

(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
2014年8月21日(21.08.2014)



(10) 国際公開番号
WO 2014/126213 A1

- (51) 国際特許分類:
G06F 19/22 (2011.01) C12N 15/115 (2010.01)
- (21) 国際出願番号: PCT/JP2014/053516
- (22) 国際出願日: 2014年2月14日(14.02.2014)
- (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 2013-027851 2013年2月15日(15.02.2013) JP
- (71) 出願人: NECソリューションイノベータ株式会社(NEC SOLUTION INNOVATORS, LTD.) [JP/JP]; 〒1368627 東京都江東区新木場一丁目18番7号 Tokyo (JP).
- (72) 発明者: 秋富 穰(AKITOMI Jou); 〒1368627 東京都江東区新木場一丁目18番7号 NECソリューションイノベータ株式会社内 Tokyo (JP). 堀井 克紀(HORII Katsunori); 〒1368627 東京都江東区新木場一丁目18番7号 NECソリューションイノベータ株式会社内 Tokyo (JP).
- (74) 代理人: 辻丸 光一郎, 外(TSUJIMARU Koichiro et al.); 〒6008813 京都府京都市下京区中堂寺南町134 京都リサーチパーク1号館301号室 Kyoto (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, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, 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, SA, 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.
- (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 (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, KM, ML, MR, NE, SN, TD, TG).

(54) Title: SELECTION DEVICE FOR CANDIDATE SEQUENCE INFORMATION FOR SIMILARITY DETERMINATION, SELECTION METHOD, AND USE FOR SUCH DEVICE AND METHOD

(54) 発明の名称: 類似判断の候補配列情報の選択装置、選択方法、およびそれらの用途

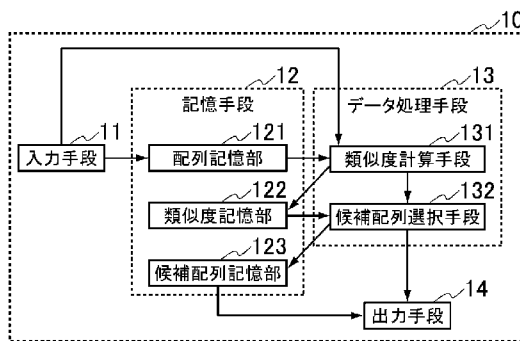


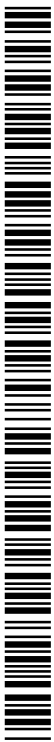
FIG. 1:
 11 Input means
 12 Storage means
 13 Data processing means
 14 Output means
 121 Sequence storage section
 122 Degree of similarity storage section
 123 Candidate sequence storage means
 131 Degree of similarity calculation means
 132 Candidate sequence selection means

(57) Abstract: Provided is a device for easily determining the similarities between sequence information items. This candidate selection device (10) is provided with an input means (11), a sequence storage section (121), a degree of similarity storage section (122), a candidate sequence storage section (123), a degree of similarity calculation means (131), a candidate sequence selection means (132), and an output means (14). The input means (11) inputs information on sequence groups and virtual sequence groups. The degree of similarity calculation means (131) selects from the sequence groups a comparison source and a comparison destination, and calculates the difference in frequency of the virtual sequences with respect to a comparison source sequence and a comparison destination sequence, as a degree of similarity of the comparison destination sequence with respect to the comparison source sequence. In a case where the degree of similarity of the comparison destination sequence with respect to the comparison source sequence fulfills the allowable conditions of the degree of similarity set to the virtual sequence group, the candidate sequence selection means (132) selects the comparison source sequence and the comparison destination sequence as a candidate sequence group for determining the similarities between sequences. For the candidate sequence group, by determining the similarities between sequences, one sequence and a sequence similar thereto are selected as a similar sequence information group.

sequence similar thereto are selected as a similar sequence information group.

(57) 要約:

[続葉有]



WO 2014/126213 A1

添付公開書類:

— 国際調査報告 (条約第 21 条(3))

配列情報間の類似を、容易に判断するための装置を提供する。本発明の候補選択装置 10 は、入力手段 11、配列記憶部 121、類似度記憶部 122、候補配列記憶部 123、類似度計算手段 131、候補配列選択手段 132、出力手段 14 を備える。入力手段 11 は、配列群および仮想配列群の情報を入力し、類似度計算手段 131 は、前記配列群から比較元と比較先とを選択し、比較元配列と前記比較先配列との前記各仮想配列の頻度の相違を、前記比較元配列に対する前記比較先配列の類似度として計算する。候補配列選択手段 132 は、前記比較元配列に対する前記比較先配列の類似度が、前記仮想配列群に設定した類似度の許容条件を満たす場合、前記比較元配列および前記比較先配列を、配列間の類似を判断する候補配列群として選択する。前記候補配列群について、配列間の類似を判断することにより、ある配列とこれに類似する配列とを類似配列情報群として選択する。

明 細 書

発明の名称：

類似判断の候補配列情報の選択装置、選択方法、およびそれらの用途

技術分野

[0001] 本発明は、配列情報群における配列情報間の類似の判断に関する発明であり、具体的には、配列情報から類似判断の候補配列情報を選択する候補選択方法、候補配列情報から類似配列情報群を選択する類似選択方法、目的の類似配列情報群の濃縮を判定する判定方法、およびこれらの方法を実行する各装置、プログラムならびに記録媒体に関する。

背景技術

[0002] 近年、抗体に代わるターゲットへの結合分子として、いわゆるアプタマーと呼ばれる核酸分子の開発が進められている。前記アプタマーは、一般に、S E L E X (Systematic Evolution of Ligands by EXponential enrichment) 法により調製されている（特許文献1、非特許文献1）。S E L E X法は、核酸ライブラリーと前記ターゲットとの接触、および、前記ターゲットに結合した核酸の増幅を、1セットの選択処理とし、複数ラウンドを繰り返す。これによって、初期のライブラリーから、ラウンド毎のライブラリーにおいて前記ターゲットに結合する核酸配列が濃縮される。そして、例えば、ライブラリー内で濃縮度合いが相対的に高い複数の核酸配列を、アプタマー候補群として選択し、さらに、前記ターゲットとの結合力等を評価することによって、最終的に前記ターゲットに結合するアプタマーを決定することができる。

