

TR-0762

配列データの多重アライメント法

石川 幹人、金久 實 (京大)

April, 1992

© 1992, ICOT

**ICOT**

Mita Kokusai Bldg. 21F  
4-28 Mita 1-Chome  
Minato-ku Tokyo 108 Japan

(03)3456-3191 ~ 5  
Telex ICOT J32964

---

**Institute for New Generation Computer Technology**

## 第21章

### 配列データの多重アライメント法 石川 幹人, 金久 寛

分子進化を塩基配列やアミノ酸配列のデータに基づいて解析するには、配列間の塩基置換数やアミノ酸置換数を推定し、その情報から系統樹を作成するのが一般的である。置換数を推定するには、配列間で置換が発生したと想定される部分の対応づけがなされている必要がある。配列間でこの対応づけを行うことをアライメントというが、精度の向上のために、複数の配列データの同時比較が良く行われており、それをとくに多重アライメントと呼ぶ。多重アライメントの結果では、置換が発生したと想定される配列の類似部分が、互いに縦に並べ合わされて示される(図21・9参照)。また、多重アライメント法を用いると、配列データの特徴的部分(配列モチーフ)が抽出でき、系統樹の作成だけでなく、配列の構造や機能を推定するのにも有効である。

多重アライメントの問題は、与えられた評価尺度のもとで最適の解が原理的に計算できるという意味では解決されている。しかし、後で述べるように、この計算量は極めて大きく、現存する最高速の計算機をもってしても対処しきれない。そのうえ、評価尺度についても完全なものを決定できていないことから、現実的な意味で、多重アライメントの問題は解決されているとは言えない。従来、多重アライメントは熟練した研究者が人手で行うことが多く、むしろ、そのほうが結果の品質が良かった。だが、塩基配列やアミノ酸配列の決定法が確立されて、その配列データベースが急速に膨らんでいる今日、多重アライメントの計算機による自動化の要求は強まっている。さらに、最近の計算機技術の向上は、品質の良い多重アライメントを実用的な時間内に算出する可能性を期待させる。

そこで、この章では、計算機を用いて多重アライメントを求める方法を5手法解説する。この分野はまだ研究途上であり、どの手法も一長一短ではあるが、それらの特徴をうまく捉えて活用すれば、多重アライメントの要求の多くに応えられるものと考えている。なお、以下の解説は、アミノ酸の配列データを扱う場合に限って話を進める。一般に、分子進化の解析には、アミノ酸配列を用いるほうが多い。それは、復帰突然変異や並行突然変異が起こる確率が、塩基サイトよりもアミノ酸サイトのほうが低く、分子進化の解析が容易になるからである。とはいえ、塩基配列を扱う場合も原理的にはアミノ酸の配列を扱う場合と同様であるので、読者はここで解説する多重アライメント法を、容易に塩基配列へ応用できるであろう。

#### 21・1 セグメント法

##### 21・1・1 概要

セグメント法とは、図21・1に示すように、複数の配列にわたって共通した、あるいは類似した部分配列(セグメント)を見つけだし、その部分をピンで止め合わせるように縦に揃えることで、多重アライメントを作る方法である。人間が手作業で多重アライメントを作る場合も、多かれ少なかれ、このセグメント法を行っているように思われる。この方法は、配列群の

なかに部分的に類似性の高い領域が存在するときには、その領域を貫くピン止めがすぐ発見でき、極めて有効である。しかし、その反面、類似性が低い領域を処理しようとする、多数のピン止めが複雑に交差しあい、どれを優先すべきかという問題が発生する。

## 2.1.1.2 方法

セグメント法による  $n$  本の配列 (配列 1 ~ 配列  $n$ ) を貫くピン止めの発見手順は、次の通り。  
(1) 配列 1 の先頭のアミノ酸から始めて、指定された長さ (たとえば 20 個) のアミノ酸を連続して取り出し、ひとつのセグメントとする。同様に先頭から 2 番目のアミノ酸から、やはり指定された長さ 20 のセグメントをとる。それを配列 (長さ  $m$  とする) の終端まで行くと全部で  $m - 19$  個のセグメントがとれる。

(2) 配列 2 の先頭のアミノ酸から、長さ 20 のセグメントをとり、それと配列 1 から得られたセグメントとをひとつずつ比較し、同じ位置に同一のアミノ酸がいくつ対応しているかを調べる。指定された割合以上の個数 (たとえば割合 70 % が指定されていると、20 個に対して 14 個) が対応していればマッチングであり、配列 1 と配列 2 が、そのセグメント対において、ピン止めされるとみなされる。一方、 $m - 19$  個のセグメントのいずれともマッチングしない場合は、この配列 2 の先頭からとられたセグメントは捨ててしまう。(3) 配列 2 の先頭から 2 番目のアミノ酸から長さ 20 のセグメントをとり、(2) と同様の処理をする。順次、配列 2 の末端まで同様な処理をする。それによって、配列 1 と配列 2 を貫くピン止めに相当するセグメント対が、複数得られる。

(4) 配列 3 の先頭のアミノ酸から、長さ 20 のセグメントをとり、それと先に得られたセグメント対とをひとつずつ比較する。セグメント対のどちらか一方とマッチングしていれば、配列 1 と配列 2 と配列 3 が、そのセグメントの 3 つ組において、ピン止めされるとみなされる。一方、配列 1 と配列 2 を貫くピン止めに相当するセグメント対のいずれともマッチングしない場合は、この配列 3 の先頭のアミノ酸からとられたセグメントは捨ててしまう。(5) 配列 3 の先頭から 2 番目のアミノ酸から長さ 20 のセグメントをとり、(4) と同様の処理をする。順次、配列 3 の末端まで同様な処理をする。それによって、配列 1 と配列 2 と配列 3 を貫くピン止めに相当するセグメントの 3 つ組が、複数得られる。

(6) 同様な処理を配列 4 から配列  $n$  まで行い、配列 1 から配列  $n$  までを貫くピン止めに相当するセグメントの  $n$  個組が、複数得られる。(この結果は、配列を処理する順番によって異なる。その意味で、この手順はヒューリスティクスである。すべての組合せを調べてピン止めを発見するのは、組合せ数が多く現実的でないため、何らかのヒューリスティクスを用いた手順にならざるをえない<sup>1)</sup>。)

配列 1 から配列  $n$  までを貫くピン止めが見つかったならば、強いピン止めから順に (各セグメント間で対応したアミノ酸の総個数が多い順に)、実際の配列群をピン止めしていく。その際に、すでにピン止めされている部分と交差するピン止めが出てきた場合は、そのピン止めは捨ててしまう。すべてのピン止めが終わったならば、多重アライメントの完成である。なお、どのピン止めにも属さない配列の部分は、ピン止めとピン止めの間に、長さ合わせのためのギャップとともに配置される。

<sup>1)</sup>D. J. Bacon, W. F. Anderson, *J. Mol. Biol.*, 191, 153-161 (1986).

