

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2013年10月31日(31.10.2013)



(10) 国際公開番号  
WO 2013/162010 A1

- (51) 国際特許分類:  
G06F 19/22 (2011.01) C12Q 1/68 (2006.01)
- (21) 国際出願番号: PCT/JP2013/062426
- (22) 国際出願日: 2013年4月26日(26.04.2013)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:  
特願 2012-101755 2012年4月26日(26.04.2012) JP
- (71) 出願人: 独立行政法人放射線医学総合研究所  
(NATIONAL INSTITUTE OF RADIOLOGICAL SCIENCES) [JP/JP]; 〒2638555 千葉県千葉市稲毛区穴川4丁目9番1号 Chiba (JP). 株式会社メイズ  
(MAZE, INC.) [JP/JP]; 〒1640011 東京都中野区中央3丁目13番11号 MGビル508 Tokyo (JP).
- (72) 発明者: 安倍 真澄(ABE, Masumi); 〒2638555 千葉県千葉市稲毛区穴川4丁目9番1号 独立行政法人放射線医学総合研究所内 Chiba (JP). 笠間康次(KASAMA, Yasuji); 〒2638555 千葉県千葉市

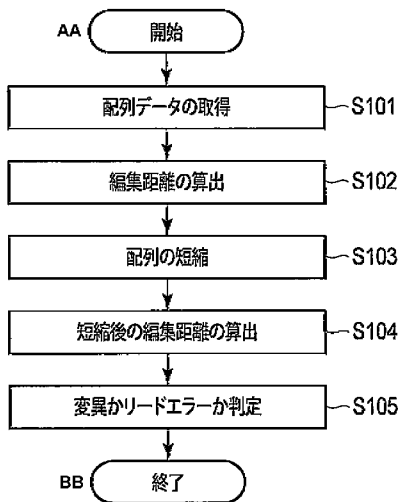
稲毛区穴川4丁目9番1号 独立行政法人放射線医学総合研究所内 Chiba (JP). 湯野川 春信(YUNOKAWA, Harunobu); 〒1640011 東京都中野区中央3丁目13番11号 MGビル508 株式会社メイズ内 Tokyo (JP). 佐藤 伸司(SATO, Shinji); 〒2994336 千葉県長生郡長生村岩沼918番2 株式会社メイズ内 Chiba (JP). 近藤 一弘(KONDO, Kazuhiro); 〒1640011 東京都中野区中央3丁目13番11号 MGビル508 株式会社メイズ内 Tokyo (JP). 日永田 隆志(HIEIDA, Takashi); 〒1640011 東京都中野区中央3丁目13番11号 MGビル508 株式会社メイズ内 Tokyo (JP).

- (74) 代理人: 蔵田 昌俊, 外(KURATA, Masatoshi et al.); 〒1050001 東京都港区虎ノ門1丁目12番9号 鈴榮特許総合事務所内 Tokyo (JP).
- (81) 指定国 (表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS,

[続葉有]

(54) Title: METHOD FOR DETERMINING READ ERROR IN NUCLEOTIDE SEQUENCE

(54) 発明の名称: 塩基配列のリードエラーを判定する方法



AA START  
 S101 OBTAIN SEQUENCE DATA  
 S102 CALCULATE EDIT DISTANCE  
 S103 SHORTEN SEQUENCES  
 S104 CALCULATE EDIT DISTANCE AFTER SHORTENING SEQUENCES  
 S105 DETERMINE WHETHER DIFFERENCE COMES FROM MUTATION OR READ ERROR  
 BB END

(57) Abstract: A method for determining whether the difference between a first reference sequence and a first sequence of interest which has homology with the first reference sequence comes from a mutation occurring in the first sequence of interest or a read error occurring during sequencing. The method comprises substituting a sequence composed of a predetermined number or more of contiguous nucleotides in each of the first sequence of interest and the first reference sequence by a sequence composed of the predetermined number of contiguous nucleotides to thereby produce a second sequence of interest and a second reference sequence (S103), wherein all of the nucleotides are the same as each other. The method additionally comprises: calculating the edit distance of the second sequence of interest relative to the second reference sequence (S104); and determining whether the difference comes from the abovementioned mutation or the above-mentioned read error on the basis of the edit distance (S105).

(57) 要約: 方法は、第1のリファレンス配列と、前記第1のリファレンス配列と相同性を有する第1の対照配列との差異が、前記第1の対照配列の変異によるものかシーケンシングのリードエラーによるものかを判定する方法である。方法は、前記第1の対照配列及び前記第1のリファレンス配列の同一塩基が所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して第2の対照配列及び第2のリファレンス配列を作成すること(S103)を具備する。方法は、さらに、前記第2のリファレンス配列に対する前記第2の対照配列の編集距離を算出すること(S104)と、前記編集距離に基づいて、前記差異が前記変異によるものであるか前記リードエラーによるものであるかを判定すること(S105)とを具備する。



WO 2013/162010 A1



LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

パ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

(84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ

添付公開書類:

— 国際調査報告 (条約第 21 条(3))

## 明 細 書

**発明の名称**：塩基配列のリードエラーを判定する方法

### 技術分野

[0001] 本発明は、塩基配列のリードエラーを判定する方法に関する。

### 背景技術

[0002] 一塩基多型（SNP）、複数塩基の置換、塩基の挿入又は欠失（Insertion/Deletion）等の変異の特定については、診断や治療、ウィルスや細菌の同定、家畜等の遺伝情報の解析等、研究及び臨床等の場において多くの需要が存在する。例えば国際公開第2006/110855号には、ピロリン酸配列決定技法を利用したシーケンシングによって、効率的にSNP等の変位を特定する技術が開示されている。

### 発明の概要

[0003] 例えば国際公開第2006/110855号に開示されているようなピロリン酸配列決定技法等を利用したシーケンシングにより取得された塩基配列データには、同一塩基が連続する部分において、同一塩基の連続数が実際と異なるリードエラーが含まれることがある。リファレンス配列とこのリファレンス配列と比較される対照配列との差異が、変異によるものなのかシーケンシングにおけるリードエラーによるものなのかを判定することは困難である。

[0004] そこで本発明は、塩基配列の差異が変異によるものなのかリードエラーによるものなのかを判定するための方法を提供することを目的とする。

[0005] 前記目的を果たすため、本発明の一態様によれば、塩基配列のリードエラーを判定する方法は、第1のリファレンス配列と、前記第1のリファレンス配列と相同性を有する複数の第1の対照配列との差異が、前記第1の対照配列の変異によるものかシーケンシングのリードエラーによるものかを判定する方法であって、各々の前記第1の対照配列の同一塩基が所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して各々の

第2の対照配列を作成することと、前記第1のリファレンス配列の同一塩基が前記所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して第2のリファレンス配列を作成することと、前記第2のリファレンス配列に対する各々の前記第2の対照配列の編集距離である複数の短縮後編集距離を算出することと、複数の前記短縮後編集距離の平均値である短縮後平均値を算出することと、前記短縮後平均値に基づいて、前記差異が前記変異によるものであるか前記リードエラーによるものであるかを判定することと、を具備する。

[0006] また、前記目的を果たすため、本発明の一態様によれば、塩基配列のリードエラーを判定する方法は、第1のリファレンス配列と、前記第1のリファレンス配列と相同性を有する第1の対照配列との差異が、前記第1の対照配列の変異によるものかシーケンシングのリードエラーによるものかを判定する方法であって、前記第1の対照配列の同一塩基が所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して各々の第2の対照配列を作成することと、前記第1のリファレンス配列の同一塩基が前記所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して第2のリファレンス配列を作成することと、前記第2のリファレンス配列に対する前記第2の対照配列の編集距離を算出することと、前記編集距離に基づいて、前記差異が前記変異によるものであるか前記リードエラーによるものであるかを判定することと、を具備する。

[0007] 本発明によれば、塩基配列の差異が変異によるものなのかリードエラーによるものなのかを判定するための方法を提供できる。

### 図面の簡単な説明

[0008] [図1]図1は、本発明の各実施形態に係る判定装置の構成例を示すブロック図である。

[図2]図2は、本発明の各実施形態に係る判定装置の機能ブロックの構成例を示すブロック図である。

[図3]図3は、第1の実施形態に係る判定装置による処理の一例を示すフロー

チャートである。

[図4]図4は、第2の実施形態に係る判定装置による処理の一例を示すフローチャートである。

[図5]図5は、第2の実施形態の変形例に係る判定装置の機能ブロックの構成例を示すブロック図である。

[図6]図6は、HiCEPを説明するための模式図である。

[図7]図7は、第3の実施形態に係るクラスタ決定装置による処理の一例を示すフローチャートである。

### 発明を実施するための形態

#### [0009] [第1の実施形態]

本発明の第1の実施形態について図面を参照して説明する。本実施形態に係る塩基配列の差異が変異によるものなのかリードエラーによるもののかを判定するため判定装置1の構成例の概略を図1に示す。判定装置1は、マザーボード11と、CPU12と、メインメモリ13と、ハードディスクドライブ(HDD)14と、入力装置15と、出力装置16と、記録媒体インターフェース(媒体I/F)17と、通信インターフェース(通信I/F)18とを備える。