[0003] このように、アプタマー候補群は、ライブラリー内における濃縮度合いによって選択できるため、S E L E X法においては、濃縮度合いの評価が必要である。濃縮度合いの評価は、通常、以下のように行われている。まず、各ラウンドのライブラリーに含まれる核酸配列をシーケンスで解読する。そして、ライブラリー内における同じ核酸配列の出現数（以下、重複度ともいう

)をカウントする。このカウント数の増減により、各核酸配列の濃縮度合いを評価する。例えば、 n 回目のラウンド (R_n)における核酸配列 X の重複度 m_n と、次のラウンド、すなわち $n+1$ 回目のラウンド (R_{n+1})における核酸配列 X の重複度 m_{n+1} とを比較して、 $m_n < m_{n+1}$ であれば、核酸配列 X は、ラウンド ($n+1$)において、ラウンド (n)よりも濃縮されていると判断できる。また、同じラウンドのライブラリー内において、核酸配列 X の重複度 m_x と核酸配列 Y の重複度 m_y とを比較して、重複度の大きい方が、他方に比べて濃縮されていると判断できる。

先行技術文献

特許文献

[0004] 特許文献1：特許第2763958号

非特許文献

[0005] 非特許文献1：Science, (1990) 249, 505-510.

発明の概要

発明が解決しようとする課題

[0006] しかしながら、濃縮度合いによってアプタマー候補群を選択しても、異なる全ての核酸配列について、前記ターゲットとの結合力を評価することは、非常に労力を有し、現実的ではない。

[0007] 一方、ライブラリー内には、ある核酸配列（以下、元配列ともいう）に対して完全に同じ塩基配列も含まれるが、前記元配列に対して数塩基程度のミスマッチを有する類似した核酸配列（以下、類似配列ともいう）が含まれる場合がある。そして、発明者らは、前記類似配列は、例えば、前記ターゲットとの結合の強さが前記元配列と異なることがあるが、前記ターゲットに対する特性等は、前記元配列と同一であることが多いとの知見を得ている。このため、核酸配列について、完全に同一か否かという分類ではなく、許容できる範囲で類似し合っている核酸配列を、同一の配列群とすることにより、アプタマーの評価を効率化できる。しかしながら、この場合、複数の核酸配

列を一個ずつ照らし合わせて類似か否かを判断することも、労力、コストおよび時間がかかる。特に、次世代シーケンサー等を用いて大量の核酸配列の情報が得られた場合等、非常に計算コストがかかる。また、このような問題は、核酸配列に特化した問題ではなく、要素が並んだ配列情報について、共通する問題である。

[0008] そこで、本発明は、容易に、配列情報間の類似を判断するための装置、方法、プログラムおよび記録媒体を提供することを目的とする。

課題を解決するための手段

[0009] 前記目的を達成するために、本発明の候補選択装置は、下記（a）、（b）、（c）および（d）手段を備えることを特徴とする、配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する候補選択装置である。

（a）配列情報群の各配列情報について、仮想配列情報群の各仮想配列情報の頻度をカウントする工程を実行する手段

（b）前記配列情報群から、比較元となる配列情報と比較先となる配列情報とを選択する工程を実行する手段

（c）前記比較元配列情報の前記各仮想配列情報の頻度と、前記比較先配列情報の前記各仮想配列情報の頻度との相違を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程を実行する手段

（d）前記比較元配列情報に対する前記比較先配列情報の類似度が、前記仮想配列情報群に設定した類似度の許容条件を満たす場合、前記比較元配列情報および前記比較先配列情報を、配列情報間の類似を判断する候補配列情報群として選択する工程を実行する手段

[0010] 本発明の類似選択装置は、下記（A）および（B）手段を備え、前記（A）手段が、前記本発明の候補選択装置であることを特徴とする、配列情報群から、相互に類似する類似配列情報群を選択する類似選択装置である。

（A）配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する工程を実行する手段

(B) 前記候補配列情報群の各候補配列情報を相互に対比し、同一および類似する配列情報を類似配列情報群 (G3) として選択する工程を実行する手段

[0011] 本発明の判定装置は、下記 (X) および (Y) 手段を備え、前記 (X) 手段が、前記本発明の類似選択装置であることを特徴とする、目的の類似配列情報群の濃縮の判定装置である。

(X) 配列情報群から、目的配列情報とそれに類似する配列情報とを目的の類似配列情報群として選択する工程を実行する手段

(Y) 前記類似配列情報群における前記目的配列情報と前記類似する配列情報との重複度の合計から、前記類似配列情報群の濃縮を判定する工程を実行する手段

[0012] 本発明の候補選択方法は、下記 (a)、(b)、(c) および (d) 工程を含むことを特徴とする、配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する候補選択方法である。

(a) 配列情報群の各配列情報について、仮想配列情報群の各仮想配列情報の頻度をカウントする工程

(b) 前記配列情報群から、比較元となる配列情報と比較先となる配列情報とを選択する工程

(c) 前記比較元配列情報の前記各仮想配列情報の頻度と、前記比較先配列情報の前記各仮想配列情報の頻度との相違を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程

(d) 前記比較元配列情報に対する前記比較先配列情報の類似度が、前記仮想配列情報群に設定した類似度の許容条件を満たす場合、前記比較元配列情報および前記比較先配列情報を、配列情報間の類似を判断する候補配列情報群として選択する工程

[0013] 本発明の類似選択方法は、下記 (A) および (B) 工程を含み、前記 (A) 工程が、前記本発明の候補選択方法を含むことを特徴とする、配列情報群から、相互に類似する類似配列情報群を選択する類似選択方法であ

る。

(A) 配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する工程

(B) 前記候補配列情報群の各候補配列情報を相互に対比し、同一および類似する配列情報を類似配列情報群 (G3) として選択する工程

[0014] 本発明の判定方法は、下記 (X) および (Y) 工程を含み、前記 (X) 工程が、前記本発明の類似選択方法を含むことを特徴とする、目的の類似配列情報群の濃縮の判定方法である。