## 2 1・1・3 検討

$n$ 本の配列を貫くピン止めを見つけることは、結果的に $n$ 次元のホモロジーマトリックスを求めることに相当している。マトリックス内の各ドットが、ひとつのピン止めに対応していると考えられる。図2 1・2に2次元のホモロジーマトリックスの例を示す。対角線の方に線が見える領域は、その2本の配列の類似性が高い部分に相当する。このように2本の配列の類似性が視覚的にわかりやすいため、人手で行う多重アライメントの補助として、よく使用される表現である<sup>2</sup>。このホモロジーマトリックスを見ながら、セグメント法にまつわる問題点を理解しよう。

類似性の低い配列群の多重アライメントをセグメント法で行う際、ピン止めの発見手順において、マッチング判定基準の割合を下げなければ、ピン止めが発見されないことが良くある。そんなとき、セグメントを長くし、判定基準を下げていくと、今度は、途端に発見されるピン止めが多量になるという事態に陥る。ホモロジーマトリックスでいうと、マトリックス上に見られる対角線方向の線が途端に増え、それらが互いに重なってしまうのである。すると、配列の特定の部分が、他の配列の部分と複数箇所にわたって対応し、複数のピン止めが複雑に交差や融合を起こしていることとなる。こうなるとピン止めの優先順位を慎重に選ばないと、得られるアライメントの信頼性は低いものになってしまう。

今のところ、この優先づけの決め手は得られていない。熟練した研究者もセグメント法を用いているとすると、彼らは生物学上の知識と経験に裏付けられたノウハウとで、こうした優先づけを効果的に解決しているように思える。熟練した研究者が多重アライメントを解く過程を計算機上を実現するには、現在確立されつつある知識情報処理の手法が、今後有望であろう<sup>3</sup>。その第一歩として、すでに知られているモチーフを、2 6・3にあるようなデータベースから参照し、モチーフと類似したセグメントのピン止めは、その優先順位を上げるという方法を活用できる。

セグメント法を利用する別な試みには、複数の配列のなかから、はっきりと類似した部分だけを発見する前処理に、セグメント法を限定して使うものがある。Argosら<sup>4</sup>は、二つの連続したアミノ酸が同一の部分（セグメント長：2、判定基準100%に相当する）に注目してピン止めし、ピン止めの間の部分のアライメントには、後に述べるような他の手法を適用している。また並列計算機の応用には、多重アライメントを並列に行うために、ピン止めによって問題を分割するという発想が有効である<sup>5</sup>。

## 2 1・2 単純組合せ法

### 2 1・2・1 概要

単純組合せ法は図2 1・3に示すように、配列1と配列2、配列2と配列3、という具合に、配列を2本ずつペアワイズにアライメントし、その結果を次々と組合せて多重アライメントを

<sup>2</sup>宮田隆，林田秀宜，菊野玲子，安永照雄，“統生化学実験講座第1巻遺伝子研究法I”，p.381，東京化学同人（1986）。

<sup>3</sup>金久實，新田克己，小長谷明彦，田中秀俊，人工知能学会誌，6，630-640（1991）。

<sup>4</sup>M. Vingron, P. Argos, *CABIOS*, 5, 115-121 (1989).

<sup>5</sup>Butler, Foster, Karonis, Olson, Overbeek, Pflunger, Price, Tuecke, "Strand: New Concepts in Parallel Programming", ed. by Foster and Taylor, p.253, Prentice-Hall, (1990).

作する方法である。ペアワイズのアライメントには、通常ダイナミックプログラミングが用いられる。組合せ法は画一的で高速ではあるが、素朴に適用してしまうと、結果の多重アライメントの品質はあまり良くない。そのため、類似配列から順に組合せるとか、類似配列ごとにクラスターに分けて各々行うとかの改良が試みられている。

## 21・2・2 方法

単純組合せ法は、ペアワイズのアライメントを求める技術がわかれば、あとは明白である。ペアワイズアライメントの問題は、ダイナミックプログラミング（以下 DP と略す）の技術で、与えられた評価基準に沿った最適解が高速に求められる。DP は、最適化問題の解法として古くから知られていた技術であるが、生物の配列マッチングに応用されたのも随分前のことである<sup>6</sup>。

ペアワイズのアライメントについて DP の概念的説明を、図 21・4 を用いて行おう。たとえば、ADHE、AHIE という 2 つの短い配列をアライメントする場合、この 2 つの配列を図のような 2 次元のネットワークの辺に対応させる。斜め方向のアーキ（矢）は、そのアーキの位置に対応する 2 つのアミノ酸の類似度が割り振られる。また、縦および横方向のアーキは、ギャップに対応し、ギャップを挿入するときのコストが割り振られる。このように問題を定式化すると、最適なアライメントを求めることは、このネットワーク上の最良の経路を求めることに対応する。たとえば、太いアーキで表された経路が最良となったとする。この太いアーキで表された経路を順に見ていくと、A と A が対応し、D に対応するもう片方のアミノ酸はなく（つまりギャップが対応し）、H には H が対応し、という具合に解釈することができる。結果として図の右側にあるアライメントが得られる。

最良の経路は、具体的には、左上の端から右下の端に向かって、各ノードに至る最良経路を段階的に決定していくことにより求めることができる。各ノードの計算を行うためには、直前の段階の各ノード（ペアワイズ DP では 3 つある）で求めた、左上端からそこまでの最良経路の評価値を参照して、今求めたいノードに至る評価値をそれぞれ計算する。そして、それらの最良値を求めて、それをそのノードまでの評価値とすればよい。直前のどのノードを選択したかという情報も記憶しておく。この操作を右下のノードまで繰返し、最後に逆向きに、選択したノードをたどれば、ネットワーク全体の最良経路を求めることができる。また、ネットワークのなかの部分的な最良経路を求め、部分的な類似性に注目したアライメントを行う DP の改良手法も開発されている<sup>7</sup>。

分子進化を考えると挿入や欠失がまとめて起こることから、ギャップコストはギャップの長さに依存した値にすることが望ましく、DP にも通常そうしたギャップコストが適用される<sup>8</sup>。ギャップコストに伴う計算効率を考えると、ギャップの長さ  $k$  に対して、 $a + kb$  のような一次式を与えるのが妥当である<sup>9</sup>。図 21・5 に、ギャップコストが  $a + kb$  のときの、各ノードの処理を図示した。  $a + kb$  の処理を実現するには、3 方向に異なる累積評価値  $D_{i,j}$  を送ると考えるとよい。斜め方向は、 $i$  番目のアミノ酸と  $j$  番目のアミノ酸の類似度  $W_{i,j}$  を最良の経路の

<sup>6</sup>S. B. Needleman, C. D. Wunsch, *J. Mol. Biol.*, 48, 443-453 (1970).

<sup>7</sup>W. B. Goad, M. I. Kanehisa, *Nucleic Acids Res.*, 10, 247-263 (1982).

<sup>8</sup>T. F. Smith, M. F. Waterman, *J. Mol. Biol.*, 147, 195-197 (1981).

<sup>9</sup>O. Gotoh, *J. Mol. Biol.*, 162, 705-708 (1982).

評価値に加えるだけでよいが、縦や横方向は、注目ギャップが第一ギャップである場合、コスト  $a$  を加味して最良経路を判定したうえで、ギャップコスト  $b$  を加えねばならない。

類似度には、Dayhoff マトリックス<sup>10</sup>を用いるのが最も一般的である。この尺度は、当時知られていたアライメントをもれなく調べ、同じ列に該当アミノ酸対が並んでいることが偶然に対してどの程度大きいかを数値化したものである(図 2 1・6)。数値は確率の対数値になっているため、それらの足し算は複合事象の共起確率を算出したことに相当する。近年、新たに知られたアライメントも増えてきたので、マトリックスの更新が叫ばれているが、それはまだ実現されてない。歴史的には、同じ Dayhoff マトリックスでも、図 2 1・6 の各値を 10 分の 1 にしたものが使われていたが、ギャップコスト  $b$  を整数のまま小さな値にできることから、図 2 1・6 の使用を薦めたい。

Dayhoff マトリックスには、ギャップコストについての指針はないが、図 2 1・6 を使用するときは、 $a = -60$ ,  $b = -5$  くらいが目安になる。もちろん、これらのパラメータを変えることによって、ギャップの入り方が変わるのであるから、問題に応じたパラメータ調整が必要となる<sup>11</sup>。アライメントの結果の両端に存在するギャップ(アウトギャップ)は、配列の内部に入るギャップとは異なったコストを与えられるようにするとなお良い。配列の類似部分が水平方向に大きくズレている配列対を効果的にマッチングするには、アウトギャップをゼロ(Dayhoff マトリックスにおける中立的な値)にしておくとい良いだろう。