[0010] 判定装置1内の各部分は、マザーボード11を介して通信を行う。CPU12は、種々の演算を行う。メインメモリ13は、例えばDRAMを含み、CPU12における演算に必要な情報等を一時記憶する。HDD14は、CPU12によって行われる演算に係るプログラム等を含む各種情報を記憶している。入力装置15は、例えばキーボードやマウス等のユーザからの入力に係るユーザインターフェースを含む。出力装置16は、例えばディスプレイやプリンタ等のユーザへの出力に係るユーザインターフェースを含む。媒体I/F17は、記録媒体101との通信を行うインターフェースである。通信I/F18は、ネットワーク102と接続するためのインターフェースである。

[0011] 本実施形態に係る判定装置1の機能ブロックの構成例の概略を図2に示す

。判定装置 1 は、CPU 1 2 及びメインメモリ 1 3 や HDD 1 4 に記録されたプログラム等によって構成される演算部 3 0 を有する。演算部 3 0 は、データ取得部 3 1 と、配列短縮部 3 2 と、編集距離算出部 3 3 と、エラー判定部 3 4 とを有する。演算部 3 0 は、記録部 4 1 やネットワーク 4 2 や出力部 4 3 と接続している。

[0012] データ取得部 3 1 は、HDD 1 4、記録媒体 1 0 1 等で構成される記録部 4 1 や、ネットワーク 4 2 から塩基配列に係るデータを取得する。このデータには、リファレンス配列と、このリファレンス配列と差異を有し、その差異が変異によるものなのかリードエラーによるものなのかが判定されるべき対照配列とが含まれている。データ取得部 3 1 が取得したリファレンス配列を第 1 のリファレンス配列と称し、対照配列を第 1 の対照配列と称することにする。

[0013] 配列短縮部 3 2 は、第 1 のリファレンス配列について、同一塩基が連続して 3 つ以上並ぶ配列をその塩基が 2 つ連続する配列に置換して、短縮された塩基配列を作成する。この短縮された塩基配列を第 2 のリファレンス配列と称することにする。配列短縮部 3 2 は、同様に第 1 の対照配列について同一塩基が連続して 3 つ以上並ぶ配列をその塩基が 2 つ連続する配列に置換して、短縮された配列を作成する。この短縮された配列を第 2 の対照配列と称することにする。

[0014] 編集距離算出部 3 3 は、第 1 のリファレンス配列に対する各々の第 1 の対照配列の編集距離を算出する。同様に、編集距離算出部 3 3 は、短縮された第 2 のリファレンス配列に対する短縮された各々の第 2 の対照配列の編集距離を算出する。エラー判定部 3 4 は、編集距離算出部 3 3 が算出した編集距離に基づいて、第 1 のリファレンス配列と第 1 の対照配列との差異が、変異によるものなのかリードエラーによるものなのかを判定する。

[0015] 本実施形態に係る判定装置 1 は、次のような処理を行う。例えばパイロシーケンス（登録商標）法等のピロリン酸配列決定技法を用いたシーケンシングにおいては、同一塩基が連続する場合に連続数を誤るリードエラーが含ま

れやすい。このようなエラーが含まれる問題をいわゆるホモポリマー問題という。判定装置 1 は、第 1 のリファレンス配列と、第 1 の対照配列との差異が、ホモポリマー問題に係るリードエラーによるものなのか、変異が存在することによるものなのかを判定する。

[0016] 本実施形態に係る判定装置 1 の動作を説明する。本実施形態において実行される処理のフローチャートを図 3 に示す。演算部 30 のデータ取得部 31 は、ステップ S101 において、配列データを取得する。配列データは、第 1 のリファレンス配列と第 1 の対照配列とを含む。ここで、配列データに含まれるのは、第 1 のリファレンス配列に対して差異がある第 1 の対照配列である。この差異は公知の手法によって特定されたものであるが、この差異がシーケンシングにおけるリードエラーなのか、変異が存在することによる差異なのかが明らかになっていない。第 1 のリファレンス配列と第 1 の対照配列とは、この差異を含む所定の塩基数の配列である。なお、第 1 のリファレンス配列には、トランスクリプトの配列やゲノム配列等を含む種々の配列が用いられ得る。配列データは、記録媒体 101 から取得してもよいし、ネットワーク 102 から取得してもよい。また、HDD 14 に記録されている第 1 のリファレンス配列と第 1 の対照配列とのデータを取得してもよい。

[0017] 本説明では、例えば第 1 のリファレンス配列が「AGCCTTTA」（以降、配列 1 と称する）であり、第 1 の対照配列が「AGCCTTTTA」（以降、配列 2 と称する）であるものとする。すなわち、配列 1 と配列 2 とでは、連続する「T」の数が異なる。

[0018] 演算部 30 の編集距離算出部 33 は、ステップ S102 において、第 1 のリファレンス配列に対する第 1 の対照配列の編集距離を算出する。ここで、編集距離とは、リファレンス配列と対象とする配列とを比較し、塩基の挿入、削除又は置換の有無を求め、挿入、削除又は置換のうち何れか 1 つが存在する毎にそれぞれ「1」を加算して求められる値である。例えば第 1 のリファレンス配列「AGCCTTTA」（配列 1）に対する第 1 の対照配列「AGCCTTTTA」（配列 2）の編集距離は、T が 1 つ挿入されているので

1である。

[0019] 演算部30の配列短縮部32は、ステップS103において、第1のリファレンス配列及び第1の対照配列について、同一塩基が3塩基以上連続する配列をその塩基が2塩基連続した配列に変換して、短縮配列を作成する。例えば、配列1は「AGCCTTA」（以降、配列3と称する）となり、配列2は「AGCCTTA」（以降、配列4と称する）となる。第1のリファレンス配列の短縮配列を第2のリファレンス配列と称し、第1の対照配列の短縮配列を第2の対照配列と称することにする。

[0020] 演算部30の編集距離算出部33は、ステップS104において、第2のリファレンス配列に対する第2の対照配列の編集距離を算出する。配列3に対する配列4の編集距離は、両配列が一致しているので0となる。

[0021] 演算部30のエラー判定部34は、ステップS105において、短縮前の編集距離と短縮後の編集距離とに基づいて、第1のリファレンス配列と第1の対照配列との差異が変異によるものなのかリードエラーによるものなのかを判定する。例えば、エラー判定部34は、短縮前の編集距離と短縮後の編集距離とが同一であるとき差異が変位によるものであると判定し、短縮前の編集距離と短縮後の編集距離とが相違するとき差異がリードエラーによるものであると判定する。また、短縮後の編集距離が0であるとき、差異がリードエラーによるものであると判定し、短縮後の編集距離が0以外であるとき、差異が変異によるものであると判定してもよい。配列1と配列2との例では、短縮前の編集距離が1であり、短縮後の編集距離が0であるので、配列1と配列2との差異はリードエラーによるものであると判定される。演算部30は、判定結果をHDD14や出力装置16や記録媒体101やネットワーク102等を含む出力部43に出力し、記録等させる。その後、処理は終了する。

[0022] ホモポリマー問題に係るリードエラーは、同一塩基が2つ連続で出現してもほとんど発生せず、同一塩基が3つ以上連続して出現すると、連続数が多い程、エラーの発生確率が上昇する。そこで、本実施形態では、同一塩基が



連続して出現するとき、その連続数を全て2つに短縮することで、解析においてシーケンシングにおけるリードエラーの影響を受けないようにしている。すなわち、ホモポリマー問題に起因するシーケンシングにおけるリードエラーがある場合には短縮されることで配列の相違がなくなり、リードエラーがなく変異が存在する場合には短縮されることで相違が維持されることが利用されている。

[0023] 第1の対照配列が別の場合の例を示す。第1のリファレンス配列が配列1であり、第1の対照配列が「AGCCGTTA」（以降、配列5と称する）であるとする。このとき、ステップS102で算出される配列1に対する配列5の編集距離は、「T」の1つが「G」に置換されているので、1となる。ステップS103において配列5が短縮されて得られる第2の対照配列は、「AGCCGTTA」（以降、配列6と称する）となる。ステップS104で算出される配列1に対する配列6の編集距離は1となる。したがって、ステップS105において、変異が存在すると判定される。

[0024] このように本実施形態によれば、ホモポリマー問題を含む配列データの解析において、ホモポリマーが短縮され、短縮された配列の編集距離に基づいて、配列の差異が変異によるものなのかリードエラーによるものなのかが判定される。すなわち、本実施形態によれば、SNP等が正確に同定され得る。

[0025] なお、図3を参照して説明した処理は一例であり、各処理の順序は変更され得るし、一部の変更や省略もされ得る。また、本実施形態では、同一塩基が3塩基以上連続する配列を2塩基に短縮しているが、いくつに短縮するようによい。例えば同一塩基が4塩基以上連続する配列を3塩基に短縮するようによい。ただし、3塩基以上連続する場合にリードエラーの発生確率が上昇するので、同一塩基が3塩基以上連続する配列を2塩基に短縮することが好ましい。

[0026] [第2の実施形態]

