

ゲノム解析における計算機科学の役割
—ショットガン・シークエンス支援システムを中心に

朝稻宏文 江口至洋
三井情報開発(株) 総合研究所

ヒト・ゲノム解析は30億塩基対にもおよぶ塩基配列の解読を目的としている。その目的を達成するためには、計算機科学が重要な役割を担っている。ここでは、それら計算機科学(Genome-informatics)の中で、特にわれわれが開発しつつあるショットガン・シークエンス支援システム（塩基配列のフラグメント・データの結合編集システム）について述べる。結合編集システムはDNAシークエンサーで読まれる多くのフラグメントを入力データとし、それらの間の重複部分配列をもとにゲノムの塩基配列を再構成するものである。大きくは、自動結合の過程と、その結果得られるフラグメントの並置を多重並置アルゴリズムによって詳細編集(refinement)する過程からなる。この方法により、高等生物に共通にみられるAlu繰り返し配列をも正しく再構成しうること、酵母第6染色体の例では、高い精度で配列決定がなされうることが示された。

The role of computer science in genome analysis
:Shotgun sequencing support system

Hirobumi Asaine Yukihiko Eguchi
Mitsui Knowledge Industry Co., Ltd., Research Institute
7-4, 3-chome Kojimachi, Chiyoda-ku, Tokyo 102, Japan

It is planned to sequence 3×10^9 base pairs in human genome project. In order to achieve this purpose, computer science plays a large part. We developed a shotgun sequencing support system. The system takes a batch of DNA fragment data supplied by a DNA sequencer, searches a sufficient length overlap between two fragment data in the batch, and joins the two data using multiple sequence alignment method subsequently. Such process continues until a whole genome sequence is reconstructed.

This process enables us to reconstruct a 36Kbp sequence of human genome with 46 Alu repeats. In case of 15Kbp sequencing of yeast chromosome VI, the process shows precision and the origin of ambiguities of the determined base sequences.

1. まえがき

ヒト・ゲノム解析は30億bpにおよぶ塩基配列を解読し、その機能をも明らかにすることを目的としている。その対象とする塩基配列の膨大さから、通常の実験室でなされている数キロbp程の塩基配列の解読とは異なった水準での計算機の利用が必要とされている。それらヒト・ゲノム解析を目的とした計算機科学はGenome-informaticsとして総称される。

Genome-informaticsは図1に示されるように、主として塩基配列からなるデータベースとそのマネージメント・システム、各種解析ソフトウェア、それらがインプリメントされているハードウェア、そして研究交流を支える通信技術からなっている。

これらの基礎技術からなるGenome-informaticsはマッピングやシークエンシング、遺伝情報解析などを適用分野としている。そのうち、われわれが開発しつつあるショットガン・シークエンス支援システム（塩基配列のフラグメント・データの結合編集システム）の概要を述べるとともに、ヒト・ゲノムや酵母染色体の塩基配列決定に適用した結果からそのシステムの評価と課題を明らかにする。

2. Genome-informaticsの範囲

Genome-informaticsの適用範囲をその構成要素であるデータベースとソフトウェアを軸に鳥瞰すると図2のようになる。ヒト・ゲノムのシークエンシングにはまず海図ともいいくべき、遺伝子地図や物理的 地図が必要とされる。物理地図の内、シークエンシングとの直接的接点をなすものは整列クローン・マップである。ヒト・ゲノムのマップ・データベースとしてはGDB (Genome Data Base) がある。なお、これらマップを効率的に作成するためにはソフトウェアが必要とされる。

マッピングされた数10Kbpの塩基配列フラグメントはある戦略に従いシークエンシングされるが、現在は数百bpのDNA断片を蛍光色素で標識し、ポリアクリルアミドゲル中を電気泳動させ、時間軸上に得られるシグナルからその塩基配列を読み取る方法が採られている。読み