以上、ペアワイズ DP について説明してきた。ペアワイズ DP が実装できれば、図 2 1・3 に示したように、配列をペアワイズにアライメントし、その結果を次々と組合せて多重アライメントを作るのは容易である。

### 2 1・2・3 検討

単純組合せ法は、DP を行うときに一度に比較する配列が 2 本であるため、一緒に比較していない配列同士の類似部分が、組合せたときにズレてしまう問題点がある。図 2 1・3 をみると、配列 2 と配列 4 とに同一な HIRA という部分があるが、配列 2 と配列 4 とは同時には比較されないので、多重アライメントにおいて RA の部分が同一カラムに揃わないという現象が起きている。処理する配列群の類似性が低いときには、とくにこの現象が顕著に起きる。これを防ぐには 2 つのアプローチがあるが、ひとつは、一度に比較する配列を増やすことであり、もうひとつは、より類似した配列から順に組合せることである。

比較する配列を増やすことは、それほど容易ではない。DP は、原理的には、一度に何本の配列でも同時に処理でき、与えられた評価値における最適な多重アライメントが得られるはずである。ところが  $n$  本の配列を同時にアライメントする  $n$  次元の DP は、概して、配列の  $n$  乗の計算量と  $n$  乗のメモリー量が必要であり、現実的に可能なのは 3 次元までである<sup>12</sup>。アライメントに大幅なギャップが入ることはあまりないことから、ネットワークの対角部分(図 2 1・4 では、右上隅と左下隅の部分)をカットすると処理をやや削減できる。それにより、問題に

<sup>10</sup>R. M. Schwartz, M. O. Dayhoff, "Atlas of Protein Sequence and Structure 5:3", p.353, Nat. Biomed. Res. Found., Washington, D. C. (1978).

<sup>11</sup>W. M. Fitch, T. F. Smith, *Proc. Natl. Acad. Sci. USA*, 80, 1382-1386 (1983).

<sup>12</sup>M. Murata, *Proc. Natl. Acad. Sci. USA*, 82, 3073-3077 (1985).

よっては4次元以上のDPが可能となる<sup>13</sup>。しかし、3次元以上のDPによって得られたアライメントを互いに組合せて、品質の良い多重アライメントを作るには複雑な処理が必要となる<sup>14</sup>。

一方、類似した配列から順に組合せれば、組合せ法によって得られる多重アライメントの精度を大きく改善できる。なぜなら、類似した配列同士のアライメントは確実で、信頼性が高いからである。配列の類似性を調べるのには、通常、すべての配列対についてペアワイズDPを行い、得られたペアワイズアライメントの評価値を互いに比べる方法がとられている。評価値を基準に類似性の高い配列群ごとにクラスターに分けて、クラスター内でそれぞれ組合せたのちに、クラスター同士を組合せる試みもなされている<sup>15</sup>。こうした発想は、ツリーベース法へとつながっていく。

## 21.3 ツリーベース法

### 21.3.1 概要

ツリーベース法は、配列間の類似性に従って描かれたツリー状の階層関係に基づいて、類似性の高い配列から順に多重アライメントを形成していく手法である<sup>16</sup>。その意味で、単純組合せ法の組合せ順をツリー状に拡張したものと捉えることができるが、ここでは、要素技術にプロファイルDPを導入したものを、とくにツリーベース法として解説する。ツリーベース法は、類似性の比較的低い配列群でも、ある程度実用的な多重アライメントが高速に得られる。しかし、単純組合せ法と同様に、初期の段階で起こったアライメントの誤りが、後々まで波及して、最後に得られる多重アライメント全体の精度を落とす危険性がある。

### 21.3.2 方法

最初にツリーの作成法を説明しよう。すべての配列対についてペアワイズDPを行い、得られたペアワイズアライメントの評価値をそれぞれの配列間の類似性とする。まず、最も類似性の高い配列対をツリーの枝として結線する。その結合した配列対の、残りの配列に対する類似性は、先の評価値の平均値とする。そしてまた、最も類似性の高い配列対を結線する。これを繰返すとツリーが完成する。(この手順は、24.1の平均距離法と同様であるから、詳しくはそこを参照されたい。)

厳密には、ペアワイズアライメントの評価値は、比較する配列対の長さやアミノ酸の構成比に依存する部分があるので、統計的な規格化を行う必要がある。しかし、その規格化された評価値を求めるには、ランダムにシャッフルした配列対の平均評価値を求めるため、ペアワイズアライメントを100回ほど行わねばならない。ツリーベース法に規格化された評価値を使用するのは、コスト高の割に十分な見返りがない<sup>17</sup>。むしろ、もっと実践的な方法で、評価値の曖昧性を吸収するのが良い。たとえば、類似性の高い配列対がいくつかあるときには該当ツリー

<sup>13</sup>H. Carrillo, D. Lipman, *SIAM J. Appl. Math.*, 48, 1073-1082 (1988).

<sup>14</sup>広沢誠, 星田昌紀, 石川幹人, 戸谷智之, 第2回公開ワークショップ「ヒトゲノム計画と情報解析技術」論文集, 120-123 (1991).

<sup>15</sup>W. R. Taylor, *CABIOS*, 3, 81- (1987).

<sup>16</sup>D. Feng, R. F. Doolittle, *J. Mol. Evol.*, 25, 351-360 (1987).

<sup>17</sup>J. G. Barton, M. J. E. Sternberg, *J. Mol. Biol.*, 198, 327-337 (1987).

を複数作っておいて、多重アライメントをする際にすべての場合を行い、その評価値の一番良い多重アライメントを採用する方法が効果的である。

ツリーに沿った多重アライメントには、プロファイル DP を用いる。図 2 1・7 に示すように、プロファイル DP は配列群 A と配列群 B とを、配列群 A, B 内の各々のアライメントはくずさずに、ペアワイズ DP する。プロファイル DP を行うときには、各アークに割当てられるアミノ酸の類似度が、そのアークの位置に対応する配列群 A のアミノ酸と、配列群 B のアミノ酸との類似度の総和になる。配列群 A が 2 本で、配列群 B が 3 本のときは、6 通りのアミノ酸対に関する類似度の和をとることになる。類似度にはペアワイズ DP と同様に、Dayhoff マトリックスを使用するとよいだろう。

ギャップコストの処理には注意を要する。ペアワイズ DP と違って、斜め方向のアークにもギャップに対応することがあるので、注目しているギャップがどの種類のギャップかが、ノードをさかのぼらないと判断できない。図 2 1・8 に示すように、配列対でギャップ同士に対応しているカラムは、普通、ギャップとはみなさない（無視する）ので、注目しているギャップが第 1 ギャップ ( $a+b$  を加える) か、第 2 ギャップ以降 ( $b$  のみを加える) かは、ネットワークの経路を戻り、対応関係を調べてから判断する。

プロファイル DP が実現できれば、ツリーに従って多重アライメントを行うことは容易である。ツリーの枝先から幹に向かって、枝が合わさる部分ごとに、ふたつの枝に相当するふたつの配列群をプロファイル DP して、多重アライメントを作っていく。そうすると最後には、全体の多重アライメントが得られる。プロファイル DP に際して、DP される配列群のなかに、近い配列対と遠い配列対とが混在するようになる（ツリーの幹に近づくにつれて、その傾向が出てくると）、アライメント結果に、近い配列対の影響が相対的に大きく現れてくる。それを防ぐには、プロファイル DP を行うときの配列群内の各配列に対し、近い配列同士はまとめて、重みづけを軽くするとよい<sup>4</sup>。