(X) 配列情報群から、目的配列情報とそれに類似する配列情報とを目的の類似配列情報群として選択する工程

(Y) 前記類似配列情報群における前記目的配列情報と前記類似する配列情報との重複度の合計から、前記類似配列情報群の濃縮を判定する工程

[0015] 本発明のプログラムは、前記本発明の候補選択方法、前記本発明の類似選択方法および前記本発明の判定方法からなる群から選択された少なくとも一つを、コンピュータ上で実行可能なことを特徴とするプログラムである。

[0016] 本発明の記録媒体は、前記本発明のプログラムを記録していることを特徴とする。

発明の効果

[0017] 本発明によれば、配列情報間の類似を判断するにあたって、まず、類似を判断するための候補配列群が選択される。このため、例えば、全ての配列情報間の類似を確認する従来の方法とは異なり、簡便に効率よく類似の判断を行うことができる。このため、例えば、アダマーの濃縮の判定等についても、労力、時間およびコストの軽減が可能となる。

図面の簡単な説明

[0018] [図1]図1は、本発明の候補選択装置の実施形態を示すブロック図である。

[図2]図2は、本発明の候補選択方法および候補選択プログラムの実施形態を示すフローチャートである。

[図3]図3は、本発明の候補選択方法および候補選択プログラムの実施形態を

示すフローチャートである。

[図4]図4は、本発明の類似選択装置の実施形態を示すブロック図である。

[図5]図5は、本発明の類似選択方法および類似選択プログラムの実施形態を説明するためのフローチャートである。

[図6]図6は、本発明の類似選択方法および類似選択プログラムの実施形態を説明するためのフローチャートである。

[図7]図7は、本発明の類似選択装置のその他の実施形態を示すブロック図である。

[図8]図8は、本発明の類似選択方法および類似選択プログラムのその他の実施形態を説明するためのフローチャートである。

[図9]図9は、本発明の類似選択方法および類似選択プログラムのその他の実施形態を説明するためのフローチャートである。

発明を実施するための形態

[0019] 本発明において、「配列情報群」は、複数の配列情報から構成される群を意味し、前記複数の配列情報は、例えば、全て、異なる配列情報でもよいし、同じ配列情報と異なる配列情報とを含んでもよい。本発明は、異なる配列情報間における類似を判断するにあたって、類似判断の候補となる候補配列情報の選択を目的とする。このため、前記複数の配列情報は、例えば、全て、異なる配列情報が好ましい。前記配列情報群に含まれる前記配列情報の個数は、特に制限されない。

[0020] 本発明において、「配列情報」は、特に制限されず、要素の並びに関する情報である。前記要素は、例えば、文字および記号の少なくとも一方があげられ、具体例として、核酸の種類を示す文字または記号、アミノ酸の種類を示す文字または記号等があげられる。核酸の種類を示す文字または記号としては、例えば、A、G、C、TおよびU等の塩基の種類を示す文字または記号があげられる。アミノ酸の種類を示す文字または記号としては、例えば、Me t等の3文字表記、M等の1文字表記の文字または記号があげられる。前記配列情報は、具体例として、核酸配列の配列情報、アミノ酸配列の配列

情報等があげられる。前記配列情報の長さは、前記配列情報を構成する要素の数ともいうことができる。前記配列情報の長さは、特に制限されず、要素が、例えば、5～200個であり、好ましくは、10～150個であり、さらに好ましくは20～120個である。

[0021] 本発明において、「仮想配列情報群」は、複数の仮想配列情報から構成される群を意味する。前記仮想配列情報は、前記配列情報を構成する要素（構成単位ともいう）から構築された仮想の配列情報である。前記要素は、前記配列情報群の配列情報の種類に応じて決定でき、具体的には、前記配列情報群における配列情報と同じ要素である。前記仮想配列情報は、例えば、前記要素を任意に並べた情報ということができ、前記仮想配列情報群は、複数の、任意の異なる並びの情報から構成される群ということができる。前記仮想配列情報の長さは、前記仮想配列情報を構成する要素の数ともいうことができる。前記仮想配列情報の長さは、特に制限されず、要素が、例えば、1～10個であり、好ましくは、1～7個であり、さらに好ましくは1～4個である。前記仮想配列情報群の各仮想配列情報は、例えば、全て同じ長さであることが好ましい。

[0022] 本発明において、前記配列情報群から選択した比較または対比し合う配列情報を、それぞれ、比較元配列情報および比較先配列情報という。ある配列情報に対して、他の配列情報を対比する場合、前者の配列情報を「比較元」ともいい、後者の他の配列情報を「比較先」ともいう。

[0023] 本発明において、「仮想配列情報の頻度」とは、対象となる配列情報において、前記仮想配列情報が出現する頻度を意味し、例えば、頻度ベクトルの要素、出現数ということもできる。また、「頻度の相違」とは、二つ以上の配列情報間の頻度の相違を意味し、例えば、比較先の配列情報の頻度と比較元の配列情報の頻度との相違である。

[0024] 本発明において、「類似度」は、比較元配列情報に対する比較先配列情報の類似の程度を示す。また、本発明において、「類似度の許容条件」は、前記比較元配列情報に対して、前記比較先配列情報が類似判断の候補となり得

ることを示す、類似度の条件である。前記類似度の許容条件は、任意に設定でき、例えば、2つの配列情報を対比した場合に許容できる要素のミスマッチの個数に基づいて設定できる。2つの配列情報の対比とは、例えば、2つの配列情報の要素の並びの対比である。前記類似度の許容条件は、例えば、2つの配列情報を対比した場合に許容できるミスマッチの個数（M）に、前記仮想配列情報の長さ（要素の個数N）を乗じた値を設定できる。

[0025] 本発明において、「重複度」とは、複数の配列情報から構成される配列情報群において、完全に同一である配列情報の個数を意味し、例えば、出現数ということもできる。また、本発明において、「類似重複度」とは、複数の配列情報から構成される配列情報群において、完全に同一である配列情報の重複度と、前記配列情報に類似する他の配列情報の重複度との合計を意味する。前記配列情報に対して、類似する他の配列情報が2つ以上存在する場合、例えば、前記配列情報と、類似する各他の配列情報との間の重複度の合計を、それぞれの類似重複度とする。

[0026] <本発明の候補選択装置および候補選択方法>

本発明の候補選択装置は、前述のように、下記（a）、（b）、（c）および（d）手段を備えることを特徴とする、配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する候補選択装置である。