本発明の第2の実施形態について説明する。ここでは、第1の実施形態と

の相違点について説明し、同一の部分については、同一の符号を付してその説明を省略する。ピロリン酸配列決定技法を用いたシーケンシングでは、一度に大量の配列データが得られる。本実施形態では、リファレンス配列と同一性があるとされた一群の配列を配列データとして取り扱う。本実施形態に係る判定装置 1 は、この一群の配列のうちリファレンス配列と差異がある部分について、その差異がリードエラーによるものなのか、ヘテロ接合型の S N P 等によるものなのかを判定する。ここで S N P 等とは、1 塩基の置換、挿入又は欠失をいう。ホモポリマー問題のように、リードエラーがランダムではなく一定の条件で発生する場合、従来の判定手法ではそのリードエラーが S N P 等と判定される可能性がある。本実施形態では、このようなホモポリマー問題を含むリードエラーを判別する。なお、ヘテロ接合型の S N P 等が存在する場合、理想的には、リファレンス配列に対する各配列の編集距離の平均値は 0.5 となり、2 種類の配列が 50% ずつ存在することになる。

[0027] 本実施形態に係る判定装置 1 の動作を説明する。本実施形態において実行される処理のフローチャートを図 4 に示す。演算部 30 のデータ取得部 31 は、ステップ S 201 において、リファレンス配列と同一性があるとされた一群の配列のうち、リファレンス配列と差異がある部分の配列を配列データとして取得する。リファレンス配列には、トランスクリプトの配列やゲノム配列等を含む種々の配列が用いられ得る。配列データの一例を表 1 に示す。表 1 の最上段に第 1 のリファレンス配列としてのリファレンス配列が示されている。その下には、第 1 の対照配列としての配列が示されている。これら配列を配列 7 乃至配列 14 と称することにする。配列 7 乃至 14 のうち、リファレンス配列と異なる部分は、シーケンシングにおけるリードエラーか S N P 等によるものか判別されていない差異である。

[表1]

表 1

リファレンス配列	AGCCT-TTA	編集距離
配列 7	AGCCT-TTA	0
配列 8	AGCCT-TTA	0
配列 9	AGCCT-TTA	0
配列 10	AGCCTTTTA	1
配列 11	AGCCTTTTA	1
配列 12	AGCC--TTA	1
配列 13	AGCC--TTA	1
配列 14	AGCCG-TTA	1

[0028] 演算部 30 の編集距離算出部 33 は、ステップ S202 において、第 1 のリファレンス配列に対する各第 1 の対照配列の編集距離を算出する。例えば表 1 に示す各配列について編集距離を算出すると次のようになる。配列 7 乃至 9 は、それぞれコンセンサス配列と同一であるので、それぞれ編集距離は 0 である。配列 10 及び 11 は、コンセンサス配列に対して T が 1 つ挿入されているので編集距離は 1 である。配列 12 及び 13 は、コンセンサス配列に対して T が 1 つ削除されているので、編集距離は 1 である。配列 14 は、コンセンサス配列に対して T が 1 つ G に置換されているので、編集距離は 1 である。演算部 30 の編集距離算出部 33 は、ステップ S203 において、全ての第 1 の対照配列の編集距離の平均値を算出する。

[0029] 演算部 30 の配列短縮部 32 は、ステップ S204 において、第 1 のリファレンス配列及び各第 1 の対照配列について、同一塩基が 3 塩基以上連続する配列をその塩基が 2 塩基連続した配列に変換して、短縮配列を作成する。例えば表 1 に示した第 1 の対照配列としての配列 7 乃至 14 は、表 2 に示すような短縮された第 2 の対照配列としての配列 15 乃至 22 に各々変換される。

[表2]

表 2

リファレンス配列	AGCC-TTA	編集距離
配列 1 5	AGCC-TTA	0
配列 1 6	AGCC-TTA	0
配列 1 7	AGCC-TTA	0
配列 1 8	AGCC-TTA	0
配列 1 9	AGCC-TTA	0
配列 2 0	AGCC-TTA	0
配列 2 1	AGCC-TTA	0
配列 2 2	AGCCGTTA	1

[0030] 演算部 3 0 の編集距離算出部 3 3 は、ステップ S 2 0 5 において、第 2 のリファレンス配列に対する各第 2 の対照配列の編集距離を算出する。この編集距離は、例えば表 2 に示すようになる。演算部 3 0 の編集距離算出部 3 3 は、ステップ S 2 0 6 において、第 2 の対照配列の編集距離の平均値を算出する。

[0031] 演算部 3 0 のエラー判定部 3 4 は、ステップ S 2 0 7 において、第 2 の対照配列の編集距離の平均値が所定の範囲内であるか否かによって、第 1 のリファレンス配列と第 1 の対照配列との差異が、SNP 等によるものである可能性があるか、リードエラーによるものであるかを判定する。例えば表 3 に示す判定基準に基づいて判定される。すなわち、第 2 の対照配列（短縮後）の編集距離の平均値が 0.25 以上 0.75 以下である場合、SNP 等によるものである可能性があるとして判定される。これは、2 倍体において、対立遺伝子の一方に SNP 等が存在するヘテロ接合の場合には、理想的には編集距離は 0.5 となることに基づく。なお、本実施形態では、編集距離の平均値の範囲を 0.25 以上 0.75 以下と 0.5 を中心とする値に設定しているが、0.5 を含む他の値としてもよい。

[0032] また、後に詳述する理由により、第 2 の対照配列（短縮後）の編集距離の

平均値が0.75より大きいときでも第1の対照配列（短縮前）の編集距離の平均値が0.25以上0.75以下である場合、SNP等によるものである可能性があるとして判定される。それ以外の場合、リードエラーによるものであるとして判定される。ステップS207においてSNP等によるものである可能性があるとして判定されたとき、処理はステップS208に進む。リードエラーによるものであるとして判定されたとき、処理はステップS210に進む。

[表3]

表 3

		短縮前		
		編集距離<0.25	0.25≤編集距離≤0.75	0.75<編集距離
短縮後	編集距離<0.25	リードエラー	リードエラー	リードエラー
	0.25≤編集距離≤0.75	SNP	SNP	SNP
	0.75<編集距離	リードエラー	SNP	リードエラー

[0033] 演算部30のエラー判定部34は、ステップS208において、SNP等が存在する可能性がある配列の数の割合が所定の範囲内であるか否かを判定し、第1のリファレンス配列と第1の対照配列との差異が、SNP等によるものである可能性があるか、リードエラーによるものであるかを判定する。例えば、SNP等が存在する可能性がある配列の数が、全配列数の33%以上67%以下のとき、差異はSNP等によるものであるとして判定される。それ以外の場合、差異はリードエラーによるものであるとして判定される。この判定は、第2の対照配列を同一配列ごとにグループ分けした場合、最も多くの配列が含まれるグループの配列数が全配列数の66%以下であり、かつ、2番目に多くの配列が含まれるグループの配列数が全配列数の33%以上である場合に、差異はSNP等によるものであるとして判定されるように変更されてもよい。

[0034] 差異はSNP等によるものであるとして判定されたとき、処理はステップS209に進む。差異はリードエラーによるものであるとして判定されたとき、処理はステップS210に進む。すなわち、ステップS207の判定の条件とス

ステップS 208の判定の条件とを共に満たすとき、処理はステップS 209に進み、それ以外の場合、処理はステップS 210に進む。

[0035] 演算部30のエラー判定部34は、ステップS 209において、SNP等が存在する配列を特定する。演算部30は、SNP等が存在する配列に係る情報を出力部43に出力し、記録等させる。その後処理は終了する。

[0036] 一方、演算部30のエラー判定部34は、ステップS 210において、第1のリファレンス配列と第1の対照配列との差異はリードエラーによるものであると結論付け、その旨を出力部43に出力し、記録等させる。その後処理は終了する。

[0037] 表1に示した第1のクラスタの例では、ステップS 205において、短縮後の各第2の対照配列の編集距離は表2に示すようになる。この場合、ステップS 206において算出される短縮後の第2の対照配列の編集距離の平均値は、0.125である。したがって、ステップS 207の判定において、リードエラーが存在すると判定され、ステップS 210においてリードエラーが存在すると結論付けられ、処理は終了する。

[0038] 表1に示す例のような場合、本実施形態によれば、配列短縮部32が配列を短縮することで、ホモポリマー問題に起因するエラーが取り除かれ、第1のリファレンス配列と第1の対照配列との差異はリードエラーによるものであると結論付けられる。

[0039] 一群の第1の対照配列の別の例を表4に示す。この場合、ステップS 202で算出される編集距離は表4に示すとおりであり、ステップS 203で算出される編集距離の平均値は0.625である。

[表4]

表4

リファレンス配列	AGCC-TTTA	編集距離
配列23	AGCC-TTTA	0
配列24	AGCC-TTTA	0
配列25	AGCCTTTTA	1
配列26	AGCCTTTTA	1
配列27	AGCC--TTA	1
配列28	AGCCGTTTA	1
配列29	AGCC-GTTA	1
配列30	AGCC-GTTA	1

[0040] ステップS204で短縮される結果、配列23乃至30は、それぞれは表5に示す第2の対照配列としての配列31乃至38のようになる。この場合、ステップS205で算出される短縮後の編集距離は表5に示すとおりであり、ステップS206で算出される短縮後の編集距離の平均値は0.375である。

[表5]