### 2 1・3・3 検討

この第 2 1 章は、分子進化解析に必要な系統樹を作成する前段階として、配列の多重アライメントを求める方法を解説している。にもかかわらず、多重アライメントを求めるのにツリーを描くとは、一見矛盾しているようにも感じられる。理想的には、多重アライメントと系統樹が同時に得られるのがよいのであろうが、それを厳密に実現するには、配列ごとの重みづけの任意性が発生するので、 $n$  次元 DP よりもさらに膨大な処理が必要である<sup>18</sup>。その大幅な近似解法も提案されてはいる<sup>19</sup>。とはいえ、ツリーを暫定的に設定すれば、多重アライメントの品質が上がり、それにつれて系統樹の精度が上がるのであるから、積極的に利用するに越したことはない。

ツリーベース法の問題点のひとつは、初期の段階で起こったアライメントの誤りが、後々まで影響し、結果の多重アライメントの品質が悪化するおそれがある点である。そうした影響を緩和するための循環的修正手法も提案されている。そのひとつは、ツリーベース法で作成した多重アライメントについて、対応する系統樹を作成し、その系統樹に沿ってツリーベース法を行って得られた結果について、また系統樹を作成するというサイクルを、もはや系統樹が更新

<sup>18</sup>D. Sankoff, *SIAM J. Appl. Math.*, 28, 35-42 (1975).

<sup>19</sup>J. Hein, "Methods in Enzymology", ed. by R. F. Doolittle, Vol.183, p.626, Academic Press (1990).

されなくなるまで続ける方法である。Kawaiら<sup>20</sup>は、この手法を並列計算機上を実現したうえで、アミノ酸配列の多重アライメントデータベースを作成した。

もうひとつの循環的修正手法は、Bartonら<sup>17</sup>の次の方法である。ツリーベース法で完成した多重アライメントから、始めの配列を抜きだし、それと残りの全部とでプロファイルDPを行い、多重アライメントを更新する。その結果から、今度は2番目の配列を抜きだし、それと残りの全部とで再びプロファイルDPを行う。その操作をすべての配列にわたって行ったならば、それを1サイクルとする。このサイクルを、もはや何の更新もされない回数まで繰返す。この修正手法は、次節で解説する反復改善法の発想のもとともなっている。

## 21.4 反復改善法

### 21.4.1 概要

ツリーベース法で得られた多重アライメントは、それをいかに修正しようとも、最初に作成したツリーの品質に大きく依存する。そこで、ツリーを始めに設定することなく、ツリーベース法と同等以上の多重アライメントを得ようとするのが、反復改善法<sup>21</sup>である。反復改善法は、内側にギャップの入っていない初期状態から、配列群をランダムにふたつのグループに分け、そのグループ同士をプロファイルDPにてアライメントし、得られた結果に対し、また配列群をランダムに分けてプロファイルDPを行うことを繰返すものである。全体の多重アライメントの評価値は繰返しごとに向上するので、乱数の系列にもよるが、ツリーベース法より良い評価値を与えることが多い。しかし、問題の規模が大きいと、収束までに多くの時間を必要とする。

### 21.4.2 方法

反復改善法の基本的アイデアは次のところにある。多重アライメント全体の評価値が、それに含まれるすべての配列対の評価値の和で与えられるという条件のもとでは、ある多重アライメントを任意にふたつの配列群に分け、互いにプロファイルDPを行うと、得られる多重アライメントは、もとの多重アライメントより、評価値がよい(少なくとも等しい)。つまり、プロファイルDPを繰返すと、多重アライメントの評価値が初期状態から次第に良くなっていくのである。

プロファイルDPについては前節で述べたので、反復改善法に固有の技術で新たに述べることはあまりない。ただ、反復改善法は乱数を使う。 $n$ 本の配列をふたつの配列群へ分割する仕方は $2^{n-1}-1$ 通りあるが、それを乱数でランダムに選ぶ。ただし、処理が進んでアライメントが整ってくると、プロファイルDPを行っても評価値がなかなか改善されない事態になる。そこで、評価値が改善されないときは、同じ分割を再び試さないように、それを記憶しておき、まだ試していない分割からランダムにひとつ選ぶと良い。 $2^{n-1}-1$ 通りの分割すべてにわたって、改善がみられなかったならば、収束と判定する。そのときの解が最終の多重アライメントになるのである。図21.9に反復改善法による結果の一例を示す。

<sup>20</sup>M. Kawai, A. Kishino, K. Naito, *FUJITSU Scientific and Technical Journal*, 27, 270-277 (1991).

<sup>21</sup>M. P. Berger, P. J. Munson, *CABIOS*, 7, 479-484 (1991).

プロフィール DP にあたっての注意点を述べよう、毎回プロフィール DP を行う前には、分割された配列群の各カラムを調べ、全配列にわたりギャップが入っているカラムがあったならば、それを除去する必要がある。そうしないと DP を行うごとに、ギャップが増える一方となつてうまくない。また、計算時間の削減に、前に述べたネットワークの右上隅と左下隅をカットするのが有効である。しかし、あまり大幅なカットをしてプロフィール DP を行うと、評価値がかえって下がる事態が発生するので注意を要する。

### 21.4.3 検討

反復改善法は、評価値をだんだん改善するが、どの程度の解に至るかは乱数の系列に左右される。また、処理する配列の本数が多いと、プロフィール DP の時間が無視できないうえに、収束までのサイクル数も大きなものとなる。

もし大規模な並列計算機があると、こうした問題を軽減できる。 $2^{n-1} - 1$ 通りの分割の仕方すべてを並列にプロフィール DP を行い、そのうちで最も評価値の良い結果を採用するのを繰返す方法を実現できる。そうすれば乱数を使うことがないので、安定して良い解が得られるし、収束までのサイクル数も少なくてすむ<sup>22</sup>。

## 21.5 シミュレーテッドアニーリング法

### 21.5.1 概要

シミュレーテッドアニーリング法（以下 SA と略す）は、計算機科学の分野でしばしば使われる最適化アルゴリズムであり、多重アライメントにも応用できる<sup>23</sup>。全体の評価値（とくにエネルギーと呼ばれる）を徐々に改善していくという点で反復改善法と似ているが、DP を使わないので、DP では扱えないような複雑な評価値を入れることができる利点がある。そのうえ、解の探索にあたっては、変形操作（オペレーションと呼ばれる）を工夫できるので、SA は柔軟性のある枠組といえる。しかしその反面、反復改善法以上に、収束にかかる時間が大きいので、現在の計算機パワーでは、小さな多重アライメントや、他の手法で得られた結果の修正に使われることになる。

### 21.5.2 方法

元来、アニーリングとは、物理系の焼きなまし過程を意味する。つまり、ある物質を高温から徐々に温度を下げることにより、結晶などの非常に安定な物質が得られる過程を指している。SA とは、この焼きなまし過程を模倣して組合せ最適化問題を解くアルゴリズムである。

図 21.10 に SA のアルゴリズムを示したが、その基本的な考え方は次の通りである。最適化問題を解くことは、エネルギー（評価値）が最良な状態を探し当てることである。素朴な方法で、初期状態から少しずつ状態を変形させて、エネルギーの良くなる（小さくなる）方向を探索していく方法があるが、それでは、ローカルミニマム（局所的にのみエネルギー最小な

<sup>22</sup>Ishikawa, Hoshida, Hirose, Toya, Onizuka, Nitta, *Proceedings of Fifth Generation Computer Systems '92*, (1992).

<sup>23</sup>石川幹人, 戸谷智之, 星田昌紀, 新田克己, 萩原淳, 金久實, 第2回公開ワークショップ「ヒトゲノム計画と情報解析技術」論文集, 112-115 (1991).