（a）配列情報群の各配列情報について、仮想配列情報群の各仮想配列情報の頻度をカウントする工程を実行する手段

（b）前記配列情報群から、比較元となる配列情報と比較先となる配列情報とを選択する工程を実行する手段

（c）前記比較元配列情報の前記各仮想配列情報の頻度と、前記比較先配列情報の前記各仮想配列情報の頻度との相違を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程を実行する手段

（d）前記比較元配列情報に対する前記比較先配列情報の類似度が、前記仮想配列情報群に設定した類似度の許容条件を満たす場合、前記比較元配列情報および前記比較先配列情報を、配列情報間の類似を判断する候補配列情報

群として選択する工程を実行する手段

- [0027] 本発明の候補選択装置において、前記仮想配列情報群が、配列情報を構成する要素から構築された仮想配列情報の群であることが好ましい。
- [0028] 本発明の候補選択装置において、前記(c)手段が、下記(c1)および(c2)工程を実行する手段であることが好ましい。
- (c1) 前記仮想配列情報ごとに、前記比較元配列情報における頻度と前記比較先配列情報における頻度との差を求める工程
- (c2) 前記各仮想配列情報の頻度の差のうち、正数の差のみの総和の絶対値または負数の差のみの総和の絶対値を求め、前記絶対値を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程
- [0029] 本発明の候補選択装置において、前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの個数に基づき設定された条件であることが好ましい。2つの配列情報の対比とは、2つの配列情報のアライメントということもできる。
- [0030] 本発明の候補選択装置において、例えば、前記配列情報が、塩基配列であり、前記配列情報を構成する要素が、A、G、C、TおよびUの塩基であることが好ましい。
- [0031] 本発明の候補選択装置において、前記仮想配列情報の塩基長が、例えば、1～10塩基長であることが好ましい。
- [0032] 本発明の候補選択装置において、前記仮想配列情報群の各仮想配列情報が、すべて同じ塩基長であることが好ましい。
- [0033] 本発明の候補選択装置において、前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの塩基数に基づき設定された条件であることが好ましい。
- [0034] 本発明の候補選択装置において、前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの塩基数(M)に前記仮想配列情報の塩基長(N)を乗じた値であることが好ましい。
- [0035] 本発明の候補選択装置は、さらに、下記(e)手段を有することが好まし

い。

(e) 前記 (b)、(c) および (d) 手段による各工程の反復を実行する手段

この場合、前記 (b) 手段は、例えば、前記工程の実行ごとに、前記配列情報群から、異なる配列情報を前記比較元配列情報として選択することが好ましい。

[0036] 本発明の候補選択方法は、前述のように、下記 (a)、(b)、(c) および (d) 工程を含むことを特徴とする、配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する候補選択方法である。本発明の候補選択方法は、特に示さない限り、前記本発明の候補選択装置における説明を援用できる。

(a) 配列情報群の各配列情報について、仮想配列情報群の各仮想配列情報の頻度をカウントする工程

(b) 前記配列情報群から、比較元となる配列情報と比較先となる配列情報とを選択する工程

(c) 前記比較元配列情報の前記各仮想配列情報の頻度と、前記比較先配列情報の前記各仮想配列情報の頻度との相違を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程

(d) 前記比較元配列情報に対する前記比較先配列情報の類似度が、前記仮想配列情報群に設定した類似度の許容条件を満たす場合、前記比較元配列情報および前記比較先配列情報を、配列情報間の類似を判断する候補配列情報群として選択する工程

[0037] 本発明の候補選択方法は、前記仮想配列情報群が、配列情報を構成する要素から構築された仮想配列情報の群であることが好ましい。

[0038] 本発明の候補選択方法は、前記 (c) 工程が、下記 (c1) および (c2) 工程を含むことが好ましい。

(c1) 前記仮想配列情報ごとに、前記比較元配列情報における頻度と前記比較先配列情報における頻度との差を求める工程

(c 2) 前記各仮想配列情報の頻度の差のうち、正数の差のみの総和の絶対値または負数の差のみの総和の絶対値を求め、前記絶対値を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程

[0039] 本発明の候補選択方法は、前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの個数に基づき設定された条件であることが好ましい。

[0040] 本発明の候補選択方法は、前記配列情報が、塩基配列であり、前記配列情報を構成する要素が、A、G、C、TおよびUの塩基であることが好ましい。

[0041] 本発明の候補選択方法は、前記仮想配列情報の塩基長が、1～10塩基長であることが好ましい。

[0042] 本発明の候補選択方法は、前記仮想配列情報群の各仮想配列情報が、すべて同じ塩基長であることが好ましい。

[0043] 本発明の候補選択方法は、前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの塩基数に基づき設定された条件であることが好ましい。

[0044] 本発明の候補選択方法は、前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの塩基数(M)に前記仮想配列情報の塩基長(N)を乗じた値であることが好ましい。

[0045] 本発明の候補選択方法は、さらに、下記(e)工程を含むことが好ましい。この場合、前記(b)工程において、前記工程の実行ごとに、前記配列情報群から、異なる配列情報を前記比較元配列情報として選択することが好ましい。

(e) 前記(b)、(c)および(d)工程を反復する工程

[0046] 本発明の候補選択方法は、前記各工程が、全て、コンピュータ上で実行されることが好ましい。本発明の候補選択方法は、例えば、前記各工程が、全て、前記本発明の候補選択装置により実行されてもよい。

[0047] 以下、図面を参照しながら本発明のさらに具体的な実施形態について説明

する。ただし、本発明は、以下の実施形態に限定されない。以下、配列情報を配列、配列情報群を配列群と示す。

[0048] [実施形態 1]

実施形態 1 は、本発明の候補選択装置および候補選択方法に関する。本実施形態は、前記配列として、核酸の塩基配列を使用する一例である。

[0049] 本実施形態によれば、複数の塩基配列からなる塩基配列群から、塩基配列間の類似の判断候補となる候補配列群を選択できる。

[0050] 図 1 に、本実施形態の候補選択装置の一例の構成を示す。図 1 に示すように、候補選択装置 10 は、入力手段 11、配列記憶部 121、類似度記憶部 122 および候補配列記憶部 123、類似度計算手段 131 および候補配列選択手段 132、ならびに出力手段 14 を備える。類似度計算手段 131 および候補配列選択手段 132 は、例えば、図 1 に示すように、ハードウェアであるデータ処理手段（データ処理装置） 13 に組み込まれてもよく、ソフトウェアまたは前記ソフトウェアが組み込まれたハードウェアでもよい。各記憶部 121、122、123 は、例えば、図 1 に示すように、ハードウェアである記憶手段 12 に組み込まれてもよい。データ処理手段 13 は、CPU 等を備えてもよい。

