

## ゲノム解析に用いる DP マッチングの分割統治法による高速化

中川 匠<sup>1)</sup> 吉原 郁夫<sup>2)</sup> 山森 一人<sup>3)</sup> 安永 守利<sup>4)</sup>A Divide-and-Conquer Method of Dynamic Programming for  
Genome InformaticsTakumi NAKAGAWA Ikuo YOSHIHARA Kunihiro YAMAMORI  
Moritoshi YASUNAGA

## ABSTRACT

Alignment based on DP-matching is used to extract unknown feature pattern embedded in genome sequence. To extract feature pattern with base length  $n$ , complete set of similarity of candidate pattern with base length  $N = 2, 3, \dots, n$  is required. When extending base length  $n$ , number of candidate pattern and execution time increase exponential order  $O(4^{n+1})$ . This paper propose a new faster method of extracting feature pattern by reusing similarity which is calculated past step. The method enable extracting feature pattern with base length  $n = 10$  to speed up as much as 9 times than conventional method.

Key Words:

Genome Informatics, Dynamic Programming, Divide-and-Conquer, Alignment, Feature Extraction

## 1 はじめに

ヒトゲノム計画をはじめとするゲノム (genome) 解析プロジェクトの進展により, さまざまな生物種のゲノム情報が大量に蓄積されている<sup>1)</sup>. 現在, これら大量のゲノム情報からその働きや意味を明らかにすることが重要な課題になっている.

遺伝子とはタンパク質をコードしている DNA (Deoxyribo Nucleic Acid: デオキシリボ核酸) 塩基配列の一部であり, 実際に遺伝子が読み取られタンパク質が作られることを遺伝子の発現という. この過程において発現調節要素と呼ばれる部分塩基配列が重要な役割を果たしている事がわかっている. 現在までにいくつかの発現調節要素が発見されているが, いまだに発見されていない発現調節要素もたくさん存在すると考えられている. これら未知の発現調節要素を発見することは遺伝子発現のメカニズムを解明するうえで重要な意味を持つ. そこでこれまで発現調節要素の候補となる特徴的な部分塩基配列 (特徴

パターン) を自動的に抽出する研究を行ってきた<sup>2)</sup>.

特徴パターンの抽出にあたっては考慮すべき問題が2つある. 1つは特徴パターンは常に同じ部分塩基配列ではなく通常は多少変化して存在するということがある. もう1つは抽出しようとする特徴パターンの長さや塩基配列, 出現位置が未知であると言う問題である. そこでアラインメントを用いた探索を総当たりで行う必要がある.

ゲノム解析の分野ではアラインメントを効率的に計算する方法として動的計画法 (Dynamic Programming: DP) による DP マッチングが用いられる<sup>3)</sup>. しかし特徴パターン抽出では総当たりでアラインメントを求めなければならず, DP マッチングの計算量が増大し特徴パターン抽出に膨大な時間がかかる. そこで分割統治法の考えを用いて DP マッチングに要する計算量を削減する方法を提案する. それは抽出をおこなう候補パターンを短い部分配列に分割し DP マッチングを行う方法である. この手法を用いることで候補パターンの長さが短くなり DP マッチングの計算量が削減される. また一度計算したアラインメントスコアを保存し再利用することで DP マッチングの実行回数自

<sup>1)</sup>情報工学専攻学生<sup>2)</sup>情報システム工学科教授<sup>3)</sup>情報システム工学科助教授<sup>4)</sup>筑波大学電子・情報工学系教授