状態)に捕まってしまう。SAは、ローカルミニマムに捕まらずに、グローバルミニマム(大局的にエネルギー最小な状態)に至ることを可能にするものである。そのためSAでは、温度パラメータ $T$ を導入し、その値が高いときには、エネルギーの悪化する(大きくなる)方向の探索をも、ある程度(確率的に)許すのである。そして、温度パラメータを、高温から低温へと段々と下げていくことにより、探索範囲を確率的に徐々に絞り込み、グローバルミニマムへ至るのを可能にしている。

ここでの温度の下げ方を温度スケジュール(温度パラメータの列 $\{T_n\}$ に相当)と呼び、それを適切に設定すれば、十分な時間ののち、最適解を求めることが理論的には可能である<sup>24</sup>。しかし、現実には、処理時間が限られているので、その時間内に準最適解を求めることとなる。一般的な温度スケジュールは、設定した最高温から始めて、一定回数の変形を繰返したならば温度を9割程度に等比的に下げることが、エネルギー値が収束するまで行うものである。設定する最高温は扱う問題ごとに異なるが、予備実験をして、解が十分ランダムになる温度から始めるのが良い。各温度での繰返し回数の決定は難しく、扱う問題ごとに十分な回数を推測する方法もあるが、通常は、許される時間の制約から決めてしまうことが多い。

SA法を用いて多重アライメントを解くには、いくつかの定式化があるが、そのひとつを紹介しよう。SAの定式化とは、エネルギーとオペレーション(変形操作)とを定義することと言いかえられる。エネルギーとしては、いろいろ複雑な評価値を使えるが、ここではとりあえず、前節と同様の、Dayhoffマトリックスに基づいた評価値を使用することにしよう。ただし、エネルギーは小さい程よいとする慣例から、評価値の符号を反転して使う。

オペレーションの定義には、ギャップをどのように扱うかが鍵となる。SAでは、オペレーションが対称である(状態Aから状態Bへ変更するオペレーションがあれば、状態Bから状態Aへ変更するオペレーションが必ずある)ことが本質的である。その対称性を保持するには、配列の頭部や尾部に、あらかじめ十分な数のギャップを付加する方法がよい。たとえば、アライメントの初期状態として、まったくギャップの入っていない状態を採用するならば、図2 1・1 1 (a)のアライメント状態が初期状態となる。

そして、状態に対するオペレーションは次のように定義する。配列群から1本の配列をランダムに選び、それに対して、任意のギャップと任意のカラム位置をそれぞれランダムに選択し、選択されたギャップを選択されたカラム位置に移動させる。そして、間の部分の配列を、移動したギャップがあった方へ1カラム分移動する。たとえば先の初期状態において、RSVの配列が選ばれ、その配列の右末端のギャップと、中央部のAなるアミノ酸のあるカラム位置が選ばれた場合、図2 1・1 1 (b)のように変形される。

さらに、多重アライメントにはギャップが固まりで入りやすいことから、ギャップを長方形の固まりで動かすブロックオペレーションを、前のギャップ単独のオペレーションに拡張するかたちで導入する。そのブロックの作り方について、水平方向と垂直方向を分けて説明する。ひとつのオペレーションにおいて、両方向を考慮するとブロックオペレーションとなる。図2 1・1 1 (c)には、(b)に続いて、ブロックオペレーションが2回行われた後の状態の例を示している。

水平方向:ギャップがランダムに選択されたときに、選択されたギャップから右方向へみていき、右側にギャップが連続している限り、水平方向のひとつのブロックとして一緒に動かす。

<sup>24</sup>S. Kirkpatrick, C. D. Gelatt, M. P. Vecchi, *Science*, 220, 671-680 (1983).

垂直方向：配列群から配列をランダムに選ぶときに、任意の本数の任意の組合せを選び、それらの配列に共通してギャップが並ぶカラムだけを、水平方向の選択対象にし、垂直方向のひとつのブロックとして一緒に動かす。配列  $n$  本から任意本数選ぶ組合せの数は、 $2^n$  通りである。全配列を取り出す場合が選択されたならば、全配列にわたってギャップが並んでいるカラムを配列内から選択し、配列先頭や末端にそのギャップ列を移動するという、調整のための例外的なオペレーションとするとよい。

以上のような定式化を用いて、多重アライメントを SA で解いたところ。図 2 1・1 1 (a) から始めて 15 分ほど（オペレーション 1 万回に相当）で、図 2 1・1 1 (d) の結果が得られた。

### 2 1・5・3 検討

SA は、エネルギーやオペレーションを工夫できることから、柔軟な手法ではあるが、収束に必要な時間がかなり多くなる。そのため、小さな多重アライメントや、他の手法で得られた結果の修正に使われることが多い。また、反復改善法と同様に、乱数を用いるので、解が安定していないという欠点もある。

ここでも並列計算機を利用した欠点の緩和が考えられる。異なった乱数系列（あるいは異なった初期状態）を使った SA を並列に行い、そのうちから最も評価値の良い結果を採用する方法が効果的である。また、他の手法で得られた結果の修正をする際には、中低温からアニーリングを始めないとならないが、並列計算機を利用した温度並列 SA<sup>23</sup> は、そうした場合に柔軟に対応できる。