表5

リファレンス配列	AGCC-TTA	編集距離
配列31	AGCC-TTA	0
配列32	AGCC-TTA	0
配列33	AGCC-TTA	0
配列34	AGCC-TTA	0
配列35	AGCC-TTA	0
配列36	AGCCGTTA	1
配列37	AGCCGTTA	1
配列38	AGCCGTTA	1

[0041] この場合、ステップS207における判定において、短縮後の編集距離の

平均値が0.25以上0.75以下であるので、SNP等が存在する可能性があると判定され、処理はステップS208に進む。表5に示した短縮後の配列によれば、配列36乃至38が、SNP等が存在する可能性のある配列である。その配列数である3は、全配列数である8の37.5%である。したがって、ステップS208においてもSNP等が存在すると判定され、処理はステップS209に進む。

[0042] ステップS209において、表4に示された第1の対照配列のうち、配列36乃至38に相当する配列28乃至30にはSNP等が存在すると特定される。この例のように、本実施形態によれば、ホモポリマー問題を有する第1の対照配列においても、SNP等が存在する配列が正確に特定され得る。

[0043] 一群の第1の対照配列の別の例を表6に示す。この場合、配列39乃至46のステップS202で算出される編集距離は表6に示すとおりであり、ステップS203で算出される編集距離の平均値は0.625である。

[表6]

表6

リファレンス配列	AGCCTCTTA	編集距離
配列39	AGCCTCTTA	0
配列40	AGCCTCTTA	0
配列41	AGCCTCTTA	0
配列42	AGCCTCTTA	0
配列43	AGCCTTTTA	1
配列44	AGCCTTTTA	1
配列45	AGCCT-TTA	1
配列46	AGCC--TTA	2

[0044] ステップS204で短縮された結果、配列39乃至46は、それぞれは表7に示す第2の対照配列としての配列47乃至54のようになる。この場合、ステップS205で算出される短縮後の編集距離は表7に示すとおりであり、ステップS206で算出される短縮後の編集距離の平均値は1である。



[表7]

表7

リファレンス配列	AGCCTCTTA	編集距離
配列47	AGCCTCTTA	0
配列48	AGCCTCTTA	0
配列49	AGCCTCTTA	0
配列50	AGCCTCTTA	0
配列51	AGCC--TTA	2
配列52	AGCC--TTA	2
配列53	AGCC--TTA	2
配列54	AGCC--TTA	2

[0045] この例では、ステップS207における判定において、短縮後の編集距離の平均値が0.75以上であり、短縮前の編集距離の平均値が0.25以上0.75以下であるので、SNP等が含まれている可能性があるとして、処理はステップS208に進む。表7に示した短縮後の第2の対照配列によれば、配列51乃至54が、SNP等が存在する可能性のある配列である。その配列数である4は、全配列数である8の50%である。したがって、SNP等が存在すると判定され、処理はステップS209に進む。ステップS209において、表6に示された第1の対照配列のうち、配列51乃至54に相当する配列43乃至46にはSNP等が存在すると特定される。

[0046] この例のように、同一の塩基が連続する配列において、間に1塩基の置換が存在する場合、短縮後の編集距離が非常に大きくなる。したがって、短縮後の編集距離が0.25以上0.75以下であるか否かの判定のみで分離すべきか否かの判定が行われると誤った結果になる。そこで本実施形態では、短縮後の編集距離が0.75以上であり、短縮前の編集距離が0.25以上0.75以下である場合も、分離すべきと判定されるようになっている。このようにして正確なSNP等の特定がなされる。

[0047] 一群の第1の対照配列の別の例を表8に示す。この場合、配列55乃至6

2のステップS202で算出される編集距離は表8に示すとおりであり、ステップS203で算出される編集距離の平均値は0.125である。

[表8]

表8

リファレンス配列	AGCTTCTTA	編集距離
配列55	AGCTTCTTA	0
配列56	AGCTTCTTA	0
配列57	AGCTTCTTA	0
配列58	AGCTTCTTA	0
配列59	AGCTTCTTA	0
配列60	AGCTTCTTA	0
配列61	AGCTTCTTA	0
配列62	AGCTTTTTA	1

[0048] 配列55乃至62がステップS204で短縮された結果である第2の対照配列はそれぞれ表9に示す配列63乃至70のようになる。この場合、ステップS205で算出される短縮後の編集距離は表9に示すとおりであり、ステップS206で算出される短縮後の編集距離の平均値は0.375である。

[表9]

表9

リファレンス配列	AGCTTCTTA	編集距離
配列63	AGCTTCTTA	0
配列64	AGCTTCTTA	0
配列65	AGCTTCTTA	0
配列66	AGCTTCTTA	0
配列67	AGCTTCTTA	0
配列68	AGCTTCTTA	0
配列69	AGCTTCTTA	0
配列70	AGC---TTA	3

[0049] この場合、ステップS207における判定において、短縮後の編集距離の平均値が0.25以上0.75以下であるので、SNP等が存在する可能性があると判定され、処理はステップS208に進む。表9に示した短縮後の配列によれば、第1の対照配列は、配列70が、SNP等が存在する可能性のある配列である。その配列数である4は、全配列数である8の12.5%である。したがって、リードエラーが存在すると判定され、処理はステップS210に進む。ステップS210において、リードエラーが存在すると結論付けられ、処理は終了する。

[0050] この例のように、短縮後の編集距離の平均値が0.25以上0.75以下でありステップS207でSNPが存在する可能性があると判定されても、配列数に偏りが大きい場合、ステップS208の判定でリードエラーであると判定され、差異はリードエラーによるものであると結論付けられる。

[0051] このように本実施形態によれば、ホモポリマー問題を含む配列データの解析において、ホモポリマーが短縮され、短縮された配列の編集距離に基づいて、第1のリファレンス配列と第1の対照配列との差異が、SNP等によるものであるのか、リードエラーによるものであるのかが正確に判定される。すなわち、本実施形態によれば、SNP等が正確に同定され得る。

[0052] なお、図3を参照して説明した処理は一例であり、各処理の順序は変更され得るし、一部の変更や省略もされ得る。また、本実施形態では、2倍体のヘテロ接合型の場合であって、第1の対照配列のおよそ半分がリファレンス配列と同一であり、残りにSNP等が存在する場合を想定している。このため、ステップS207の判定において、判定の基準となる範囲が表3に示すように0.5を中心とした範囲に設定されている。しかしながら、これに限らない。ホモ接合型である場合、ほぼ全ての第1の対照配列がリファレンス配列と1塩基異なる配列となることが想定される。したがって、SNP等がある場合、編集距離の平均値は1近くになると想定される。したがって、この場合には表3に相当するステップS207における判定基準は、例えば表10のようになる。また、この場合、ステップS208における判定基準は、例えば75%以上の配列がSNP等を含むと考えられるとき、第1の対照配列にはSNP等が含まれると判定され、75%未満のとき、配列の差異はリードエラーによるものと判定されるように設定される。

[表10]

表10

		短縮前		
		編集距離<0.75	0.75≤編集距離≤1.25	1.25<編集距離
短縮後	編集距離<0.75	リードエラー	リードエラー	リードエラー
	0.75≤編集距離≤1.25	SNP	SNP	SNP
	1.25<編集距離	リードエラー	SNP	リードエラー

[0053] また、ヘテロ接合型であって、第1の対照配列のおよそ半分がリファレンス配列と異なるある第1のSNPを含み、残りが別の第2のSNPを含む場合も、表3に相当するステップS207における判定基準は、例えば表10のようになる。また、ステップS208の判定では、第1のSNPを含む第1の対照配列の数と、第2のSNPを含む第2の対照配列の数とが、それぞれ全配列数の33%以上67%以下である場合、SNP等が含まれ、それ以外のとき、配列の差異はリードエラーによるものであると判定されるように

設定される。

[0054] また、SNP等の同定以外の用途に用いられるのであれば、変異によるグループ分けが3つ以上となるように構成されてもよい。この場合、ステップS207における判定基準や、ステップS208における判定基準は異なるものになる。例えば、4つのグループに分けられる場合は、ステップS207における判定基準の範囲は、0.25を中心とした範囲に設定される。

[0055] また、本実施形態では、リファレンス配列に対する編集距離を求めているが、一群の配列のコンセンサス配列に対する編集距離を求めるようにしてもよい。ここでコンセンサス配列とは、一群の第1の対照配列の全てに対する同一性が最も高くなるように決定された配列、すなわち第1の対照配列の多数に共通する配列のことをいう。

[0056] また、判定装置1の演算部30は、図5に示すようにクラスタ決定部35をさらに備えてもよい。このクラスタ決定部35は、エラー判定部34が特定したSNP等に基づいて、第1の対照配列を2つのクラスタに分離する。例えば、クラスタ決定部35は、表4に示す一群の配列について、配列23乃至27を含む第1のクラスタと、配列28乃至30を含む第2のクラスタとにグループ分けする。この処理は、例えばステップS209の後に行われる。

[0057] このように、本実施形態によれば、配列の差異が変異によるものなのかシーケンシングにおけるリードエラーによるものなのかが判定される。この判定を用いれば、ゲノムやcDNA配列のバリエーションの真偽が判定され得る。また、配列のアセンブルにより作成されたクラスタの変異部分の真偽が判定され得る。また、判定結果によりクラスタが分割され得るし、クラスタを作成する際の前処理方法となり得る。