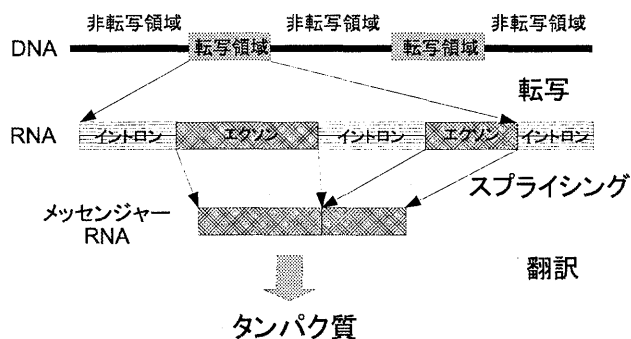


図.1 遺伝子の発現

体を削減することができる。

## 2 遺伝子の発現と特徴パターン

### 2.1 ゲノムの構造と特徴

ゲノムとはある一つの生物が持っている完全な染色体のことで、染色体は一連の DNA 分子から構成されている<sup>4)</sup>。DNA は2本のポリヌクレオチド鎖がらせん状に巻きあつた構造である。ポリヌクレオチド鎖はヌクレオチド (nucleotide) と呼ばれる分子が鎖状につながってできている。ヌクレオチドは1個の糖、1個のリン酸、1個の塩基から成り立っており、塩基の部分だけが異なる4種類が存在する。異なる塩基の部分はアデニン (adenine; A), シトシン (cytosine; C), グアニン (guanine; G), チミン (thymine; T) の4種類があり、これら4つの塩基の配列によって DNA は表される。

DNA は多数の異なる遺伝情報を担う区画 (遺伝子) にわかれており、実際に遺伝子が読み取られタンパク質が作られることを遺伝子の発現という。

遺伝子発現は図1に示すメカニズムによって行われる<sup>5)</sup>。まず DNA の塩基配列が RNA に移し換えられる。この過程は転写と呼ばれる。次に、スプライシングという作業によってイントロンが取り除かれメッセンジャー RNA (mRNA) がつくられる。mRNA の塩基配列はアミノ酸をコードしており、そのアミノ酸が順に結合してタンパク質が生成される。この過程は塩基配列から質的にまったく異なるアミノ酸配列への変換であり、翻訳と呼ばれる。

### 2.2 発現調節要素と特徴パターン

DNA 塩基配列の中で転写される領域を転写領域といい、その始まりの位置を転写開始点という。遺伝子の近傍にはいくつかの特徴的な短い配列が存在し、これらの適切な組合せがシグナルとなり転写開始点が決定されている<sup>6)</sup>。これらは発現調節要素とよばれ、TATA ボックスや CCAAT ボックスなどが良く知られ

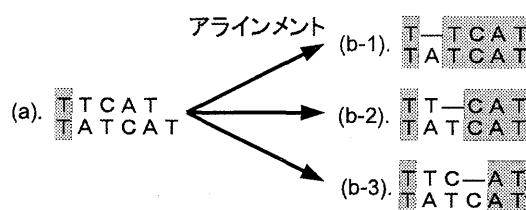


図.2 アラインメント

ており、まだ発見されていないものもたくさん存在すると考えられる。

このような未知の発現調節要素の発見は遺伝子発現のメカニズムを解明する上で重要である。しかし大量に存在する DNA 塩基配列データの中から未知の発現調節要素を発見する事は非常に困難である。

これまでに見付かっている発現調節要素の特徴から、転写開始点を基準として一定の位置付近で頻出する塩基部分配列が発現調節要素の候補パターンと考えられる。そこで発現調節要素の候補となるパターンを計算機を用いて抽出し、それを元に生物学的な考察を加えて発現調節要素を発見する。この発現調節要素の候補となるパターンを特徴パターンと呼ぶ。

発現調節要素は常に同じ塩基配列で存在するとは限らず、通常は多少変化した塩基配列として存在することもある。そのため類似した塩基配列も特徴パターンとして抽出できるようなパターンマッチングを用いる必要がある。

また特徴パターンを抽出する上でその長さや出現位置に関する先験的知見は全くなく、塩基配列のどの位置にどんな特徴パターンが現れるかわからない。そのため総当たりのパターンマッチングを行わなければならない。長さ  $N$  までの特徴パターンを抽出するには、長さ  $N$  までの全通りの塩基配列を候補パターンとして用意し DNA 塩基配列データの全ての位置でパターンマッチングする必要がある。

### 2.3 アラインメント

類似した塩基配列を特徴パターンとして抽出するため、アラインメントを用いたパターンマッチングを利用する。アラインメントとは2つの塩基配列に対して塩基の間にギャップと呼ばれるすき間を入れ、塩基配列の塩基の対応関係を求めたものである。<sup>7)</sup>