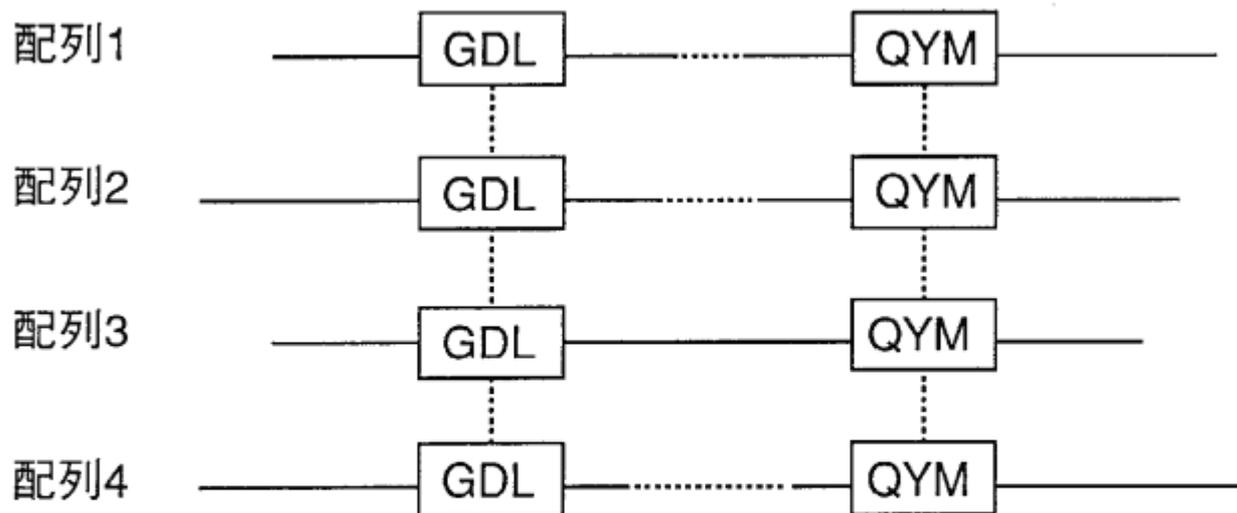
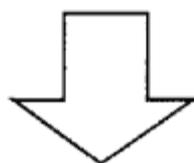
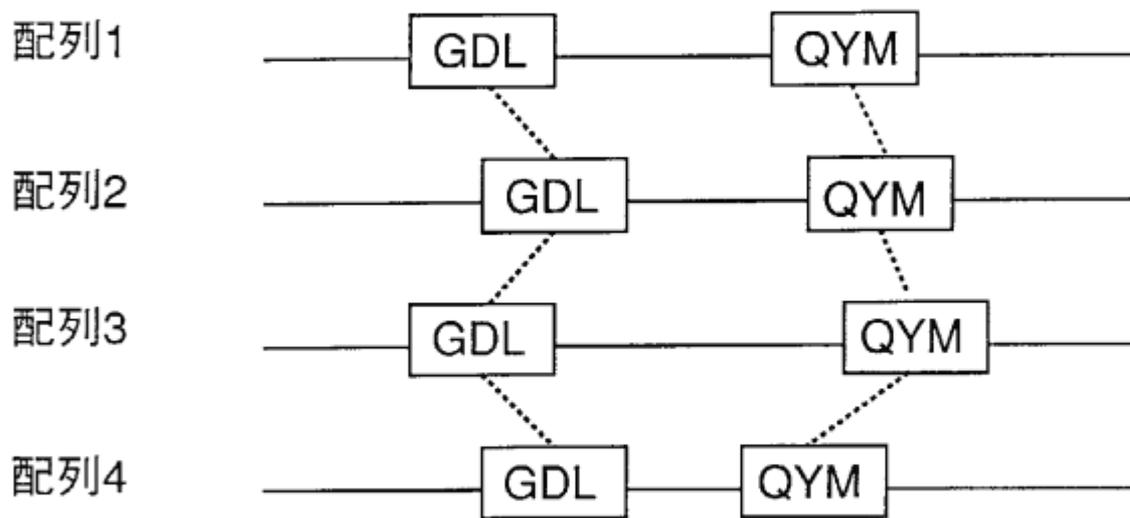
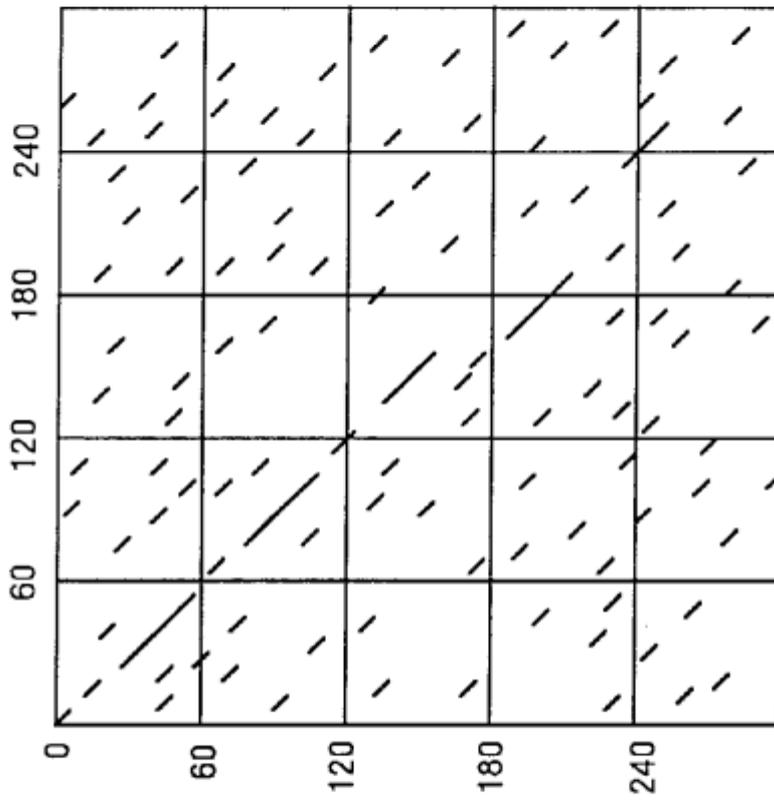


図 2 1 ・ 1 セグメント法による多重アライメント

配列Aのアミノ酸の位置番号



配列Bのアミノ酸の位置番号

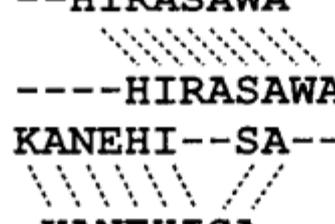
図21・2 ホモロジーマトリックスの例

(a)

配列1 : ISHIKAWA  
配列2 : HIRASAWA  
配列3 : KANEHISA  
配列4 : IKEHIRA

(b)

配列1 : ISHIKA--WA  
配列2 : --HIRASAWA  
配列2 : ----HIRASAWA  
配列3 : KANEHI--SA--  
配列3 : -KANEHISA  
配列4 : IK--EHIRA



(c)

配列1 : ----ISHIKA--WA  
配列2 : -----HIRASAWA  
配列3 : -KANE-HI--SA--  
配列4 : IK--E-HI--RA--

図21・3 単純組合せ法による多重アライメントの手順

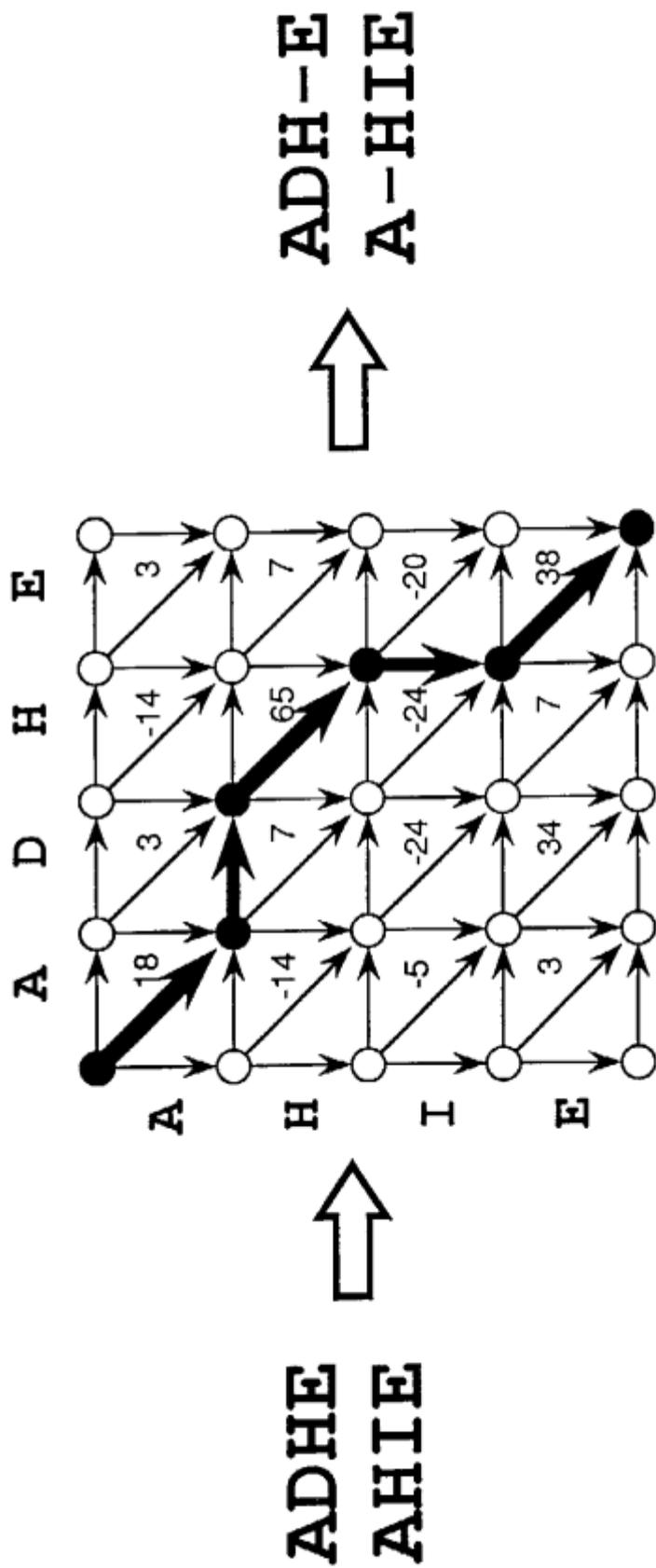
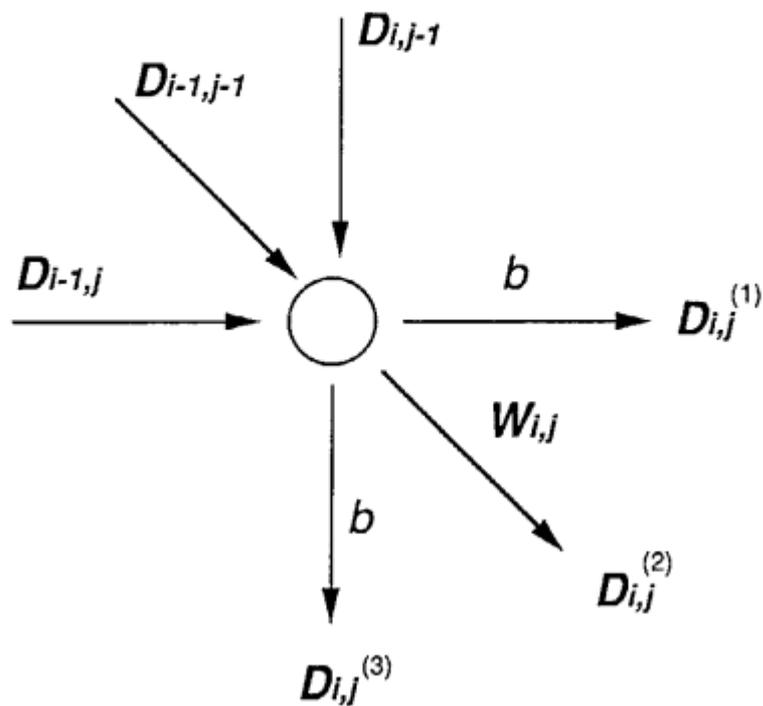