取りソフトウェアはシグナルのピーク位置やその幅から塩基配列を推計するが、GC塩基を多く含む領

Genome-informatics = database
+ DBMS
+ software
+ hardware
+ communication

図1 Genome-informaticsの構造

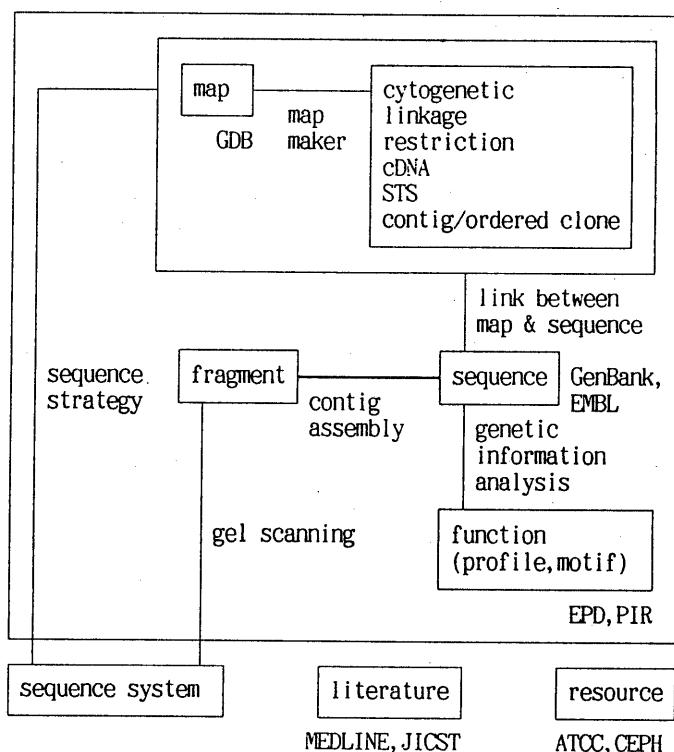


図2 ヒト・ゲノム解析を目的としたGenome-informatics

域でのバンドの圧縮現象等により、読み取り誤差が発生する。

長鎖のシークエンシング方法の1つであるショットガン法では、読み取られた数百bpの塩基配列フラグメントを繋ぎ合わせる操作 (contig assembly) が必要である。得られた塩基配列はマップ上に位置付けられ、ゲノム全体がシークエンスされるまで繰り返される。各塩基配列は GenBank や EMBL データベースとして蓄積されている。それら個々のデータの機能解析が遺伝情報解析としてなされるが、タンパク質のコーディング領域ではそのタンパク質の構造や機能の推計もなされる。この時になされる最も基本的な解析にホモジーニティ解析がある。

これらの全データベースやソフトウェアを各研究室単位に維持することは非効率的であり、コンピュータネットワークを介した共同研究が一般的になされる。

3. ショットガン・シークエンスにおける結合編集アルゴリズム

3. 1 ショットガン・シークエンス法

ショットガン法は、精度と速さの点から、長鎖の DNA をシークエンシングするのに適した方法であると考えられている (図 3)。ショットガン法では、ゲノムから切り出された 10K から 50K bp の 2 本鎖の DNA 塩基配列が超音波等によりランダムに切断され、500 から 1000 bp の多くのフラグメントが調製される。各フラグメントは入ファージ等にサブクローニングされ、その内の数百クローンがサンガー法等を用いた DNA シークエンサーによりシークエンシングされる。この時、各フラグメントの全長がシークエンシングされることではなく、シークエンシングされたフラグメントの長さは 350 から 500 bp である。このフラグメントにはサブクローニングに用いたファージ等の配列が混在していたり、3' 端のシークエンシング精度が悪い等から、結合編集操作の前段階として各フラグメントの調製段階が必要とされる。

調製されたフラグメント間の重複部分配列をもとにゲノムの塩基配列を再構成する。最初の結合過程では、2つのフラグメント間の重複部分配列のみを探索し、重複部分が充分に長くかつその領域での一致塩基比率が充分に高いという 2つの条件が満たされる場合、それらを結合し、新しいフラグメントとする。この操作を、もはやどのフラグメント間でも結合条件が満たされなくなるまで繰り返す。

結合操作が終了した段階で、決定すべき配列が全て 1 つに結合されたのか、いまだ結合は完全でなく部分的に結合された多くの配列が孤立した島のように散在しているのかの判断がなされる。後者の場合、島と島の間の塩基配列をシークエンスするためのギャップ・クロージングがなされる。その方法には 2 つ考えられる。1 つは全く新たなクローンを抽出してシークエンスする方法で、島の数が多い場合に用