[0058] [第3の実施形態]

本発明の第3の実施形態について説明する。ここでは、第2の実施形態との相違点について説明し、同一の部分については、同一の符号を付してその説明を省略する。本実施形態では、第2の実施形態に係るSNP等の同定手

法とそれを用いたクラスタの作成手法とをHiCEP (High Coverage Expression Profiling) 法に適用する。したがって、本実施形態に係る判定装置1は、図5に示す構成を有する。

[0059] HiCEPの概要を図6を参照して説明する。HiCEPでは、サンプル内のmRNA 201が逆転写され種々のcDNA 202が合成される。このcDNAの一部である断片DNA 203が切り出され、その各々の端に既知の配列であるアダプタ配列204が付加される。このアダプタ配列204が付加された種々の断片DNA 203を含むサンプルは、256分割される。その後、アダプタ配列を用いたPCR法が利用されることで、断片DNA 203とアダプタ配列204との間に2塩基のセレクション塩基205が挿入されたDNAが増幅され、そこに蛍光色素206が付加される。256分割された分注サンプルは、分注サンプルごとにセレクション配列が異なる。ここで、セレクション配列は、両端に2塩基ずつ、すなわち計4塩基が付加されており、その組み合わせは256通りである。HiCEPでは、256種類の分注サンプルのそれぞれが電気泳動されて、その分注サンプルに含まれる断片DNA 203が塩基長ごとに分離され、その量が蛍光色素206を用いて定量される。

[0060] 上記のようなHiCEPによれば、サンプル中の3万乃至6万種類のmRNAの断片について、それぞれの分子数に係る情報が取得され得る。HiCEPによれば、遺伝子資源を必要とせずに網羅的に高感度に発現プロファイルが再現性高く取得され得る。

[0061] HiCEPによって、例えば発現が変動する断片が検出されたらその遺伝子を同定することが求められる。網羅的なDNAシーケンシングには、高スループットな次世代シーケンサが用いられる。次世代シーケンサのうち特にピロリン酸配列決定技法を用いたシーケンシングで発生するホモポリマー問題は、アダプタ配列に隣接したセレクション塩基部分において発生すると、HiCEPを用いた解析に大きな悪影響を与える恐れがある。そこで本実施形態では、セレクション塩基部分において、第2の実施形態に係るSNP等

の同定手法を用いる。

[0062] 本実施形態において実行される処理のフローチャートを図7に示す。演算部30のデータ取得部31は、ステップS301において、配列データを取得する。ここで取得される配列データは、例えば表11に示すようなデータである。この例では、第1の対照配列として配列71乃至配列78が含まれている。また、リファレンス配列としては、配列71乃至78のコンセンサス配列が用いられる。このコンセンサス配列を第1のコンセンサス配列と称することにする。表12において、アダプタ配列の最もリード配列側の2塩基（「TT」）が表中の「アダプタ配列」の列に記載され、アダプタ配列よりもリード配列側の配列が、表中の「リード配列、及びセレクション塩基」の列に記載されている。

[表11]

表 1 1

	リード配列、及び セレクション塩基	アダプ タ配列	編集距離
コンセンサス配列	CAGCCT-	TT	
配列71	CAGCCT-	TT	0
配列72	CAGCCT-	TT	0
配列73	CAGCCTT	TT	1
配列74	CAGCCTT	TT	1
配列75	CAGCCCT	TT	1
配列76	CAGCCG-	TT	1
配列77	CAGCCGT	TT	1
配列78	CAGCCGT	TT	1

[0063] 演算部30の配列短縮部32は、ステップS302において、第1のコンセンサス配列及び各第1の対照配列について、同一塩基が3塩基以上連続する配列を同一塩基が2塩基連続する配列に短縮する。表11に示した配列の短縮後の配列を表12に示す。配列71乃至78は、それぞれ配列79乃至

86のように短縮される。

[表12]

表12

	リード配列、及び セレクション塩基	アダプ タ配列	編集距離
コンセンサス配列	CAGCC	TT	
配列79	CAGCC	TT	0
配列80	CAGCC	TT	0
配列81	CAGCC	TT	0
配列82	CAGCC	TT	0
配列83	CAGCC	TT	0
配列84	CAGCCG	TT	1
配列85	CAGCCG	TT	1
配列86	CAGCCG	TT	1

[0064] 演算部30の編集距離算出部33は、ステップS303において、短縮前の第1の対照配列について、第1のコンセンサス配列との編集距離を算出する。演算部30の編集距離算出部33は、ステップS304において、短縮前の第1の対照配列の編集距離の平均値を算出する。演算部30の編集距離算出部33は、ステップS305において、短縮後の第2の対照配列について、第1のコンセンサス配列が短縮されて生成された第2のコンセンサス配列に対する編集距離を算出する。演算部30の編集距離算出部33は、ステップS306において、短縮後の第2の対照配列の編集距離の平均値を算出する。

[0065] ここで、編集距離は、第2のコンセンサス配列におけるリード配列及びセレクション塩基の最もアダプタ配列に近い5塩基を対象とし、この5塩基に対応する第2の対照配列について算出される。すなわち、表11及び表12の「リード配列、及びセレクション塩基」に示した塩基配列に関して編集距離が算出される。コンセンサス配列と対応する配列に注目するので、例えば



配列 84 乃至 86 のように、対象となる配列は 5 塩基とは限らない。編集距離は、表 11 及び表 12 に示すとおりとなる。

[0066] 演算部 30 のエラー判定部 34 は、ステップ S307 において、編集距離の平均値が所定の範囲内であるか否かに応じて、第 1 の対照配列に SNP 等が含まれている可能性があるか否かを判定する。この判定には、例えば上述の表 3 に示す判定基準が用いられる。すなわち、短縮後の編集距離が 0.25 以上 0.75 以下である場合、第 1 の対照配列に SNP 等が含まれている可能性があるとして判定される。また、短縮後の編集距離が 0.75 より大きいときでも短縮前の編集距離が 0.25 以上 0.75 以下である場合、第 1 の対照配列に SNP 等が含まれている可能性があるとして判定される。それ以外の場合、差異はリードエラーによるものであるとして判定される。第 1 の対照配列に SNP 等が含まれている可能性があるとして判定されたとき、処理はステップ S308 に進む。差異はリードエラーによるものであるとして判定されたとき、処理はステップ S311 に進む。

[0067] 演算部 30 のクラスタ決定部 35 は、ステップ S308 において、リード配列の最もアダプタ配列側 2 塩基に基づいて第 1 の塩基配列を分離した場合にどのように分離されるかを仮定する。本実施形態では、リード配列の最もアダプタ配列側 2 塩基のみに注目する。これは、この部分にセレクション塩基があり、このセレクション塩基を重要視しているためである。例えば表 12 の例では、配列 79 乃至 83 の「CC」と配列 84 乃至 86 の「CG」との 2 種類があるので、クラスタ決定部 35 は、配列 79 乃至 86 は配列 79 乃至 83 が含まれるクラスタと配列 84 乃至 86 が含まれるクラスタとの 2 つに分離されると仮定する。

[0068] 演算部 30 のエラー判定部 34 は、ステップ S309 において、第 1 の対照配列が 2 つに分離された場合に各クラスタに含まれる配列の数が所定の範囲内であるか否かに応じて、第 1 の対照配列に SNP 等が含まれている可能性があるか否かを判定する。例えば、2 つのクラスタに分離されたときに、両クラスタ内の配列数がそれぞれ全配列数の 33% 以上 67% 以下であると

き、第1の対照配列にSNP等が含まれていると判定される。この判定は、第1の対照配列を同一配列ごとにグループ分けした場合、最も多くの配列が含まれるグループの配列数が全配列数の66%以下であり、かつ、2番目に多くの配列が含まれるグループの配列数が全配列数の33%以上である場合に、SNP等が含まれていると判定されるように変更されてもよい。それ以外の場合、配列の差異はリードエラーによるものであると判定される。第1の対照配列にSNP等が含まれていると判定されたとき、処理はステップS310に進む。配列の差異はリードエラーによるものであると判定されたとき、処理はステップS311に進む。

[0069] 表12に示す例では、配列79乃至83が含まれるクラスタ内の配列数は5であり、全配列数である8の62.5%である。一方、配列84乃至86が含まれるクラスタ内の配列数は3であり、全配列数である8の37.5%である。したがって、第1の対照配列にSNP等が含まれていると判定され処理は、ステップS310に進む。

[0070] 演算部30のクラスタ作成部35は、ステップS311において、第1の対照配列を、ステップS308で仮定したクラスタに分離し、2つのクラスタを作成する。演算部30は、作成したクラスタに係る情報を出力部43に出力し、記録等させる。その後処理は終了する。

[0071] 演算部30のクラスタ決定部35は、ステップS311において、第1のコンセンサスと第1の対照配列との差異は、リードエラーによるものであると結論する。演算部30は、第1のコンセンサスと第1の対照配列との差異はリードエラーによるものである旨を出力部43に出力し、記録等させる。その後処理は終了する。

