin, Sean R. Eddy, Anders Krogh, and Graeme Mitchison: Chapter5 Profile HMMs for sequence families, "Biological sequence analysis - Probabilistic models of proteins and nucleic acids," Cambridge University Press, pp.100-133 1998.
- [32] Lee A. Newberg et al.: A phylogenetic Gibbs sampler that yields

- centroid solutions for *cis*-regulatory site prediction, *Bioinformatics*, Vol.23, No.14, pp.1718–1727, 15 July 2007.
- [33] William A. Thompson et al. : Using the Gibbs Motif Sampler for phylogenetic footprinting, *Methods Mol. Biol.*, Vol 395, pp.403-424. 2007.
- [34] 北上 始, 斎藤成也, 太田聡史: ビッグデータ時代のゲノミクス情報処理, コロナ社, 2014年10月.
- [35] Anders Krogh, Michael Brown, I. Saira Mian, Kiminen Sjolander, and David Hausder: Hidden Markov Models in Computational Biology Applications to Protein Modeling, *Journal of Molecular Biology*, Vol.235, pp.1501-1531, 1994.
- [36] Sean R. Eddy: Multiple alignment using hidden Markov models, *Proceedings of International Conference on Intelligent Systems for Molecular Biology (ISMB-95)*, AAAI/MIT Press, pp.114-120 (1995).
- [37] Richard Hughey and Anders Krogh: "Hidden Markov models for sequence analysis: extension and analysis of the basic method," *Computer Applications in the Biosciences (CABIOS)*, Vol.12, No.2, pp.95-107, 1996.
- [38] David E. Goldberg: *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley Longman Publishing Co., Inc. Boston, MA, USA, 1989.
- [39] Kazuhito Shida: "GibbsST: a Gibbs sampling method for motif discovery with enhanced resistance to local optima," *BMC Bioinformatics*, Supplement 5, Vol.7, p.486, November 2006.
- [40] Kazuki Shida: "Hybrid Gibbs-Sampling Algorithm for Challenging Motif Discovery: GibbsDST," *Genome Informatics*, Vol.17, No.2, pp.3-13, 2006.
- [41] Scott Kirkpatrick, C. Daniel Gelatt and Mario P. Vecchi: Optimization by simulated annealing, *Science*, Vol.220, pp.671-680, 1983.
- [42] Richard Bellman: *Dynamic Programming*, Princeton University Press, 1957.
- [43] Dan Gusfield: "Algorithms on Strings, Tree, and Sequences: Computer Science and Computational Biology, Cambridge University Press, 1977.
- [44] Needleman, B. Saul and Wunsch, D. Christian: "A general method applicable to the search for similarities in the amino acid sequence of two proteins," *Journal of Molecular Biology*, Vol.48, Issue.3, pp.443–53, March 1970.
- [45] 福本 翔平, 北上 始, 森 康真: アラインメントされた配列集合からモチーフを抽出する方法, 電子情報通信学会 第5回データ工学と情報マネジメントに関するフォーラム (DEIM2013) 論文集, E5-2, Online Proceedings, 2013.
- [46] PROSITE : <http://prosite.expasy.org/>
- [47] Kazuho Ikeo, Kei Takahashi, and Takashi Gojobori: Evolutionary origin of numerous kringles in human and simian apolipoprotein(a), *Federaton of European Biochemical Societies*, Vol.287, No.1-2, 146-148, 1991.
- [48] Walter J. Gehring, Yasushi Hiromi: Homeotic genes and the homeobox, *Annu. Rev. Genet.*, Vol.20, pp.147-173, 1986.
- [49] Pieter W. Postma, Joseph W. Lengeler, Gary R. Jacobson: Phosphoenolpyruvate: carbohydrate phosphotransferase systems of bacteria, *Microbiology and Molecular Biology Reviews*, vol.57, No.3, pp.543-594, 1993.
- [50] Debra Aker Willins, Christopher W. Ryan, Jill V. Platko, Joseph M. Calvo: Characterization of Lrp, and Escherichia coli regulatory protein that mediates a global response to leucine, *Journal of Biological Chemistry*, Vol.266, No.17, pp.10768-10774, 1991.
- [51] Susanne Beck von Bodman, G.Tomas Hayman, and Stephen K. Farrand: Opine catabolism and conjugal transfer of the nopaline Ti plasmid pTiC58 are coordinately regulated by a single repressor, *Proceedings of the National Academy of Sciences of the United States of America*, Vol.89, pp.643-647, 1992.
- [52] Kazutaka Katoh and Daron M. Standley: MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability, *Molecular Biology and Evolution*, Volume 30, Issue 4, pp.772-780, 2013.