図2の(a)は2つの類似した塩基配列を上下に並べたものである。また同図の(b-1), (b-2), (b-3)はアラインメントの例を示しており、塩基配列中に挿入されている'-'はギャップを表す。このとき網掛けが施されている部分は、上下の塩基配列で塩基が一致している

表.1 類似度計算テーブル

状態	類似度
一致	+1
不一致	-1
ギャップ	-2
連続ギャップ	-1

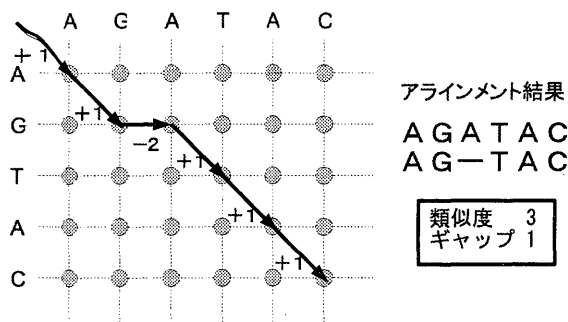


図.3 DP マッチング

ことを表す。塩基が一致している部分が多いほどその2つの塩基配列はより似ている塩基配列であると考えられる。同図の(a)では塩基配列の上下で一致している塩基は1つに過ぎないが、ギャップを挿入する事で(b-1)では5つの塩基が縦に一致している。このようにギャップを挿入することで二つの塩基配列がどれだけ一致しているかをもとめることができる。

ここで対応づけられた塩基の一致・不一致を表1のように定義し、これをもとに類似度を求めることにする。

表1の Match・Mismatch はそれぞれ対応づけられた塩基が一致・不一致であるときの類似度である。また Open gap はギャップが挿入されたときの類似度であり、Extended gap はギャップが延長されたときの類似度である。これを用いると図2の(b-1)の2つの塩基配列で対応する各塩基の類似度の和は3である。この類似度の和をアラインメントスコアと呼ぶ。

図2の(b-1)以外にも(b-2)(b-3)のように、アラインメントはギャップの位置を変えることでいくつも生成できる。このときもっともアラインメントスコアが高くなるアラインメントを最適アラインメントと呼ぶ。特徴パターンの抽出には最適アラインメントによるパターンマッチングを用いる。

### 3 特徴パターン抽出とその高速化

#### 3.1 DP マッチングによる最適アラインメントの導出

最適アラインメントを効率的に計算する方法としてDP マッチングが一般的に用いられる。

図3は二つの塩基配列AGATACとAGTACの最適アラインメントをDP マッチングにより求める例である。ア

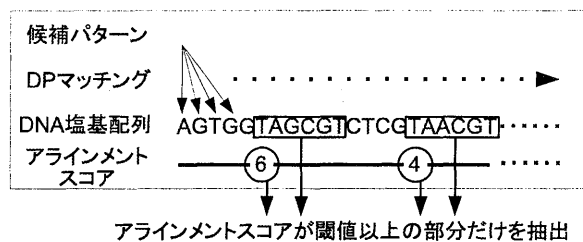


図.4 アラインメントを用いた部分配列の抽出

ラインメントスコアの計算には表1を用いている。DP マッチングでは最適アラインメント問題を図3のような2次元の格子状にして表す。この図の横方向は1つ目の塩基配列に対応し、縦方向は2つ目の塩基配列に対応する。また右下方向の矢印は対応する塩基の類似度が対応づけられていて、横方向と縦方向の矢印にはギャップの類似度が対応づけられている。そして、左上の点から右下の点までの類似度の和が最大となる経路を求めれば、最適アラインメントが得られる。この例では最適アラインメントとしてAGATACとAG-TACが得られ、そのときのアラインメントスコアは3となる。