[0072] 本実施形態によれば、HiCEPを用いた解析においても問題となるホモポリマー問題が解消され、HiCEPにおいて重要であるセレクション塩基部分のクラスタリングが正確に行われる。なお、本実施形態は、HiCEPに限らず、アダプタ配列を用いる他の解析においても同様に用いられ得る。

[0073] [第4の実施形態]

本発明の第4の実施形態について説明する。ここでは、第1の実施形態との相違点について説明し、同一の部分については、同一の符号を付してその説明を省略する。本実施形態では、第1の実施形態におけるステップS103の処理と同様に、同一塩基が3塩基以上連続する配列をその塩基が2塩基連続した配列に短縮変換された第2の塩基配列を作成し、この第2の塩基配列に基づいて、既知のクラスタリング処理を行う。

[0074] 本実施形態によれば、同一塩基が連続する配列が短縮され、塩基長が短くなるのでクラスタリングの効率が向上する。また、同一塩基が連続する配列が短縮され、ホモポリマー問題が解消されるので、クラスタリングの精度が向上する。

[0075] なお、本実施形態は、HiCEPで得られたサンプルに係るシーケンシング結果に適用してもよい。HiCEPに係るデータの場合、アダプタ配列によって断片長が揃っていないので、よいクラスタリング結果が得られ、本実施形態は特に効果を奏する。

[0076] なお、本発明は上記実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記実施形態に開示されている複数の構成要素の適宜な組み合わせにより、種々の発明を形成できる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除しても、発明が解決しようとする課題の欄で述べられた課題が解決でき、かつ、発明の効果が得られる場合には、この構成要素が削除された構成も発明として抽出され得る。さらに、異なる実施形態にわたる構成要素を適宜組み合わせてもよい。

## 請求の範囲

- [請求項1] 第1のリファレンス配列と、前記第1のリファレンス配列と相同性を有する複数の第1の対照配列との差異が、前記第1の対照配列の変異によるものかシーケンシングのリードエラーによるものかを判定する方法であって、
- 各々の前記第1の対照配列の同一塩基が所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して各々の第2の対照配列を作成することと、
- 前記第1のリファレンス配列の同一塩基が前記所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して第2のリファレンス配列を作成することと、
- 前記第2のリファレンス配列に対する各々の前記第2の対照配列の編集距離である複数の短縮後編集距離を算出することと、
- 複数の前記短縮後編集距離の平均値である短縮後平均値を算出することと、
- 前記短縮後平均値に基づいて、前記差異が前記変異によるものであるか前記リードエラーによるものであるかを判定することと、
- を具備する方法。
- [請求項2] 前記第1のリファレンス配列に対する各々の前記第1の対照配列の編集距離である複数の短縮前編集距離を算出することと、
- 複数の前記短縮前編集距離の平均値である短縮前平均値を算出することと、
- をさらに具備し、
- 前記判定することは、前記短縮前平均値と前記短縮後平均値との関係に基づいて、前記差異が前記変異によるものであるか前記リードエラーによるものであるかを判定することとである、
- 請求項1に記載の方法。
- [請求項3] 前記判定することは、前記短縮後平均値が所定の平均値範囲である

とき、又は、前記短縮後平均値が前記平均値範囲よりも大きくて且つ前記短縮前平均値が前記平均値範囲であるとき、前記差異が前記変異によるものであると判定することである、請求項2に記載の方法。

[請求項4] 前記判定することは、

第1の条件である、前記短縮後平均値が所定の平均値範囲であること、又は、前記短縮後平均値が前記平均値範囲よりも大きくて且つ前記短縮前平均値が前記平均値範囲であることと、

第2の条件である、全ての前記第1の対照配列の数に対する、前記第2のリファレンス配列と前記第2の対照配列とに差異がある前記第2の対照配列の数の割合が所定の配列数範囲であることと、

を共に満たすとき、前記差異が前記変異によるものであると判定することである、請求項2に記載の方法。

[請求項5] 前記所定塩基数は2である、請求項1乃至4のうち何れか1項に記載の方法。

[請求項6] 前記平均値範囲は、0.5を含む範囲である、請求項3又は4に記載の方法。

[請求項7] 前記配列数範囲は、50%を含む範囲である、請求項4に記載の方法。

[請求項8] 前記第1のリファレンス配列は、前記第1の対照配列の同一性に基づいて得られるコンセンサス配列である、請求項1乃至4のうち何れか1項に記載の方法。

[請求項9] 前記差異が前記変異によるものであると判定されたとき、複数の前記第1の対照配列を前記差異に応じてクラスタリングすることをさらにコンピュータに実行させる、請求項1乃至4のうち何れか1項に記載の方法。

[請求項10] 前記第1の対照配列は、HiCEPで用いられる2つのアダプタ配列に挟まれる塩基配列であり、

前記短縮後編集距離を算出することは、前記アダプタの端を基準と

して所定の数の塩基について前記第2のリファレンス配列に対する各々の前記第2の対照配列の編集距離である複数の短縮後編集距離を算出することであり、

前記短縮前編集距離を算出することは、前記短縮後編集距離を算出した塩基配列に対応する前記第1のリファレンス配列に対する各々の前記第1の対照配列の編集距離である複数の短縮前編集距離を算出することである、

請求項2乃至4のうち何れか1項に記載の方法。

[請求項11]

第1のリファレンス配列と、前記第1のリファレンス配列と同一性を有する第1の対照配列との差異が、前記第1の対照配列の変異によるものかシーケンシングのリードエラーによるものかを判定する方法であって、

前記第1の対照配列の同一塩基が所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して第2の対照配列を作成することと、

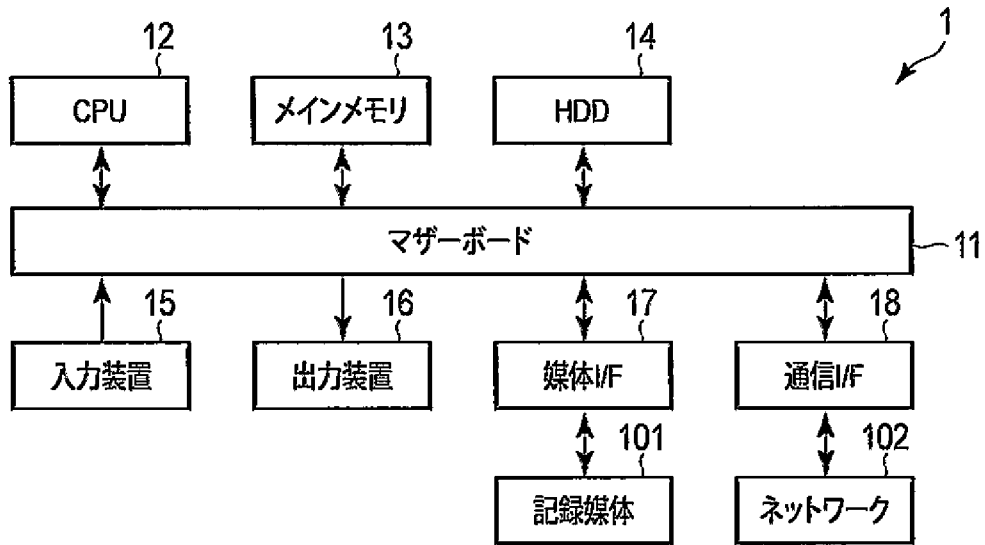
前記第1のリファレンス配列の同一塩基が前記所定塩基数以上連続する配列をこの同一塩基が前記所定塩基数連続する配列に置換して第2のリファレンス配列を作成することと、

前記第2のリファレンス配列に対する前記第2の対照配列の編集距離を算出することと、

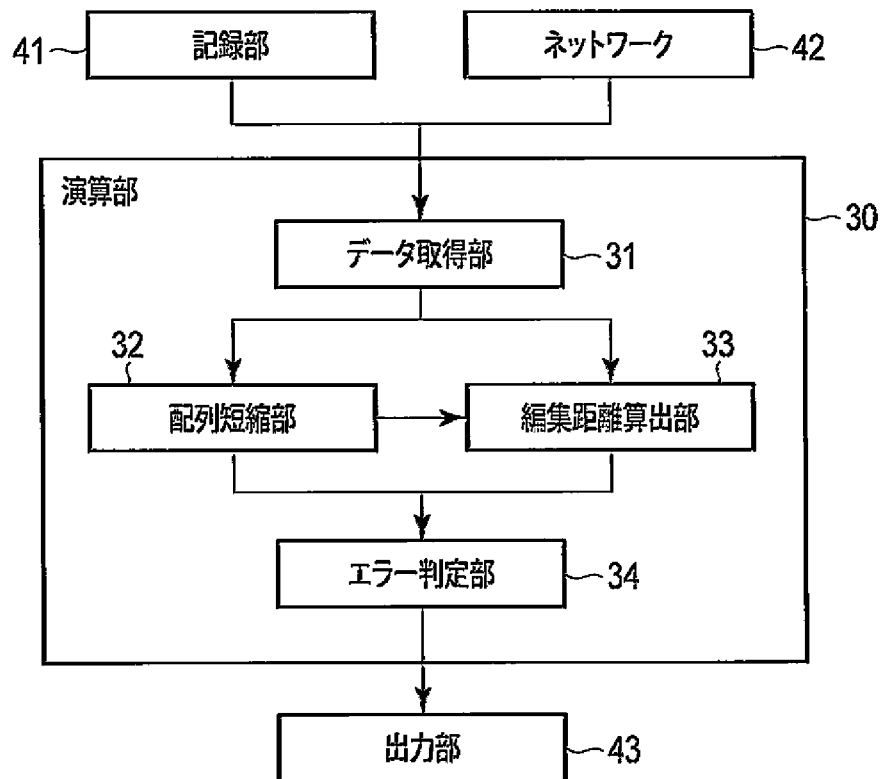
前記編集距離に基づいて、前記差異が前記変異によるものであるか前記リードエラーによるものであるかを判定することと、