[0051] 配列記憶部 121 は、入力手段 11 および類似度計算手段 131 と、類似度記憶部 122 は、類似度計算手段 131 および候補配列選択手段 132 と、候補配列記憶部 123 は、候補配列選択手段 132 および出力手段 14 と、それぞれ電氣的に接続されている。また、入力手段 11 は、類似度計算手段 131 と、類似度計算手段 131 は、候補配列選択手段 132 と、候補配列選択手段 132 は、出力手段 14 と、それぞれ電氣的に接続されてよい。候補選択装置 10 は、例えば、情報を記憶手段 12 に記憶させ、記憶させた情報をデータ処理手段 13 に出力してデータ処理を行ってもよいし、前記情報をデータ処理手段 13 に入力してデータ処理を行ってもよい。

[0052] 入力手段 11 は、配列群および仮想配列群の情報を入力する手段（入力装置）である。入力手段 11 は、特に制限されず、例えば、キーボード、マウ

ス等のコンピュータに備わる通常の入力手段、入力ファイルおよび他のコンピュータ等を用いることができる。入力手段 11 は、例えば、データベースに格納された、前記配列群および仮想配列群の情報を読み出す手段でもよい。この場合、例えば、予めサーバに格納された配列情報が、回線網を通じて、入力手段 11 に呼び出される。また、入力手段 11 は、例えば、通信インターフェースを含んでもよい。

[0053] 前記配列群における入力する配列の数は、特に制限されず、下限は、例えば、5 個、好ましくは 10 個であり、上限は、例えば、1000 万個、好ましくは 100 万個である。入力する配列の情報項目は、例えば、配列を構成する要素の順序、すなわち塩基の並びである。前記配列の長さは、特に制限されず、例えば、5~200 塩基長であり、好ましくは、10~150 塩基長であり、さらに好ましくは 20~120 塩基長である。

[0054] 前記仮想配列群の仮想配列の数は、特に制限されず、前記仮想配列の塩基長に応じて適宜決定できる。前記塩基長は、その下限が、例えば、1 塩基長であり、好ましくは 2 塩基長であり、より好ましくは 3 塩基長であり、その上限が、例えば、10 塩基長であり、好ましくは 9 塩基長であり、より好ましくは 8 塩基長であり、さらに好ましくは 7 塩基長である。前記仮想配列群において、前記各仮想配列の長さは、全て同じ長さが好ましい。

[0055] 前記仮想配列を構成する要素が 4 つの塩基 (A、C、G、および T または U) であり、前記仮想配列の塩基長が n (正数) の場合、前記仮想配列群における前記仮想配列の数は、例えば、4 の n 乗個 (4^n 個) である。具体例として、前記要素が 4 つの塩基 A、C、G および T の場合、前記 1 塩基長の仮想配列の数は、4 の 1 乗、つまり、A、C、G および T の 4 個であり、前記 2 塩基長の仮想配列の数は、4 の 2 乗、つまり、AA、AC、AG、AT、CC、CA、CG、CT、GG、GA、GC、GT、TT、TA、TC、TG の 16 個である。

[0056] 類似度計算手段 131 は、前記 (a) 工程として、前記配列群の各配列について各仮想配列群の頻度のカウント、前記 (b) 工程として、前記配列群

からの比較元配列と比較先配列との選択、前記(c)工程として、前記比較元配列に対する前記比較先配列の類似度の計算を行う。前記(a)、(b)および(c)工程の順序は、特に制限されず、順不同である。

[0057] 前記(c)工程における前記類似度の計算は、前述のように、前記(c1)として、前記仮想配列ごとに、前記比較元配列における頻度(S_n)と前記比較先配列における頻度(T_n)との差($S_n - T_n$)を求め、前記(c2)工程として、前記頻度の差($S_n - T_n$)のうち、正数の差のみの総和の絶対値または負数の差のみの総和の絶対値を求めることを行える。すなわち、前記総和の絶対値を、前記類似度とする。

[0058] 候補配列選択手段132は、前記比較元配列に対する前記比較先配列の類似度と、前記仮想配列群に設定した類似度の許容条件とに基づいて、配列情報間の類似を判断する候補配列の選択を行う。ここで選択された複数の候補配列が、候補配列群となる。

[0059] 前記類似度の許容条件は、2つの配列を対比した場合に許容できるミスマッチの塩基数に基づき設定でき、具体例として、前記許容できるミスマッチの塩基数(M)に前記仮想配列の塩基長(N)を乗じた値($N \times M$)があげられる。例えば、塩基長 $N = 1$ の前記仮想配列(A、C、GおよびT)であって、前記許容できるミスマッチの塩基数 $M = 2$ に設定した場合、許容条件($N \times M$)は、 $1 \times 2 = 2$ となる。そして、前記類似度が2以下の場合、許容条件の数値以下となり許容条件を満たすため、前記比較元配列および前記比較先配列は、配列情報間の類似を判断する候補配列として選択する。他方、前記類似度が2を超える場合、許容条件の数値を超え許容条件を満たさないため、前記比較先配列は、前記比較元配列との類似を判断する候補配列として選択しない。

[0060] 前記許容条件の一例として、前記許容できるミスマッチの塩基数(M)に前記仮想配列の塩基長(N)を乗じた値($N \times M$)を設定するのは、以下の理由による。例えば、以下の2つの配列をアラインメントした場合、大文字の1塩基がミスマッチである。これらの配列について、塩基長 $N = 2$ の仮想

配列の頻度をカウントした場合、対象元配列 Seq 1 において、下線部が c g および g g とカウントされるのに対し、対象先配列 Seq 2 において、下線部が c A および A g とカウントされる。つまり、許容できるミスマッチの塩基数が 1 であっても、1 つミスマッチの存在によって、カウントされる仮想配列は、最大 2 つが変動することになる。このため、前記許容できるミスマッチの塩基数 (M) に、前記仮想配列の塩基長 (N) を乗じることで、カウントへの影響を補正できる。