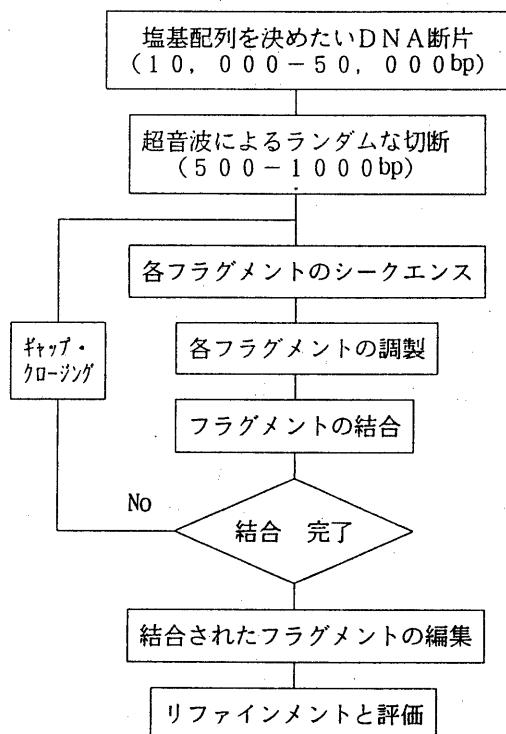


図 3 ショットガン・シークエンスの流れ

いられる。島の数が少ない場合には、各フラグメントの結合状況をみて、各島の端に位置するフラグメントを用いた合成プライマー法が用いられる。

1つに連結されたフラグメントの各塩基は、本来多くのフラグメントの同時並置を基に決定されるものである。ここでは多重並置のアルゴリズムを基に多くのフラグメントの同時並置を行い、並置された各塩基位置ごとに多数決原理に従い塩基を決定する。その後、不確定塩基やあいまい塩基のリファインメントと、各塩基ごとの決定精度に関する評価を行う。

3. 2 フラグメントの結合のアルゴリズム

2つのフラグメント間の重複部分配列の探索は、大域的ホモロジー解析〔1〕に対応する。しかし、構造や機能の予測で主に対象となる一致塩基比率30から50%のホモロジー解析とは異なり、ショットガン法では一致塩基比率80%以上の重複領域が対象である。また、ホモロジー解析の回数は対象となるフラグメントの数の2乗以上であり、できるかぎり高速な方法が求められる。

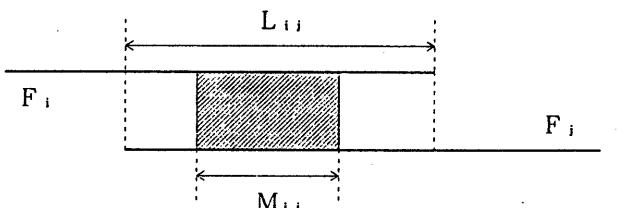
ここでは、結合条件に①完全一致の塩基配列が充分に長い、②その領域を両端に拡張した重複部分配列での一致塩基比率が充分に高い、という2つの条件を用いている（図4）。2つのフラグメント間の完全一致塩基配列を探索するアルゴリズムはDumas-Ninioらの方法〔2〕を用い、フラグメントの長さに比例する計算時間で探索している。ついで、その

塩基配列が充分に長いと、その領域を両端に拡張するが、その時DNAシークエンサーの読み取り誤差を修復するため、塩基の挿入・欠失を許した並置を行うことが必須である。ここでは、Korn-Queenの方法〔3〕を修正した方法を用いている。この結果得られた重複部分配列での一致塩基比率が充分に高いと2つのフラグメントは結合される。