を具備する方法。

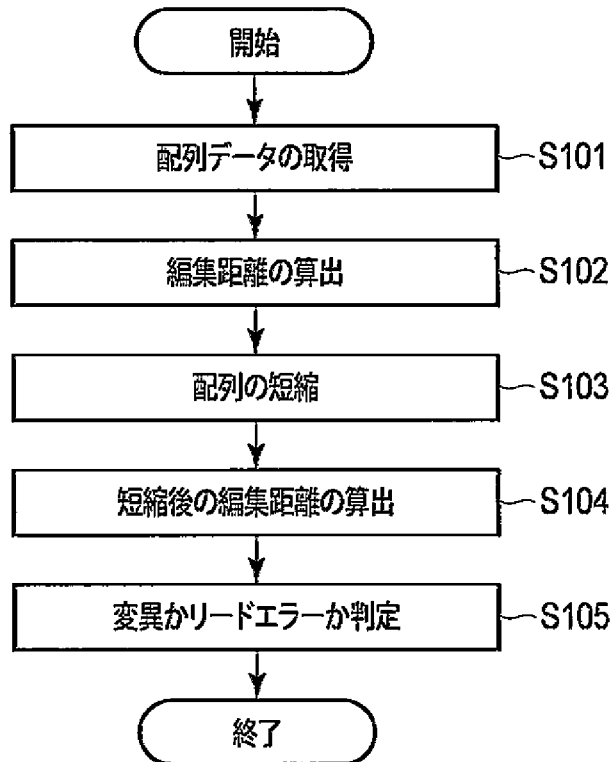
[図1]



[図2]

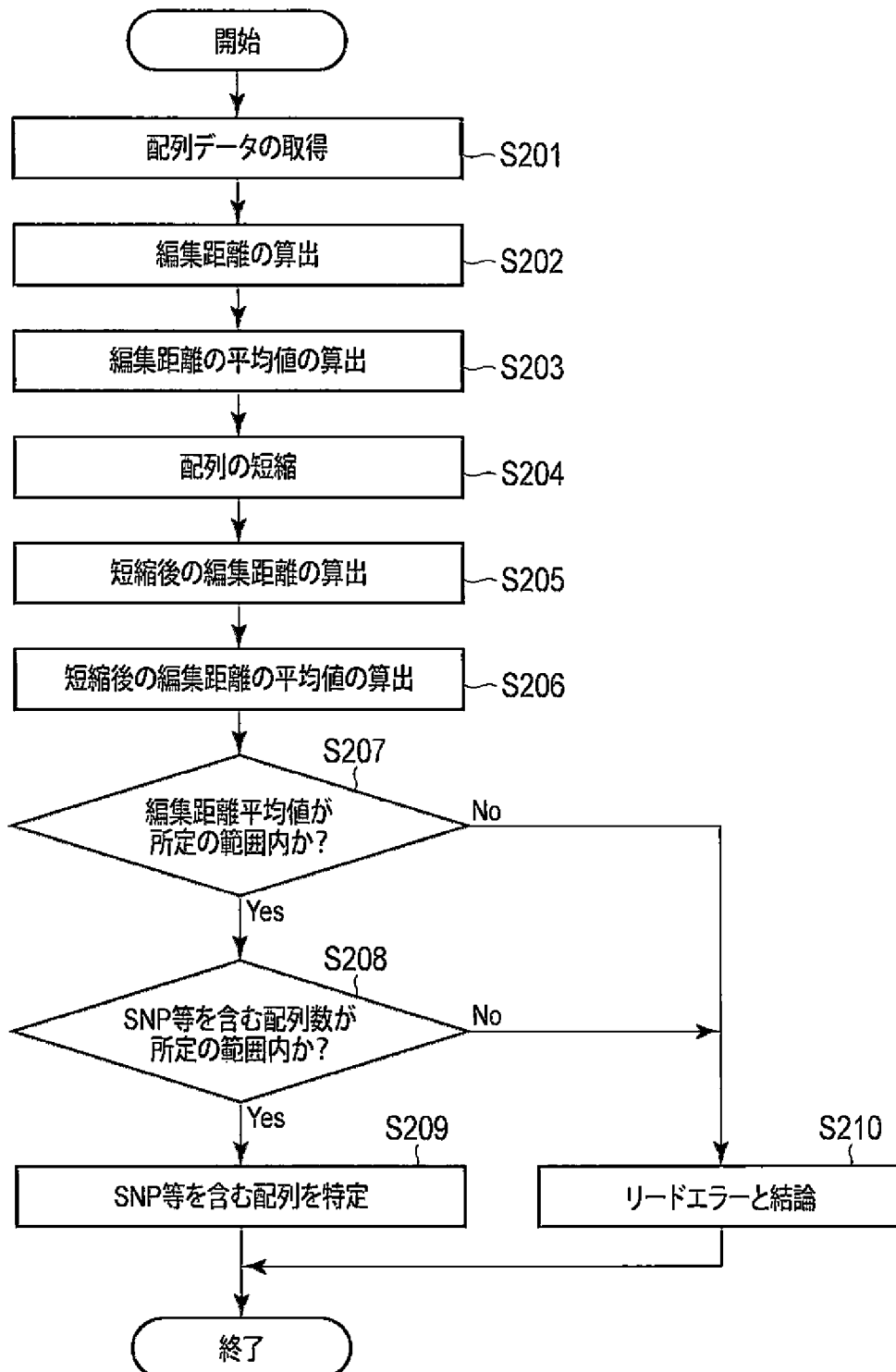


[図3]