#### 3.2 従来法による特徴パターンの抽出

従来の研究で用いられてきた特徴パターンの抽出法について説明する。

DNA 塩基配列中に含まれている候補パターンを抽出するために、候補パターンとDNA 塩基配列のDP マッチングを行う。図4に示すように、DP マッチングはDNA 塩基配列の先頭から1塩基ずつずらしながら、DNA 塩基配列の終端に達するまで局所的に行う。その中からアラインメントスコアが一定の閾値 *Threshold* を越えた部分の位置とスコアを保存する。このようにしてDNA 塩基配列から候補パターンと同じか非常に似ている部分配列が含まれている位置とそのアラインメントスコアを抽出する。長さ *N* の特徴パターンを抽出するときは、長さ *N* のありうる全ての塩基配列を候補パターンとして以上の操作を行う。

*Threshold* は候補パターンを特徴パターンとして抽出する最低のアラインメントスコアであり、以後は特徴抽出閾値と呼ぶ。いま候補パターンの長さを *l*、アラインメントに含むことが許されるギャップの数を *G* とする。またアラインメントにおける塩基の一致による類似度を *Similarity(match)*、ギャップ挿入による類似度を *Similarity(opengap)* とする。このとき *Threshold* は式(1)のように表される。

$$\begin{aligned} \text{Threshold} &= \text{Similarity}(\text{match}) \times (l - G) \\ &\quad - \text{Similarity}(\text{opengap}) \times G \end{aligned} \quad (1)$$

DNA 塩基配列データは1つではなく大量に存在し、特徴パターンを抽出するには他の DNA 塩基配列データにおいても類出する部分配列を見つけ出さねばならない。そこで全ての DNA 塩基配列データに対し候補パターンかそれに類似した部分配列が含まれている位置とそのアラインメントスコアを抽出する。そして抽出された情報をもとに特定の位置に類出するパターンを見つけ出し、特徴パターンとする。

しかし先に述べたように生物進化の過程で塩基の挿入・欠損・変異などが起こるため、なんらかの意味があるパターンが存在したとしてもそれが全く同じ位置には存在せず、少しずれた位置に存在することもある。そこで一定の範囲内に類似度を相対的に表す2種類のスコア、 $Score_M$  と  $Score_T$  が提案されている<sup>2)</sup>。M個の DNA 塩基配列データが存在する時、塩基配列の位置  $i$  におけるそれぞれのスコアを求める式を以下に示す。

$$\begin{aligned} \text{Score}_M(i) &= \max_{i-M \leq j \leq i+M} (\text{Similarity}(j)) \\ \text{Score}_T(i) &= \sum_{j=i-M}^{i+M} (\text{Similarity}(j)) \end{aligned}$$

これらのスコアをもとに塩基配列の特徴を抽出する。

### 3.3 特徴パターン抽出に必要な計算量

塩基長  $L$  の DNA 塩基配列データ  $M$  個から長さ  $l$  の特徴パターン抽出を行うことを考える。比較する2本の文字列の長さが  $I, J$  のとき DP マッチングの実行に必要な計算量  $c_1$  は式(2)のようになる<sup>8)</sup>。

$$c_1 = \frac{I}{2}(J+1)(J+2) \quad (2)$$

ここでアラインメントに含まれるギャップの数を  $g$  とした時、 $0 \leq g \leq G$  を特徴パターンとして抜き出す条件とする。候補パターンの長さが  $l$  の時、高々  $l+G$  の DNA 塩基配列と比較すれば良い。そこで式(2)の  $I = l+G$ ,  $J = l$  として、長さ  $l$  の候補パターンの DP マッチングに必要な計算量  $c_2$  を式(3)に表す。

$$c_2 = \frac{l+G}{2}(l+1)(l+2) \quad (3)$$

候補パターンの長さが  $l$  の時、考えられる候補パターンの数は  $4^l$  である。長さ  $L$  の塩基配列データの全て

の位置で全ての候補パターンと DP マッチングを行うため、DP マッチングの実行回数は  $4^l L$  となる。また  $M$  個の DNA 塩基配列データについてこれを行うので、DP マッチングを実行する回数は  $4^l LM$  となる。これらのことから長さ  $l$  の特徴パターンを抜き出すために必要な計算量  $c_3$  は式(4)のようになる。