42から339フラグメントを対象とした試行実験では、結合された島の数にはばらつきはあったものの、結合に要する時間はほぼフラグメント数の二乗に比例している。

3. 3 フラグメントの編集のアルゴリズム

1つに連結された塩基配列のある部分領域は、並置された多くのフラグメントによって決定されたものであるが、結合の段階では多重並置は行われていない。編集段階では、あいまいさの高い領域を限定し、そこに含まれるフラグメント間のホモロジー・スコアから系統樹を作成し、それに従い、最も近縁の2つフラグメント間で最適並置を繰り返し、全体の多重並置を構成している（図5）。なお、最終的な多重並置段階での挿入位置の整合性を保つため、ここで累進的になされる最適並置の各段階で生じた挿入は常に保存されるという規則（once a gap, always a gap rule）が適用されている〔4〕。また、系統樹を構成する段階ではUPG（unweighted pair-group clustering）法〔5〕を、最適並置では後藤の方法〔6〕を用いている。塩基対のホモロジー・スコアを示す置換パラメータには、4種の塩基を4ビット表現し、



if $M \leq M_{ij}$ and $R \leq m_{ij} / L_{ij}$, then fuse F_i and F_j ;

m_{ij} is the number of matches in the domain L_{ij} .

designate an extent

make phylogenetic tree by UPG method

consecutive alignment according to the tree

make consensus sequence

図5 多重並置法による編集

塩基間のビット単位の排他的論理和をとり、4つの各ビットのオンの数を採用した。したがって、置換パラメータは0から4の値をとっている。挿入・欠失の塩基数をkとしたときのギャップ・スコアg(k) = $\alpha k + \beta$ の最適パラメータ α と β は、塩基置換パラメータや問題となる並置によって異なるが、5.で述べる酵母染色体の例では、 $\alpha = 2$ 、 $\beta = 0.5$ を用いた。

3.4 コンセンサス配列のリファインメントと評価

多重並置法により得られたコンセンサス配列は、現在の所、単純な多数決原理（図6）に従い決定されたものであり、さらに各フラグメントを構成する塩基の削除はなされていない。そのため、より詳細な編集と、各塩基の信頼度を示す評価情報が必要とされる。

リファインメントは2段階に分けて行われる。一つは、結合編集段階で削除されずに残されてきた挿入塩基の削除であり、図6の例では9番目の塩基tがコンセンサス配列から削除される。次いで、あいまい塩基の処理がある。一つの基準として+鎖と-鎖の一一致を、

同種の鎖間に一致以上に重視すると、図6の1~4番目の塩基は、多数決原理の結果と同じGと判断される。しかし、ここでの処理方法に確定的な方法はないため、リファインメントされた各塩基の信頼度が必要とされる。

各塩基の信頼度は、

- ① 挿入塩基「-」やあいまい塩基「n」、「y」等の評価
- ② +鎖と-鎖の組み合わせの評価
- ③ 5'端か3'端か、等の位置情報（個々の塩基の信頼性）の評価
- ④ 各フラグメント・データそのものの信頼性の評価

に依存するが、単純にはその塩基決定に参加したフラグメントの数nと、その内多数を占めた塩基の数mの関数と考えられる。

ここでは、nとmの関数としての信頼度Rとして $R(m, n) = m^n / n$ なる関数形を考え、Rが満たすべき条件としては、

- ① $R(1, 1) = 1$
- ② $R(i, i) \leq R(j, j)$, if $i \leq j$
- ③ $R(n/2, n) < 1$

が、20以下のnについて成立とした。これらの条件を満たすaは $a < 1.301$ であるため、以下では $a = 1.3$ なる関数Rで信頼度を表現することとした。

4. ヒト・ゲノムのシークエンシング

開発された結合用ソフトウェアを評価するため、染色体凝縮に関連する遺伝子を含むヒト・ゲノム約36Kbpの評価用ショットガン・シークエンス実験を行った。ヒト・ゲノムでは多くの繰り返し配列がみられるが、それらの配列を正しく再構成しうるかが課題である。

用いたフラグメント・データは551個、平均鎖長304bpである。シークエンス実験の結果を表1に示す。ケース1から5の実験では、結合条件を

$$M \geq 50 \text{ bp}, \quad R \geq 85\%$$

	1	5	10
consensus sequence	GAYkCCmATTTGG		
fragment 1 (+)	ggc--ccca-ttt		
fragment 2 (+)	gat-ccca-tttgg		
fragment 3 (-)	-c--nca-ttttg		
fragment 4 (+)	-a-g-taa-tttgt		
fragment 5 (+)	t--t-a		

