

DNA 配列におけるプローブの順序付けに必要な最小フラグメント集合

Minimum Fragments

LA-3

for Deciding Probe Sequences for DNA Strands

田村武幸*

Takeyuki Tamura

土田大輔*

Daisuke Tsuchida

伊藤大雄*

Hiro Ito

岩間一雄*

Kazuo Iwama

1 はじめに

ゲノム解析とは、1 個体の持つ DNA 全体の塩基配列を決定することである。ゲノムの塩基数が膨大(ヒトゲノム約 30 億の塩基対からなる)である一方、シーケンサで DNA の塩基配列を高精度に解読できるのは数百塩基程度である。そこで、ゲノム全体を数百塩基程度のフラグメントに分解し、それぞれのフラグメントをシーケンサにかけ塩基配列を決定したあと、それらを繋ぎ合わせるという方法が取られる。フラグメントとは DNA の部分配列のことである。

実験により、塩基の文字列がそのフラグメントに含まれるかどうか調べることができる。この塩基の文字列のことをプローブという。例えば図 1 では、フラグメント 2 にプローブ B とプローブ G が、フラグメント 3 にプローブ G とプローブ E が存在する。よってプローブの順序は BGE(あるいは EGB) であることがわかる。

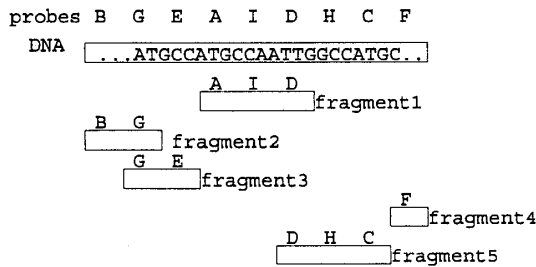


図 1: プローブとフラグメントによるゲノム解析

どのプローブがどのフラグメントに存在するかによって、プローブやフラグメントの順序が決定されていく。例えば図 1 の例では、BGE や AIDHC がそれぞれ近くにあることが決定される。しかし正確なフラグメントやプローブの順序までが必ずしも一意に決定されるわけではない。例えばプローブ A, I はプローブ D を狭んで、プローブ H, C とは反対側に存在するということは決定されるが、AI は IA であってもかまわないし、HC は CH であってもかまわない。また BGE や AIDHC や F の位置関係も決定されないし、BGE は EGB かもしれない。

フラグメントの順序を決定する問題は古くから論じられている [1]。しかしフラグメントの順序だけでなくプローブの順序までも一意に決定することができれば、その後の解析における効率化につながると考えられる。我々はプローブの順序を一意に決定するために、あとどのくらいのフラグメント数が必要かについて議論する。プローブの順序を一意に決定するために必要なフラグメント集合の性質を解析し、必要な最小フラグメント数を示した。

*京都大学情報学研究所

{tamura, tsuchida, itohiro, iwama}@kuis.kyoto-u.ac.jp

2 定式化

図 1 においてプローブを列、各フラグメントを行、フラグメントの存在するところを 1、存在しないところを 0 とすれば、(1) のような 0/1-行列で表現することができる。フラグメントの順番を決定するためには (2) のような、各行の 1 が連続している行列を導かなければならない。

$$\begin{pmatrix} A & B & C & D & E & F & G & H & I \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \quad (1)$$

$$\begin{pmatrix} B & G & E & A & I & D & H & C & F \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{pmatrix} \quad (2)$$

行列の各行の 1 を全て連続させる列の並び替えが存在するかどうかを判定し、存在する場合は、その列の並び順を示す線形時間のアルゴリズムが存在する (Kellogg ら [2])。

Kellogg ら [2] のアルゴリズムは、解さえあれば各行の 1 を連続させられる。ところが 1 が連続したからといって、プローブの順序まで一意に決定されるわけではない。(2) の行列は明らかに AI や HC を IA や CH にできるし、BGE-AIDHC-F を BGE-F-AIDHC にしてもよい。ある 0/1-行列が与えられた時にプローブの順序を一意に決定するためには、最低いくつのフラグメントを追加しなければならないだろうか。

3 解法

Kellogg ら [2] のアルゴリズムでは PQ 木と呼ばれるデータ構造が用いられる。PQ 木とは、並び替えを表すデータ構造で、P 節点と Q 節点と葉節点から構成される。P 節点は子節点の順番を任意に変更してよいことを表す節点であり、図で表すときは丸で表し、Q 節点は子節点の全体の順番を前後ひっくり返してもよいことを表す節点であり、図で表すときは矩形で表す。例えば行列 (2) を PQ 木で表現すると図 2 のようになる。

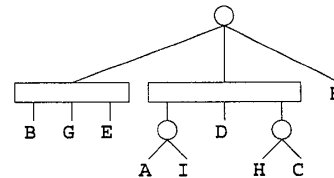


図 2: 行列 (2) に対応する PQ 木

図2のPQ木は図3のように、(A,B,C,D,E,F,G,H)に(1,0,0,1,1,0,0,1,1)と(0,0,1,1,0,1,0,1,1)という2つのフラグメントを追加すれば、木が変形しプローブの順序を一意に決定できる。1つのQ節点のみからなるPQ木を1Q木と呼ぶ。あるPQ木にフラグメント集合を追加して、プローブの順序が一意に決定できるようになれば、図3のように1Q木へと変形しているはずである。図3の変形において、1つ