$$c_3 = 4^l LM \frac{l+G}{2}(l+1)(l+2) \quad (4)$$

長さ  $N$  までの特徴パターンを抽出するためには、長さ2から  $N$  までの特徴パターン抽出を行わなければならない。よって長さ  $L$  の DNA 塩基配列データ  $M$  個から長さ  $N$  の未知の特徴パターンを抽出するのに必要な計算量  $c_4$  は式(5)で表される。

$$c_4 = \sum_{l=2}^N 4^l LM \frac{l+G}{2}(l+1)(l+2) \quad (5)$$

### 3.4 分割統治法を用いた高速化

従来法では抽出したい特徴パターンの長さが延びるに従い DP マッチングにかかる計算量が指数関数的に増大する。そのため長い特徴パターンの抽出には膨大な時間が必要となる。そこで分割統治法の考えを用い計算量を削減する方法を提案する。

#### 3.4.1 候補パターンの分割

提案手法は候補パターンを分割し、分割された部分配列について DP マッチングを行う方法である。長さ  $l$  の候補パターンを途中で分割し、 $l = a + b$  となる長さ  $a$  と  $b$  の部分配列にする。分割してできた部分配列についてそれぞれ DP マッチングを行い、それぞれのアラインメントスコアを計算する。

候補パターンを分割して DP マッチングを行う時の計算量を求める。長さ  $l$  の候補パターンの DP マッチングに必要な計算量は式(3)で表される。そこで候補パターンを途中で2分割し、 $l = a + b$  となる長さ  $a$  と  $b$  の部分配列  $A, B$  にする。提案手法は分割されてきた部分配列について DP マッチングを行い、その結果をもとにアラインメントスコアを計算する方法である。この時の計算量  $c_5$  は式(6)で表される。

$$c_5 = \frac{a+G}{2}(a+1)(a+2) + \frac{b+G}{2}(b+1)(b+2) \quad (6)$$

候補パターンの数は全部で  $4^l$  個存在し、長さ  $L$  の DNA 塩基配列データの全ての位置で DP マッチングを行う必要がある。よって長さ  $l$  の特徴パターン抽出全体で必要となる計算量  $c_6$  は式(7)で表される。

$$c_6 = 4^l LM \left\{ \frac{a+G}{2}(a+1)(a+2) + \frac{b+G}{2}(b+1)(b+2) \right\} \quad (7)$$

### 3.4.2 候補パターン全体の DP マッチング

部分配列  $A$  と  $B$  のアラインメントスコアをそれぞれ  $score(A)$  と  $score(B)$  とする。

特徴パターン抽出の条件として  $G$  個のギャップ挿入が許されているとしたとき、分割した候補パターンに許されるギャップの数も  $G$  となる。そこでギャップ挿入によるペナルティを  $P(G)$ 、部分配列  $A, B$  がとりうる最大のアラインメントスコアを  $MAXscore(A), MAXscore(B)$  とすると、式 (8) の条件が満たされたとき候補パターン全体の DP マッチングを実行する。

$$\begin{cases} score(A) \geq MAXscore(B) - P(G) \\ score(B) \geq MAXscore(A) - P(G) \end{cases} \quad (8)$$

### 3.4.3 アラインメントの保存と再利用

候補パターンを部分配列に分割しアラインメントスコアを求めた後、次の候補パターンについても同様に分割してアラインメントを求める。

このとき過去に計算したアラインメントを再利用できる場合がある。いま長さ 6 の特徴パターン抽出を行うことを考える。このとき候補パターンを **AAAAAA**, **AAAAAG**, **AAAAAC**, ..., **TTTTTC**, **TTTTTT** と変化させ、前後 3 塩基の部分配列に分割して DP マッチングを行うとする。このとき前半の部分配列について何度も同じアラインメントを計算することになる。そこで一度計算した部分配列のアラインメントを保存し再利用することにする。

アラインメントの保存と再利用による計算量の変化を求める。長さ  $l$  の候補パターンを  $l = a + b$  となる長さ  $a$  と  $b$  の部分配列に分割して DP マッチングを行うとする。このときの計算量  $c_7$  は式 (9) のように求められる。

$$c_7 = 4^a LM \frac{a+G}{2}(a+1)(a+2) + 4^b LM \frac{b+G}{2}(b+1)(b+2) \quad (9)$$