図6 コンセンサス配列決定に用いられる多数決原理
多数決原理を主に示すため、最適な多重並置から変形した並置となっている。
()内の記号は+鎖か、-鎖かを示す。

とし、ランダムに100個のフラグメントを抽出し、以前の実験結果に追加するかたちで結合操作を行った。ケース7の実験では新たなフラグメントの追加ではなく、結合条件を $M \geq 50\text{ bp}$ 、 $R \geq 80\%$ とゆるめることにより、島の数を78から68に減少させていく。

ケース7の実験結果では全てが完全に結合されてはいないが、ギャップ・クロージング実験を踏まえた詳細編集の結果との比較から、

- ①ヒト・ゲノムなど高等生物に共通にみられるA1u繰り返し配列（約280bp）が46個含まれているが、それらをも正しく再構成できている。
- ②ヒト・ゲノム以外のプラスミドDNAフラグメントは、ヒト・ゲノムと混在し島を形成することなく、正しく別個の島として再構成されている。

ことが明らかにされた。この結果、2つのフラグメントの結合アルゴリズム（図4）が、繰り返し配列を多く含む高等生物のシークエンシングに適用しうることが確認された。なお、上記①の結果を得るためにには結合条件を当初は厳しく（ $R \geq 85\%$ ）することが必須条件であった。

5. 酵母染色体ゲノムのシークエンシング

酵母第6染色体のうち約15Kbpのショットガン・シークエンシングを行った。得られたフラグメントは156個であるが、その内の5つのフラグメントは合成プライマー法により伸長反応を行った結果である。平均鎖長400bp、全体の塩基数は62Kbpであるので、約4倍量のシークエンスを行

(a) 自動結合により得られたコンセンサス配列

13230	13290
consensus TGAwtCAAGAmtoCGCGACTCTAGCCGCGkTCGgAACATGTCAGATACTTGGCGTCTGCGgAAAGgTT	
(-)06115 tgat-caagaa--cgcgactctagcgcggtcg-aacatgtcaga	
(+)06072 tgaatcaagaa-cccgactctagcgcggtcggaaacatgtcagatacttggcgtcgtgcggaaaaggtt	
(+)06114 tgat-caagaa-cccgactctagcgcgt-cggaaacatgtcagatacttggcgtcgtgcg-aaag-tt	
(-)06193 tgattcaagactccg--actctagcgcgtt	
(-)06160 tg-atcaagaa-acgcgactctagcgcgg-gcgaacatgtcagatacttggcgtcgtgc-gaaa-gtt	
(-)06096 agcgcggtc-gaacatgtcagatacttgcgcgtcgtgc-gaaa-gtt	

(b) 多重並置と多数決原理によって得られたコンセンサス配列

13230	13290
consensus TCAATCAAGAACCGGACTCTAGCCGCGTCGGAACATGTCAGATACTTGGCGTCTGCGgAAAGgTT	
(-)06115 tg-atcaagaa-cgcgactctagcgcggtc-gaacatgtcaga	
(+)06072 tgaatcaagaaccgcgactctagcgcggtcggaaacatgtcagatacttggcgtcgtgcggaaaaggtt	
(+)06114 tg-atcaagaaccgcgactctagcgcg-tcggaaacatgtcagatacttggcgtcgtgcg-aaag-tt	
(-)06193 tgattcaagactc-cgactctagcgcgtt	
(-)06160 tg-atcaagaaacgcgactctagcgcggc-gaacatgtcagatacttggcgtcgtgcg-aaag-tt	
(-)06096 agcgcggtc-gaacatgtcagatacttgcgcgtcgtgc-gaaa-gtt	

図7 結合・編集により得られたコンセンサス配列

ったことになる。

156個のフラグメントを結合した結果と、そのあと多重並置し、多数決原理でコンセンサス配列を得た結果を図7に示す。自動結合だけから得られた結果に比べ、多重並置により得られたコンセンサス配列では、あいまい塩基(wやm、k)が正しく補正されている、並置段階で生じた余分な塩基が削除されているといった効果により、信頼度が増加していることがわかる。なお、156個のフラグメントの自動結合に要した時間はSUN4で、約2時間であった。