対象元配列 Seq 1 : a a c c g g t t

対象先配列 Seq 2 : a a c c A g t t

[0061] 出力手段 (出力装置) 14 は、候補配列選択手段 132 の結果を出力する手段であればよい。また、前記出力手段 14 は、候補配列記憶部 123 に記憶された情報を出力する手段でもよい。前記出力手段 14 は、特に制限されず、例えば、ディスプレイ装置、印刷装置等のコンピュータに備わる通常の出力量、出力ファイル、および、他のコンピュータ等を使用できる。

[0062] つぎに、図 2 および図 3 のフローチャートを参照し、本実施形態の候補選択方法を説明する。本実施形態の候補選択方法は、A1 ステップ (配列入力)、A2 ステップ (類似度計算) および A3 ステップ (候補配列選択) を含む。

[0063] (A1) 配列入力

配列群の各配列および仮想配列群の各仮想配列を、それぞれ入力し、配列記憶部 121 に記憶させる。前記配列群および前記仮想配列群の情報項目は、例えば、配列における塩基の順序があげられる。

[0064] (A2) 類似度計算

前記配列群から、新しい比較元配列のセット (A21) および新しい比較先配列のセット (A22) を行い、セットした前記比較元配列と前記比較先配列について、それぞれ、前記各仮想配列の頻度をカウントする。そして、各仮想配列について、前記比較元配列の頻度と前記比較先配列の頻度との差を求め、正数の差のみの総和または負数の差のみの総和を計算する。具体的

には、 n 個 (n は正数)の仮想配列が存在する場合、前記比較元配列について、各仮想配列の頻度として n 個の頻度 (S_1, \dots, S_n)、前記比較先配列について、 n 個の頻度 (T_1, \dots, T_n)が得られる。そして、各仮想配列の頻度について、前記比較元配列と前記比較先配列との差、すなわち、 $(S_1 - T_1), \dots, (S_n - T_n)$ を求め、正数の差のみの総和または負数の差のみの総和を計算し、総和の絶対値を求める。前記総和の絶対値が、前記比較元配列に対する前記比較先配列の類似度である。

[0065] (A3) 候補配列選択

そして、前記類似度が、類似度の許容値を満たすか否か、つまり、許容値よりも大きいか否かを判断する(A31)。NOの場合、つまり、前記類似度が許容値よりも小さい場合、前記比較先配列は、前記比較元配列に対して許容できる数のミスマッチを有すると判断して、前記比較元配列と前記比較先配列が類似判断の候補配列であるとの結果を出力する(A32)。他方、YESの場合、つまり、前記類似度が許容値よりも大きい場合、前記比較先配列は、前記比較元配列に対して許容できない数のミスマッチを有すると判断して、前記比較先配列が類似候補配列ではないとの結果を出力する(A33)。

[0066] その後は、未比較の比較先配列の有無を確認する(A34)。YESの場合、つまり、未比較の比較先配列がある場合、A22ステップから同様の処理を行う。そして、NOの場合、つまり、未比較の比較先配列がない場合、さらに、未比較の比較元配列の有無を確認する(A35)。YESの場合、つまり、未比較の比較元配列がある場合、A21ステップから同様の処理を行い、NOの場合、つまり、未比較の比較元配列がない場合、終了する。なお、ある配列を比較元配列とし他の配列を比較先配列として比較済みである場合、前者を比較先配列とし後者を比較元配列とする比較は、省略し、比較済みの結果を使用してもよい。

[0067] 前記A2ステップおよびA3ステップについて、さらなる具体例として、前記仮想配列が塩基長1の場合を例にあげて説明する。

[0068] 塩基長 $N = 1$ の仮想配列を下記 4 種類、比較元配列を Seq 3、比較先配列を Seq 4 と仮定する。そして、2 つの配列をアラインメントした場合に、類似の判断候補として許容できるミスマッチの塩基数を M とし、許容値を $N \times M = 1 \times M = M$ とする。

仮想配列：A、C、G および T

比較元配列 Seq 3：ACGTACGT

比較先配列 Seq 4：AAGAACA T

[0069] 比較元配列 Seq 3 および比較先配列 Seq 4 における各仮想配列 (A、C、G、T) の頻度 $\{f_A, f_C, f_G, f_T\}$ は、それぞれ、SEQ 1 が $\{2, 2, 2, 2\}$ および Seq 2 が $\{5, 1, 1, 1\}$ となり、各頻度 $\{f_A, f_C, f_G, f_T\}$ の差は、A が $(2 - 5 = -3)$ 、C が $(2 - 1 = 1)$ 、G が $(2 - 1 = 1)$ 、T が $(2 - 1 = 1)$ となる。負数の差の総数 $(-3 + 0 + 0 + 0 = -3)$ の絶対値は 3 であり、正数の差の総数 $(0 + 1 + 1 + 1 = 3)$ の絶対値は 3 である。この絶対値 3 が、比較元配列 Seq 3 に対する比較先配列 Seq 4 の類似度であり、比較先配列 Seq 4 が、比較元配列 Seq 3 とアラインメントした際に、少なくとも 3 つのミスマッチを有することを示す。前記許容できるミスマッチの上限塩基数 M を、例えば、2 とした場合、許容値は $N \times M = 1 \times 2 = 2$ である。このため、計算した類似度と許容値とを対比すると、類似度 $3 >$ 許容値 2 であるため、比較先配列 Seq 4 は、比較元配列 Seq 3 の類似判断の候補配列からはずす。他方、前記許容できるミスマッチの上限塩基数 M を、例えば、3 とした場合、許容値は $N \times M = 1 \times 3 = 3$ である。このため、計算した類似度と許容値とを対比すると、類似度 $3 =$ 許容値 3 であるため、比較先配列 Seq 4 は、比較元配列 Seq 3 の類似判断の候補配列として選択する。

[0070] このようにして、前記比較先配列が前記許容条件を満たす場合には、前記比較先配列は、前記比較元配列と共に、類似判断の候補配列として選択する。つまり、候補配列群として選択する。他方、前記比較先配列が前記許容条件を満たさない場合には、前記比較先配列は、類似判断の候補配列として選