図21・4 ペアワイズDPによるアライメント



$$D_{i,j}^{(1)} = \text{Max}[(D_{i,j-1} + a), (D_{i-1,j-1} + a), D_{i-1,j}] + b$$

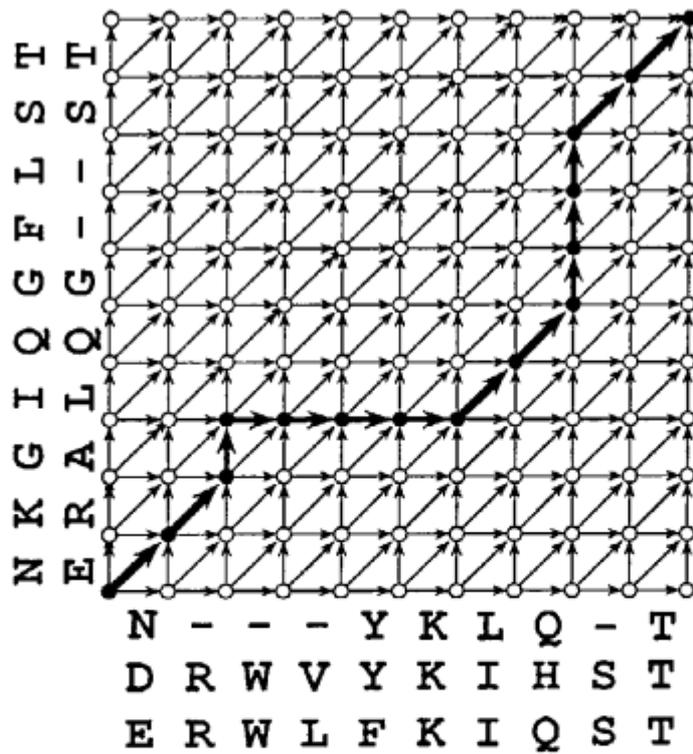
$$D_{i,j}^{(2)} = \text{Max}[D_{i,j-1}, D_{i-1,j-1}, D_{i-1,j}] + W_{i,j}$$

$$D_{i,j}^{(3)} = \text{Max}[D_{i,j-1}, (D_{i-1,j-1} + a), (D_{i-1,j} + a)] + b$$

図 2 1 ・ 5 ペアワイズDPにおける各ノードの処理  
(ギャップコスト  $a+bk$  の場合)



配列群 A



配列群 B

NKG---IQGFLST  
 ERA-----LQG--ST  
 N-----YKLQ----T  
 DR-WVYKIH---ST  
 ER-WLFKIQ---ST

図21・7 プロファイルDPによるアライメント

0 0 1 2      0 0  
 配列群Aのなかの1本    **IPERA-----LQG--STF**  
  
 1 2 0 0      1 0 0 2  
 配列群Bのなかの1本    **IPN-----YKLQ----TF**

- 0 : ギャップとはみなされない
- 1 : 第1ギャップとみなされる
- 2 : 第2ギャップとみなされる

図21・8 プロファイルD Pにおけるギャップ判定

```

HIV-1 :-----IEQLIK-----KEKVYLAWVPAHKGI---GGNEQVDKLVSAKIRKILF-LDGIDKAQDEHEKXYSNWRA--MASD-----F--
SIVagm:-----IALMIQ-----KQIYLQWVPAHKGI---GGNEIDKLVSKGIRRVLF-LEKIEEAQEKHERYHNNWKN--LADT-----Y--
EIAV  :-----GQKFA-----QLIILQHHSNSRQPW--DENKISQRGDKGFGSTGVFVVENIQEAQDEHENWHTSPKI--LARN-----Y--
SRV-1 :-----YNRSIP-----FYIGHVRAHSGPLGPI--AHGNQKADLATKTVASNIN--TNLESAQNAHTLHHLNAQT--LKLM-----F--
MPMV  :-----YNRSIP-----FYIGHVRAHSGPLGPI--AQGNQRADLATKIVASNIN--TNLESAQNAHTLHHLNAQT--LRLM-----F--
MMTV  :-----LQRLIHKRQEKFYIGHIRGHTGLPGPL--AQGNAYADSLTRILT-----ALESAQESHALHHQNAAA--LRFQ-----F--
IAPH18:-----LNRRFP-----VFITHVRAHSGPLGPM--SLGNDLADKATKLVATALS--THAQAAKEFHKRFHVTAET--LRRR-----F--
RSV   :--EDALSQRSAM-----AAVLHVRSHSEVPVGF--TEGNDVADSQATFQA---Y---PLREAKDLHTALHIGPRA--LSKA-----C--
BLV   :PSYALLYKSLLR-----HPAIVVGHVRSRSHSSAS--HPIASLNHYVDQLL-----PLETPEQWHLTHCNRA--LSR-----W--
HTLV-1:-----ALLPRLS-----RKVVYLHHVRSHTNLP--DPISRLNALDALLITPVL---QL-SPAELHSFTHCGQTA--LTLQ-----F--
HTLV-2:-----AALPPLL-----GKTIYLHHVRSHTNLP--DPISFNEYTDSLILAPLV---PL-TPQGLHGLTHCNQRA--LVS-----F--
M-MULV:-----YTSEHFH-----YTVTDIKDLTKLGAIV--DKTKKYVWVYQKPVMPDQF---TFELLDPLHQLTHLSFSK--MKALLERSHSPYY--
AKV   :-----YTPAYFH-----YTETDLKKLRELGATY--NQSKEYVWFQKPVMPDQF---VFELLDLHRLTHLGYQK--MKALLDRGESPPY--
FoamyV:-----LDQLLQ-----GH--YIKGYPKQYTYFLEDGKVKVSRPEGVKIIPPQS--DRQKIVLQAHNLAHTGREATLLKIA-----NLYW