## 3.5 並列処理による高速化

第 3.4 節では候補パターン長が延びるに従い指数関数的に増加する DP マッチング呼び出し回数をアルゴリズム的に減少させ、高速化する方法を提案した。本研究ではこれに加え処理の並列化を行うことでさらなる高速化を図る。

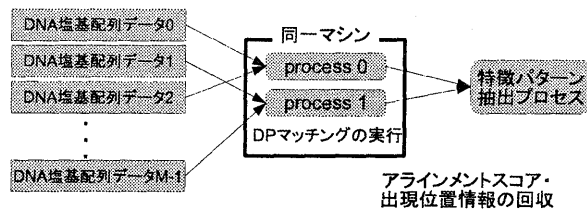


図.5 DP マッチングの並列処理

表.2 実験に使用した計算機の仕様  
ハードウェア

CPU	Intel Xeon 3.2GHz DualCPU
2次キャッシュ	2MB
メインメモリ	4GB

ソフトウェア

OS	Redhat Enterprise Edition Linux WS 3.0
コンパイラ	Intel C++ Compiler Linux Version 9.0

処理は図 5 のように行われる。

まず DP マッチング専用のプロセスである `process[0]` と `process[1]` を生成する。これらのプロセスに均等に DNA 塩基配列データを与え、候補パターンとの DP マッチングを行う。

特徴パターン抽出プロセスは全ての DP マッチング専用プロセスから DP マッチングの結果を受け取り、特徴パターン抽出処理を行う。

## 4 提案手法の性能評価

### 4.1 実験データと実験条件

特徴パターン抽出実験に用いる DNA 塩基配列データには細胞性粘菌の非転写領域を用いる。細胞性粘菌は生活様式が特徴的な生物で、培養が容易であることなどから細胞分化・形態形成のメカニズム、細胞集合、細胞間通信などの生理・生化学、遺伝学、微細構造学の有力な研究材料として用いられている<sup>9)</sup>。なお使用する細胞性粘菌のデータは筑波大学の漆原研究室によって独自に作成されたものである。今回は細胞性粘菌の塩基配列から転写開始点の上流 2000 塩基を切り出したデータ 332 個を用意した。特徴パターンに含まれるギャップの数  $G$  は 2 とし、特徴抽出閾値  $Threshold$  は式 (1) より求める。また類似度の計算には表 1 を用いる。実験に用いた計算機の詳細は表 2 のとおりである。

### 4.2 実験

細胞性粘菌の塩基配列データから長さ 10 の特徴パターン抽出実験を行う。特徴パターン抽出は従来手法を用いたプログラムと提案手法を用いたプログラムの両方で行い、その実行時間と DP マッチング処理の呼び出し回数を比較する。なお提案手法を用いたプログラムでは候補パターンの分割数を 2 に設定し、部分配

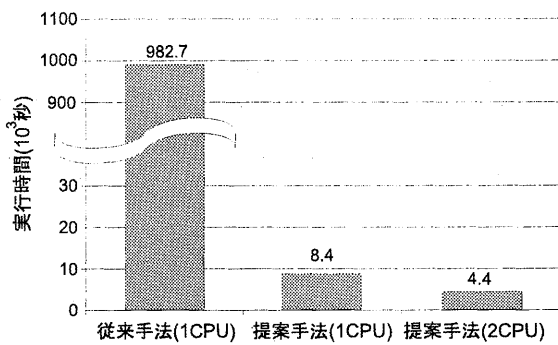


図.6 手法による実行時間の変化

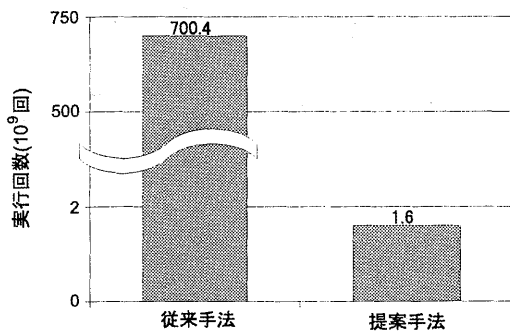


図.7 提案手法による DP マッチング実行回数の変化

列の長さは5とする。

図6は特徴パターン抽出に要した時間を示すグラフである。なお提案手法による実行時間は並列化を行わずに1CPUで実行した場合と、プロセスを2つに分け2CPUを用いた場合の2つを測定している。