択しない。また、前記比較元配列に対して、前記許容条件を満たす比較先配列が存在しない場合は、前記比較元配列も、類似判断の候補配列として選択しない。

[0071] 本実施形態における候補選択装置 10 において、入力手段 11 と類似度計算手段 131、類似度計算手段 131 と候補配列選択手段 132 が、それぞれ電氣的に接続されてもよい。また、候補選択装置 10 は、例えば、各種記憶部を備えてもよいし、備えていなくてもよい。この場合、例えば、入力手段 11 により入力された各配列について、類似度計算手段 131 により類似度を計算し、計算された類似度について、候補配列選択手段 132 により候補配列の選択を行ってもよい。

[0072] <本発明の類似選択装置および類似選択方法>

本発明の類似選択装置は、前述のように、下記 (A) および (B) 手段を備え、

前記 (A) 手段が、前記本発明の候補選択装置であることを特徴とする、配列情報群から、相互に類似する類似配列情報群を選択する類似選択装置である。

(A) 配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する工程を実行する手段

(B) 前記候補配列情報群の各候補配列情報を相互に対比し、同一および類似する配列情報を類似配列情報群 (G3) として選択する工程を実行する手段

[0073] 本発明の類似選択装置において、前記 (A) 手段は、前記本発明の候補選択装置であればよく、前記本発明の候補選択装置の記載を援用できる。

[0074] 本発明の類似選択装置は、前記配列情報群が、同一の配列情報および異なる配列情報からなる配列情報群 (G) から選択された前記異なる配列情報の群であることが好ましい。

[0075] 本発明の類似選択装置は、前記 (B) 手段が、下記 (B1)、(B2)、(B3)、(B4) および (B5) 工程を実行する手段であることが好まし

い。

(B 1) 前記候補配列情報群から、比較元となる候補配列情報と比較先となる候補配列情報とを選択する工程

(B 2) 前記比較元候補配列情報に対する前記比較先候補配列情報の類似の有無を決定する工程

(B 3) 前記比較元候補配列情報の重複度と、前記比較元候補配列情報に類似する前記比較先候補配列情報の重複度とを合計し、得られた合計値を、前記比較元候補配列情報の類似重複度とする工程

(B 4) 前記候補配列情報群から、異なる候補配列情報を、新たな比較元となる候補配列情報として選択し、前記(B 1)、(B 2)および(B 3)工程を反復する工程

(B 5) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群(G 3)として選択する工程

[0076] 前記(B 2)工程において、前記比較元候補配列と前記比較先候補配列との類似の有無は、特に制限されず、公知の方法で決定できる。具体的には、配列と配列とをアラインメントして、許容できるミスマッチ(異なる要素)の数に基づき、類似と非類似とを判断できる。具体例として、例えば、前記両配列をアラインメントした際、ミスマッチの数が、前記許容できるミスマッチの数を超える場合は非類似、前記許容できるミスマッチの数以下の場合には類似と判断できる。前記許容できるミスマッチの個数は、特に制限されず、任意に決定できる。

[0077] 重複度は、後の工程が繰り返される間に、0に再設定される。そこで、前記(B 3)工程における重複度は、各配列の初期の情報であることから、「初期重複度」ともいう。また、後の工程において再設定した重複度0は、「重複度0」または「再設定重複度」ともいう。

[0078] 本発明の類似選択装置は、前記(B)手段が、さらに、下記(B 6)、(B 7)および(B 8)工程を実行する手段であることが好ましい。類似重複

度の再算出とは、例えば、すでに得られた類似重複度をリセットし、新たに類似重複度を算出することを意味する。

(B 6) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報の重複度および前記候補配列情報に類似する候補配列情報の重複度を 0 に再設定する工程

(B 7) 重複度が 0 以外である他の候補配列情報について、類似重複度を再算出する工程

(B 8) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群として再選択する工程

[0079] 本発明の類似選択装置は、前記 (B) 手段が、さらに、下記 (B 9) の工程を実行する手段であることが好ましい。

(B 9) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報の重複度を 0 に再設定し、前記 (B 7) および (B 8) 工程を反復する工程

[0080] このように、最も大きな類似重複度に基づく類似候補群の選択と、類似重複度の再計算とを繰り返すことによって、複数の類似配列情報群が選択できる。前記類似配列情報群の再選択は、例えば、全ての候補配列について重複度が 0 に再設定されるまで行うことが好ましい。

[0081] 本発明の類似選択装置は、前記 (B) 手段が、前記 (B 1) 工程における前記比較元補配列情報と前記比較先候補配列情報との組合せとして、すでに実行した組合せの除外を実行することが好ましい。

[0082] 本発明の類似選択装置において、配列情報の情報項目として、例えば、配列を構成する要素の順序の他に、前記各配列の重複度を含んでもよい。この場合、前記配列群に含まれる配列は、全て、異なる配列であることが好ましい。また、配列情報の情報項目として、前記重複度を含まない場合、例えば、前記重複度をカウントする工程を実行する、下記 (B') 手段を含んでもよい。この場合、前記配列群に含まれる配列は、例えば、異なる配列の他に

、完全に要素の順序が同じである配列を含んでもよい。

(B') 前記配列情報群について、完全に同一な配列情報の数を重複度としてカウントする工程を実施する手段

[0083] 本発明の類似選択方法は、前述のように、下記(A)および(B)工程を含み、

前記(A)工程が、前記本発明の候補選択方法を含むことを特徴とする、配列情報群から、相互に類似する類似配列情報群を選択する類似選択方法である。

(A) 配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する工程

(B) 前記候補配列情報群の各候補配列情報を相互に対比し、同一および類似する配列情報を類似配列情報群(G3)として選択する工程

[0084] 本発明の類似選択方法は、前記(B)工程が、下記(B1)、(B2)、(B3)、(B4)および(B5)工程を含むことが好ましい。

(B1) 前記候補配列情報群から、比較元となる候補配列情報と比較先となる候補配列情報とを選択する工程

(B2) 前記比較元候補配列情報に対する前記比較先候補配列情報の類似の有無を決定する工程

(B3) 前記比較元候補配列情報の重複度と、前記比較元候補配列情報に類似する前記比較先候補配列情報の重複度とを合計し、得られた合計値を、前記比較元候補配列情報の類似重複度とする工程