```

```

-NLPPVVA---KEIVASCDKCKQLKGEAMH---G---QVDCS--P-GIWQ---LD-CTHLEG-----KV-ILVAVHVASGYIEAEVIPAET
-GLPQIVA---KEIVAMCPKCQIKGEPVH---G---QVDAS--P-GTWQ---MD-CTHLEK-----KV-VIVAVHVASGFIEAEVIPRET
-NIPREQA---RQIVRQCPCICATYLPVPHL---G---VNPRGLLPNHIWQ---MD-VTHYSE--FGNLKY-IHVSIDTFSGFLL-----
-NIPREQA---RQIVKQCPCICVYLPVPHL---G---VNPRGLFPNHIWQ---MD-VTHYSE--FGNLKY-IHVSIDTFSGFLL-----
-HITREQA---REIVKLCPCNCPDWGHAPQL---G---VNPRGLKPRVLWQ---MD-VTHVSE--FGNLKY-VHVTVDTYSHFTF-----
-ALSRKEA---REIVTQCQCCEFLPTPHM---G---INPRGIRPLQMWQ---MD-VTHIPS--FGRLQY-VHVSVDTCSGVMF-----
-NISMQA---REVVQTCPCNCS-APALEA---G---VNPRGLGPLQIWQ---TD-FTLEPR--MAPRSW-LAVTVDTASSAIVV-----
-PNPRISAWDPRSPATLCETCQKLNPTGG---GKMRITIQRGWAPNHIWQ---AD-ITHYK---YKQFTYALHVFVDTYSGA-----
-GATTTEA---SNILRSCHACRGGNPQHQM-PRG---HIRRGLLPNHIWQ---GD-ITHFK---YKNTLYRLHVWVDTFSGA-----
-GATPREA---KSLVQTCCTCQTINSQHMM-PRG---YIRRGLLPNHIWQ---GD-VTHYK---YKKYKYCLHVWVDTFSGA-----
-MLNRDRTL--KNITETCKACAQVNASKSAVKQG---TRVRGHRPGTHWE---ID-FTEIKPGLYG-YKY-LLVFIDT-----
-MLNRDKTL--QYVADSTVCAQVNASKAKIGAG---VRVRGHRPGSHWE---ID-FTEVVPGLYG-YKY-LLVFVDT-----
PNMRKDVV---KQLGR-CQQCLITNASNKA--SG---PILRDRPQKPFDKFFIDYIGPLPP--SQGYLY-VLVVVVVGH-----

```

図21・9 反復改善法による多重アライメント結果の例  
(レトロウイルスのエンドヌクレアーゼの一部)

```

begin
   $X_0$  := 初期状態;
   $\{T_n\}_{n=0,\dots,N-1}$  := 温度スケジュール;
  for  $n := 0$  to  $N-1$  do
    begin
       $X'_n$  := 変形操作を施した  $X_n$ ;
       $\Delta E := E(X'_n) - E(X_n)$ ;
      if  $\Delta E < 0$  then
         $X_{n+1} := X'_n$ ;
      else
        if  $\exp(-\Delta E/T_n) \geq$  区間 $[0,1]$ の乱数 then
           $X_{n+1} := X'_n$ ;
        else
           $X_{n+1} := X_n$ ;
      end;
    解出力  $X_n$ ;
  end;
end;

```

図21・10 シミュレーテッドアニーリングのアルゴリズム

(a) 初期状態

```
SMRV :-----GFILATPQTGEASKNVISHVIHJLATIGPKHTIKTDNGPGYTGKNFQDFCQKLQI-----
MMTV :-----YSHFTFATARTGEATKDVLQHLAQSFAYMGIPQKIKTDNAPAYVSRSIQEFLARW-----
IAP  :-----GVMFATTLTGEKASYVIQHCLEAWSAWGKPRIKTDNGPAYTSQKFRQFCRQMDVT-----
RSV  :-----IVVTQHGRVTSVAVQHHWATAIAVLGRPKAIKTDNGSCFTSKSTREWLARWGIAH-----
HTLV-1:-----SGAISATQKRKETSSEAISLLQAI AHLGKPSYINTDNGRAYISQDFLNMCTSLA-----
HTLV-2:-----DTFSGAVSVSCKKETS CETISAVLQAI SLLGKPLHINTDNGPAFLSQEFQEFCT-----
BLV  :-----HASAKRGLTTQTTIEGLLEAIVHLGRPKKLN TDQGANYTSKTFVRFCCQFGVSL S-----
```

(b) 1回変形後

```
SMRV :-----GFILATPQTGEASKNVISHVIHCLATIGPKHTIKTDNGPGYTGKNFQDFCQKLQI-----
MMTV :-----YSHFTFATARTGEATKDVLQHLAQSFAYMGIPQKIKTDNAPAYVSRSIQEFLARW-----
IAP  :-----GVMFATTLTGEKASYVIQHCLEAWSAWGKPRIKTDNGPAYTSQKFRQFCRQMDVT-----
RSV  :-----IVVTQHGRVTSVAVQHHWATAIAVLGRPK-AIKTDNGSCFTSKSTREWLARWGIAH-----
HTLV-1:-----SGAISATQKRKETSSEAISLLQAI AHLGKPSYINTDNGRAYISQDFLNMCTSLA-----
HTLV-2:-----DTFSGAVSVSCKKETS CETISAVLQAI SLLGKPLHINTDNGPAFLSQEFQEFCT-----
BLV  :-----HASAKRGLTTQTTIEGLLEAIVHLGRPKKLN TDQGANYTSKTFVRFCCQFGVSL S-----
```

(c) 3回変形後

```
SMRV :-----GFILATPQTGEASKNVISHVIHCLATIGPKHTI-----KTDNGPGYTGKNFQDFCQKLQI-----
MMTV :----YSHFTFATARTGEATKDVLQHL---AQSFAYMGIPQ-----KIKTDNAPAYVSRSIQEFLARW-----
IAP  :-----GVMFATTLTGEKASYVIQHCLEAWSAWGKPRIKTDNGPAYTSQKFRQFCRQMDVT-----
RSV  :-----IVVTQHGRVTSVAVQHHWATAIAVLGRPK-AIKTDNGSCFTSKSTREWLARWGIAH-----
HTLV-1:-----SGAISATQKRKETSSEAISLLQAI AHLGKPSYINTDNGRAYISQDFLNMCTSLA-----
HTLV-2:----DTFSGAVSVSCKKETS CETIS---AVLQAI SLLGKPLHINTDNGPAFLSQEFQEFCT-----
BLV  :----HASAKRGLTTQTTIEGLLEAIV---HLGRPKKLN TDQGANYTSKTFVRFCCQFGVSL S-----
```

(d) 最終状態

```
SMRV :-----GFILATPQTGEASKNVISHV-IHCLATIGPKHTIKTDNGPGYTGKNFQDFCQKLQI-----
MMTV :----YSHFTFATARTGEATKDVLQHL-AQSFAYNGIPQKIKTDNAPAYVSRSIQEFLARW-----
IAP  :-----GVMFATTLTGEKASYVIQHC-LEAWSAWGKPR-IKTDNGPAYTSQKFRQFCRQMDVT-----
RSV  :-----IVVTQHGRVTSVAVQHHWATAIAVLGRPKAIKTDNGSCFTSKSTREWLARWGIAH-----
HTLV-1:----SGAISATQKRKETSSEAISSL-LQAI AHLGKPSYINTDNGPAYISQDFLNMCTSLA-----
HTLV-2:--DTFSGAVSVSCKKETS CETISAV-LQAI SLLGKPLHINTDNGPAFLSQEFQEFCT-----
BLV  :-----HASAKRGLTTQTTIEGL-LEAIVHLGRPKKLN TDQGANYTSKTFVRFCCQFGVSL S-----
```

図21・11 シミュレーテッドアニーリング法による多重アライメント