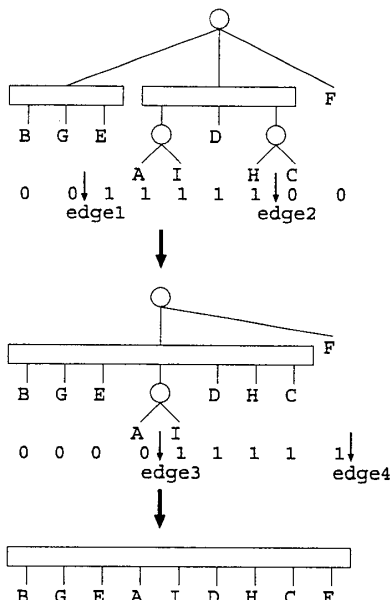


図3: 図2のPQ木の1Q木への変形

目のフラグメントではGE間とHC間、2つ目のフラグメントではAI間とFの右側に0と1の境界があると考えられる。この境界をエッジと呼ぶ。任意の隣接するプローブ間はエッジになる可能性がある。また最左のプローブの左側、最右のプローブの右側もエッジになる可能性がある。よって2つのエッジを指定すれば、1つのフラグメントを特定することができる。

与えられたPQ木を与えられた葉の順序を変えずに1Q木に変形する操作を**固定操作**と呼ぶことにする。

補題1 固定操作には以下のエッジを端にもつフラグメントの追加が必要であり、両端ともそれ以外のエッジであるフラグメントの追加は不要である。

- タイプ1: P節点の子で隣接する葉の間のエッジ。
- タイプ2: 子節点をもたない根でないQ節点の下のどこかひとつのエッジ。
- タイプ3: 与えられたPQ木の左端(右端)が根P節点を親にもつ葉である時、その左(右)側のエッジ。

例えば図3のエッジ1はタイプ2、エッジ2とエッジ3はタイプ1、エッジ4はタイプ3である。補題1で選ばれたエッジのことを**必須エッジ**と呼び、それ以外のエッジのことを**不要エッジ**と呼ぶことにする。次の補題はフラグメントの追加によって、必須エッジの数を確実に減少させることができることを示している。

補題2 間に必須エッジを含むような、2つの必須エッジに挟まれたフラグメントを追加すれば、変形後のPQ木において、それらのエッジは不要エッジになる。また変形前に不要エッジであれば、変形後も不要エッジのままである。

例えば図2において必須エッジは4つあるが、図3でまず、GE間とHC間のエッジから構成されるフラグメントが追加され、ついでAI間とFの右側のエッジから構成されるフラグメントが追加されることにより、必須エッジはすべて不要エッジとなり、1Q木が導かれている。また、それら以外のエッジが木の変形によって、必須エッジに変化することはない。

固定操作の際には、まずPQ木の葉の順番を目的とする1Q木の葉の順番と同じにしてから固定操作をほどこすものと考えてよい。葉の順番を入れ替えた後の必須エッジの数が少ないほど追加すべきフラグメント数も少ないと思われる。このことを示したのが補題3である。なお1つのP節点に着目した時、その子節点のうちPQ節点の数を v 、葉の数を l で表すことにする。

補題3 固定操作における必須エッジの最小値は以下の値の合計で求められる。

1. 子節点を持たないQ節点の数。
2. 根以外のP節点に対する、 $\max\{|l| - |v| - 1, 0\}$ の合計。
3. 根節点がP節点であれば、 $\max\{|l| - |v| + 1, 0\}$ 。

例えば図2のPQ木において、プローブFをBGEとAIDHCの間に移動させれば、必須エッジは1つ減って3個にできる。

定理1 与えられたPQ木において補題3で求められる必須エッジの数を e とする。固定操作に必要な最小フラグメント数は以下のようなものである。

1. $e \geq 3$ の時、 $\lceil \frac{e}{2} \rceil$ 個。
2. $e = 2$ の時
 - 2個の必須エッジの間の葉が1つの時、1個
 - 2個の必須エッジの間に葉が2つ以上ある時、2個

4 今後の課題

以上で最小フラグメント数が得られたが、組合せるエッジによっては、とても長いフラグメントが生成されることもある。しかし、あまりにも長いフラグメントは非現実的である。長いフラグメントを短くするためには、以下の補題を用いることができる。

補題4 長さ k のフラグメントは長さ $\lceil \frac{k+1}{2} \rceil$ の2つのフラグメントに置き換えることができる。

しかしこの補題を用いると、フラグメント数は増加してしまう。フラグメントの長さに上限がある時の最小フラグメント数は今後の課題である。

参考文献

- [1] 21世紀の医療・福祉を支える科学技術特集 電子情報通信学会誌 Vol.84 No.5 pp.341-367 2001年5月
- [2] Kellogg S. Booth and George S. Lueker. Testing for the Consecutive Ones Property, Interval Graphs, and Graph Planarity Using PQ-Tree Algorithms: Journal of Computer and System Sciences 13, pp.335-379, 1976.