この段階まではフラグメントの各塩基を削除する操作を行っていないため、図7(b)では、詳細編集の段階では削除すべき2つの塩基gが小文字で示されている。そのような塩基を全配列にわたって削除する必要がある。ここでは、いかなる塩基を削除し、いかなる塩基をあいまい塩基とするかについての、シークエンシングをおこなっている研究者の判断を得た。削除すべき塩基をdとすると、その多重並置でのパターンは「d - - -」か「d d - - -」かのどちらかであった。dの個数をm、挿入「-」の個数とmの和をnとして、各パターンの出現頻度を表2に示す。この結果、決定された全塩基数は14,870bpであった。この内、多数決原理だけではあいまいさの残る塩基を並置のパターンごとに分類すると、表3のごとくなる。

この結果、あいまい塩基の比率は3.7%となる。そのあいまいさの大部分は、ある領域が1つのフラグメントだけで決定されていることからきている。その領域によるあいまいさは2.3%であり、次いで「xxy」の0.4%、「x-」の0.3%となっている。

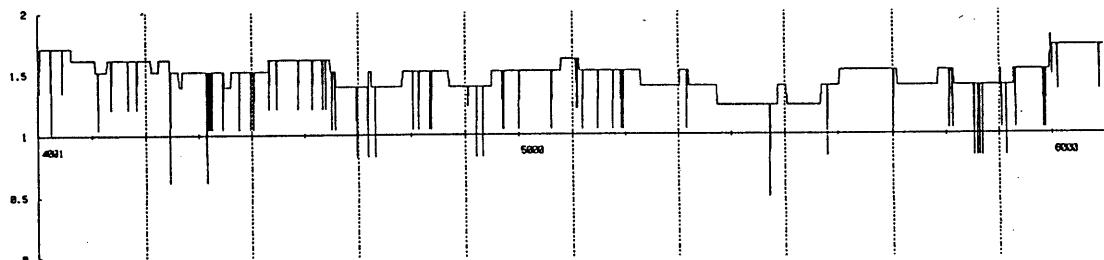


図8 シークエンスされた配列の信頼度関数

表2 削除される塩基の並置パターン

n	m		計
	1	2	
3	3 6	-	3 6
4	5 3	-	5 3
5	3 1	-	3 1
6	2 8	2	3 0
7	1 4	7	2 1
8	5	-	5
その他	-	-	6
計			1 8 2

表3 あいまいさの残る塩基の並置のパターン

並置のパターン	塩基数
x (1つのフラグメント)	3 4 8
x -	4 3
x y	2 8
x n	4
x x y	5 7
x x -	2 3
x x n	6
x x y y	4
x x - -	4
x x - - -	6
x x x y -	4
その他	3 0
計	5 5 7

x, y, zはある塩基を、nは任意の塩基を、-は挿入を示す。

表3に示されたあいまいさの基準を信頼度関数Rでみると、 $R \leq 1$ なる並置に位置する塩基があいまいと評価されている。得られた14, 870 bpの塩基配列の信頼度Rを図8に示すが、 $R \leq 1$ なる領域があいまい領域と定義されうる。

6. おわりに

genome-informaticsは現在、ヒト・ゲノム解析のボトル・ネックになっているとの見解〔7〕もあり、世界的に研究の進められている分野である。その範囲は広く、1研究室や1研究グループだけでカバーしうるものではない。今後、わが国においても、より多くの研究者の協力が図られることを期待する。

最後に本研究は理化学研究所のヒト・ゲノム解析システム・プロジェクトの一環として進められていくものであり、特に分子生物学の分野からご指導を頂いている添田栄一博士、村上康文博士に感謝いたします。

文 献

- 〔1〕江口至洋「タンパク質工学の物理・化学的基礎」、共立出版(1991)
- 〔2〕Dumas, J-P., Ninio, J.: Nucleic Acids Res., 10, 197(1982)
- 〔3〕Korn, L. J., Queen, C. L., Wegman, K. N.: Proc. Natl. Acad. Sci. USA, 74, 4401(1977)
- 〔4〕Feng, D. F., Doolittle, R. F.: J. Mol. Evol., 25, 351(1987)
- 〔5〕根井正利「分子進化遺伝学」(五條堀、斎藤訳)、培風館(1990)
- 〔6〕Gotoh, O.: J. Mol. Biol., 162, 705(1982)
- 〔7〕Mavournin, K. H., Mansfield, B. K.: Human Genome news, 2, 8(1990)