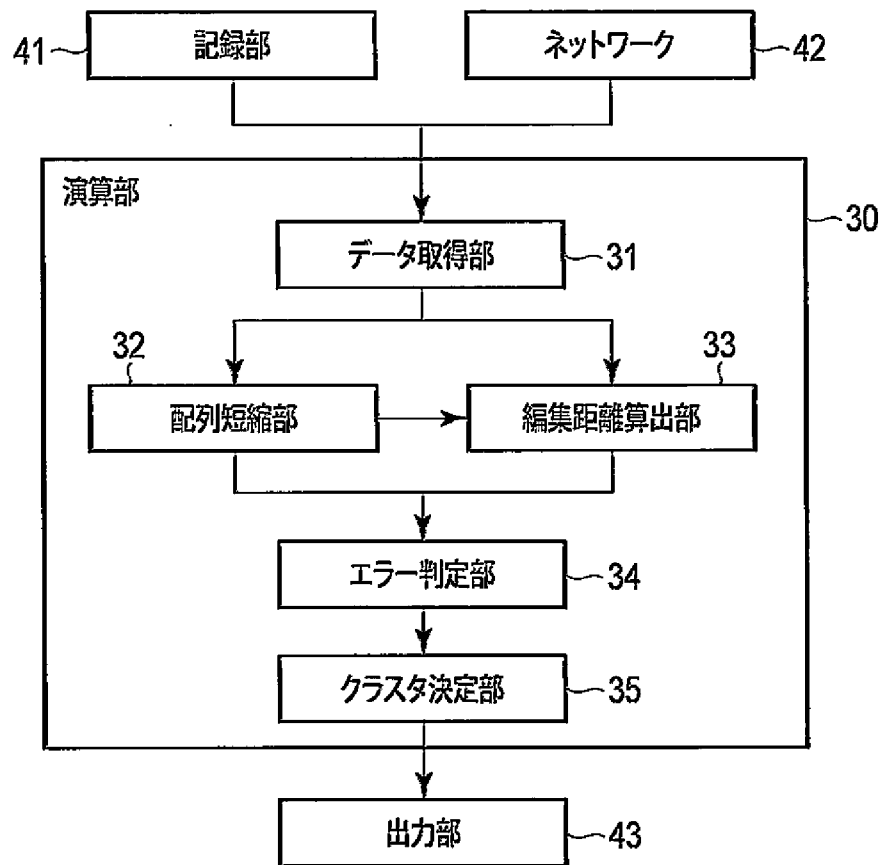




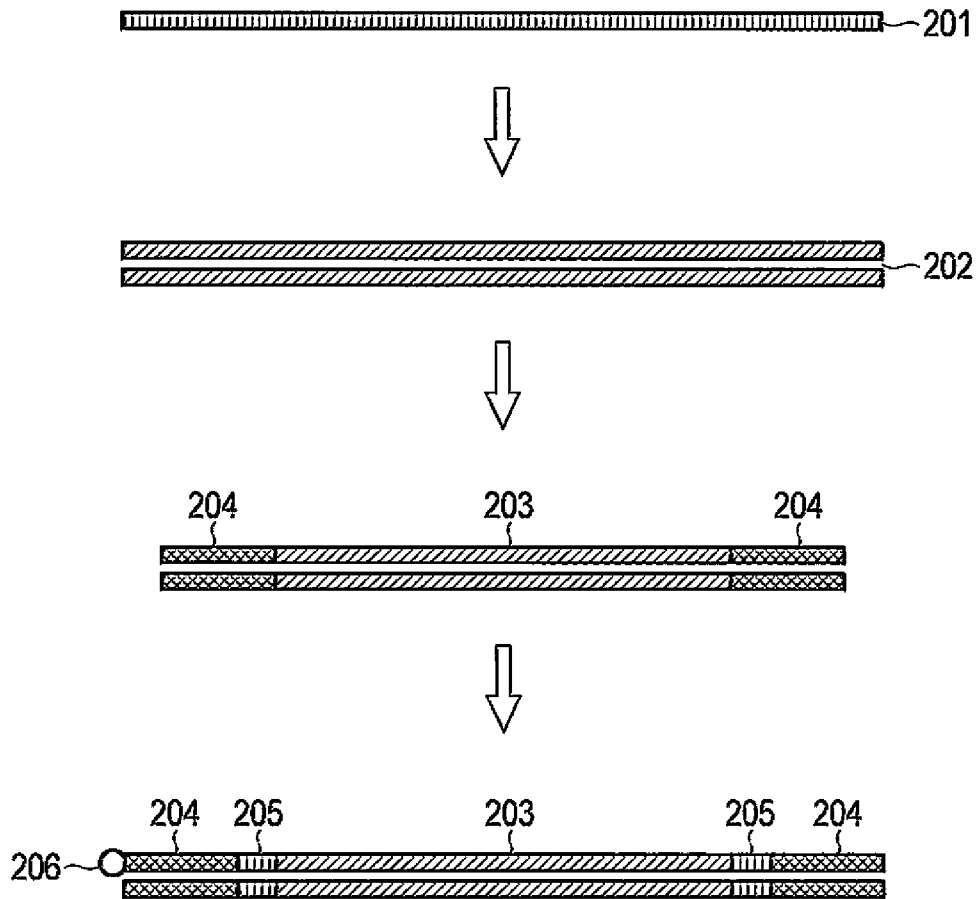
[図4]



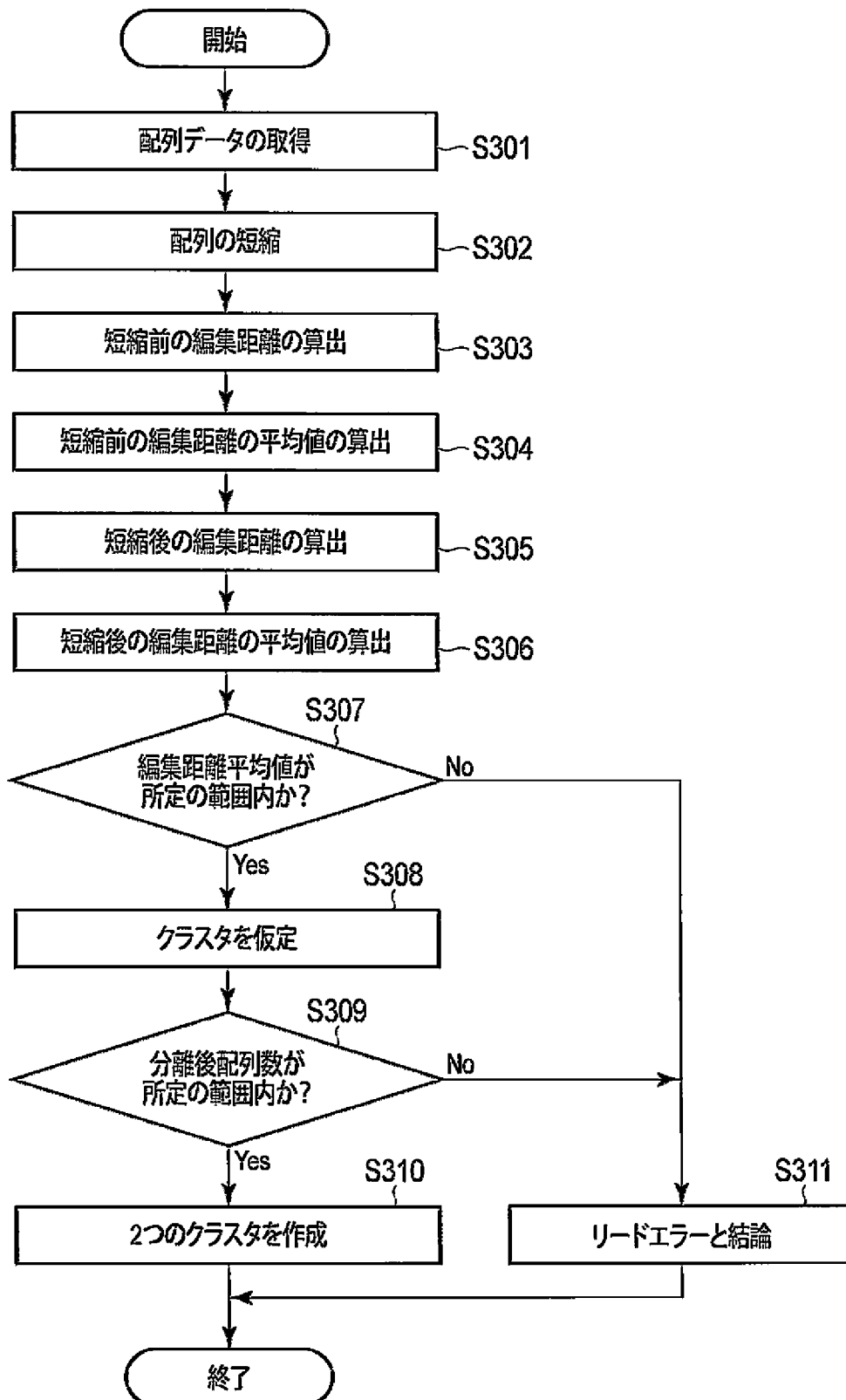
[図5]



[図6]



[図7]



**INTERNATIONAL SEARCH REPORT**

International application No.

PCT/JP2013/062426

A. CLASSIFICATION OF SUBJECT MATTER  
G06F19/22(2011.01) i, C12Q1/68(2006.01) i

According to International Patent Classification (IPC) or to both national classification and IPC

**B. FIELDS SEARCHED**

Minimum documentation searched (classification system followed by classification symbols)  
G06F19/22, C12Q1/68

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched  
 Jitsuyo Shinan Koho 1922-1996 Jitsuyo Shinan Toroku Koho 1996-2013  
 Kokai Jitsuyo Shinan Koho 1971-2013 Toroku Jitsuyo Shinan Koho 1994-2013

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	WO 2008/108297 A1 (Inter-University Research Institute Corporation Research Organization of Information and Systems), 12 September 2008 (12.09.2008), entire text; all drawings & US 2010/0205204 A1 & EP 2133807 A1	1-11
A	HiCEP Peak Database - NGSeq (PeakDB NGS), [online], 28 March 2012 (28.03.2012), [retrieval date 15 May 2013 (15.05.2013)], Internet, <URL:http://hicepweb.nirs.go.jp/about/ about_clustering_db.html>	1-11

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents:	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
"A" document defining the general state of the art which is not considered to be of particular relevance	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"E" earlier application or patent but published on or after the international filing date	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&" document member of the same patent family
"O" document referring to an oral disclosure, use, exhibition or other means	
"P" document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search  
15 May, 2013 (15.05.13)

Date of mailing of the international search report  
28 May, 2013 (28.05.13)

Name and mailing address of the ISA/  
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int.Cl. G06F19/22(2011.01)i, C12Q1/68(2006.01)i

B. 調査を行った分野

調査を行った最小限資料 (国際特許分類 (IPC))

Int.Cl. G06F19/22, C12Q1/68

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報	1922-1996年
日本国公開実用新案公報	1971-2013年
日本国実用新案登録公報	1996-2013年
日本国登録実用新案公報	1994-2013年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	WO 2008/108297 A1 (大学共同利用機関法人情報・システム研究機構) 2008.09.12, 全文、全図 & US 2010/0205204 A1 & EP 2133807 A1	1-11
A	HiCEP Peak Database - NGSeq (PeakDB NGS), [online], 2012.03.28, [検索日 2013年5月15日], インターネット <URL:http://hicepweb.nirs.go.jp/about/about_clustering_db.html>	1-11

☐ C欄の続きにも文献が列挙されている。

☐ パテントファミリーに関する別紙を参照。

\* 引用文献のカテゴリー

「A」特に関連のある文献ではなく、一般的な技術水準を示すもの  
 「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの  
 「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)  
 「O」口頭による開示、使用、展示等に言及する文献  
 「P」国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献

「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの  
 「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの  
 「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの  
 「&」同一パテントファミリー文献

国際調査を完了した日

15.05.2013

国際調査報告の発送日

28.05.2013

国際調査機関の名称及びあて先

日本国特許庁 (ISA/J P)  
 郵便番号100-8915  
 東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

岡北 有平

電話番号 03-3581-1101 内線 3562

5 L

4677