図7は従来手法と提案手法のプログラム実行時に呼び出されたDPマッチングの回数を示す。

#### 4.3 考察

従来手法を用いたプログラムの実行時間は $982.7 \times 10^3$ 秒であるが、1CPUを用いた提案手法では $8.4 \times 10^3$ 秒で実行を終了した。このことから分割統治法を用いたアルゴリズムにより117倍の高速化を達成したと考えられる。また並列化処理により2CPUを用いた提案手法による実行時間は $4.4 \times 10^3$ 秒でありおよそ223倍の高速化を達成した。これらのことから提案手法の有効性が確認された。またプログラム実行時に呼び出されるDPマッチングの回数は従来手法では7000億回以上であったが、提案手法ではおよそ16億回に削減されている。DPマッチング呼び出し回数の削減率はおよそ $\frac{1}{430}$ である。DPマッチングの呼び出し回数が $\frac{1}{430}$ となったにもかかわらず実行時間がおよそ117倍にしかならなかった。これは、提案手法においてアラインメントの保存のための処理が追加されたためだと考えられる。

なお並列化により実行時間は約 $\frac{1}{2}$ に短縮されており、利用CPUの増加がそのまま実行時間の短縮につながっていることがわかる。これは提案手法で用いた並列化がプロセス複製によるものであり、一般的な並列計算で起こる通信によるオーバーヘッドが発生しないことが理由と考えられる。

#### 5 おわりに

本研究は特徴パターン抽出に必要な処理の主要部分であるDPマッチングの計算量を削減するため、候補パターンを分割したDPマッチングとアラインメントスコアの再利用による手法を提案した。またアルゴリズムによる高速化に加えて、並列処理による高速化手法を提案した。

これにより長さ10の特徴パターン抽出実験では、従来手法に対して117倍の高速化を達成した。さらに処理の並列化を用いることで、従来手法に対しておよそ223倍の高速化を達成した。このことから本手法を用いることで特徴パターン抽出が高速に行われることを示した。

今後の課題として高速化により可能となった長い特徴パターンの抽出を行い、なんらかの特徴パターンを抽出することが挙げられる。

#### 参考文献

- [1] 林田秀宜: “DNA データベース”, 情報処理学会誌, Vol. 31, No. 7, pp. 875-877 (1990).
- [2] T. ONITANI, I. YOSHIHARA, K. YAMAMORI and M. YASUNAGA: “Extraction of feature patterns embedded in non-transcribed region of dictyostelium discoideum”, The 5th International Conference on Simulated Evolution and Learning(SEAL'04), Vol. 123, (2004).
- [3] D. W. Mount, 岡崎康司, 坊農秀雅: バイオインフォマテイクスゲノム配列から機能解析へ, メディカル・サイエンス・インターナショナル 2002.
- [4] 菊池韶彦, 榊佳之, 水野猛, 伊庭英夫 (編): 遺伝子 第7版, 東京化学同人 2002.
- [5] 美宅成樹, 金久實: ヒトゲノム計画と知識情報処理, 培風館 1995.
- [6] 村松正實, 木南凌 (編): ヒトの分子遺伝学 第3版, メディカル・サイエンス・インターナショナル 2005.
- [7] 岡崎康司, 坊農秀雅: ゲノム情報はこう活かせ!, 羊土社 2005.
- [8] 上坂吉則, 尾関和彦: パターン認識と学習のアルゴリズム, 文一総合出版 1990.
- [9] G. Glöckner, L. Eichinger, K. Szafranski, et al.: “Sequence and analysis of chromosome 2 of dictyostelium discoideum”, Nature, Vol. 418, (2002).