(B4) 前記候補配列情報群から、異なる候補配列情報を、新たな比較元となる候補配列情報として選択し、前記(B1)、(B2)および(B3)工程を反復する工程

(B5) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群(G3)として選択する工程

[0085] 本発明の類似選択方法は、前記(B)工程が、さらに、下記(B6)、(

B 7) および (B 8) 工程を含むことが好ましい。

(B 6) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報の重複度および前記候補配列情報に類似する候補配列情報の重複度を 0 に再設定する工程

(B 7) 重複度が 0 以外である他の候補配列情報について、類似重複度を再算出する工程

(B 8) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群として再選択する工程

[0086] 本発明の類似選択方法は、前記 (B) 工程が、さらに、下記 (B 9) 工程を含むことが好ましい。

(B 9) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報の重複度を 0 に再設定し、前記 (B 7) および (B 8) 工程を反復する工程

[0087] 本発明の類似選択方法は、前記 (B) 工程において、前記 (B 1) 工程における前記比較元補配列情報と前記比較先候補配列情報との組合せとして、すでに実行した組合せを除外することが好ましい。

[0088] 本発明の類似選択方法は、前記各工程が、全て、コンピュータ上で実行されることが好ましい。本発明の類似選択方法は、例えば、前記各工程が、全て、前記本発明の類似選択装置により実行されてもよい。

[0089] 以下、図面を参照しながら本発明のさらに具体的な実施形態について説明する。ただし、本発明は、以下の実施形態に限定されない。また、本実施形態において、前記候補配列群の選択は、前記実施形態 1 の記載を援用できる。以下、配列情報を配列、配列情報群を配列群と示す。

[0090] [実施形態 2]

実施形態 2 は、本発明の類似選択装置および類似選択方法に関する。本実施形態は、前記配列として、核酸の塩基配列を使用する一例である。本実施形態は、特に示さない限り、実施形態 1 の記載を援用できる。

- [0091] 本実施形態によれば、複数の塩基配列からなる塩基配列群から、塩基配列間の類似の判断候補となる候補配列を選択し、複数の前記候補配列からなる候補配列群から、相互に類似する類似配列を類似配列群として選択できる。
- [0092] 図4に、本実施形態の類似選択装置の一例を示す。図4において、図1の候補選択装置10と同じ箇所には、同じ符号を付している。図4に示すように、類似選択装置20は、入力手段11、配列記憶部121、類似度記憶部122、候補配列記憶部123および類似配列記憶部124、類似度計算手段131、候補配列選択手段132および類似配列選択手段133、ならびに出力手段14を備える。類似度計算手段131、候補配列選択手段132および類似配列選択手段133は、例えば、図4に示すように、ハードウェアであるデータ処理手段13に組み込まれてもよく、ソフトウェアまたは前記ソフトウェアが組み込まれたハードウェアでもよい。各記憶部121、122、123、124は、例えば、図4に示すように、ハードウェアである記憶手段12に組み込まれてもよい。データ処理手段13は、CPU等を備えてもよい。
- [0093] 候補配列記憶部123は、さらに、類似配列選択手段133と電気的に接続されており、類似配列記憶部124は、類似配列選択手段133および出力手段14と、それぞれ電気的に接続されている。また、候補配列選択手段132は、類似配列選択手段133と、類似配列選択手段133は、出力手段14と、それぞれ電気的に接続されてよい。類似選択装置20は、例えば、情報を記憶手段12に記憶させ、記憶させた情報をデータ処理手段13に出力してデータ処理を行ってもよいし、前記情報をデータ処理手段13に入力してデータ処理を行ってもよい。
- [0094] 本実施形態において、入力する配列の情報項目は、前述のような、配列を構成する要素の順序の他に、前記各配列の重複度を含むことが好ましい。前記情報項目として、前記重複度を含む場合、前記配列群を構成する配列は、全て、異なる配列であることが好ましい。
- [0095] また、前記情報項目として、前記重複度を含まない場合、例えば、前記（

B') 手段を含んでもよい。前記 (B') 手段により、前記配列群について、完全に同一な配列情報の数を重複度としてカウントできる。

[0096] つぎに、図5および図6のフローチャートを参照し、本実施形態の類似選択方法を説明する。本実施形態の類似選択方法は、A1ステップ（配列入力）、A2ステップ（類似度計算）、A3ステップ（候補配列選択）およびA4ステップ（類似配列選択）を含む。図5において、図2と同じステップには同じ符号を付している。

[0097] 前記A1ステップ、前記A2ステップおよび前記A3ステップは、前記実施形態1と同様に行うことができ、具体的には、前述した図3のフローチャートに従って行うことができる。前記配列入力において、前記配列群の情報項目は、例えば、配列における塩基の順序および配列の重複度があげられ、前記仮想配列群の情報項目は、例えば、配列における塩基の順序があげられる。

[0098] (A4) 類似配列選択

前記A3ステップで選択された候補配列群から、新しい比較元候補配列のセット (A41) および新しい比較先候補配列のセット (A42) を行い、セットした前記比較先候補配列が、前記比較元候補配列に類似するか否かを判断する (A43)。そして、NOの場合、つまり、前記比較先候補配列が、前記比較元候補配列に類似していない場合、前記比較先候補配列は、前記比較元候補配列との類似配列群ではないとの結果を出力する (A44)。他方、YESの場合、つまり、前記比較先候補配列が、前記比較元候補配列に類似している場合、前記比較先候補配列は、前記比較元候補配列との類似配列群であるとの結果を出力する (A45)。

[0099] その後は、前記比較元候補配列に対して、未比較の比較先候補配列の有無を確認する (A46)。YESの場合、つまり、未比較の比較先配列がある場合、A42ステップから同様の処理を行う。そして、NOの場合、つまり、未比較の比較先候補配列がない場合、さらに、未比較の比較元候補配列の有無を確認する (A47)。YESの場合、つまり、未比較の比較元候補配

列がある場合、A 4 1 ステップから同様の処理を行い、N O の場合、つまり、未比較の比較元候補配列がない場合、終了する。なお、ある配列を比較元候補配列とし他の配列を比較先候補配列として比較済みである場合、前者を比較先候補配列とし後者を比較元候補配列とする比較は、省略し、比較済みの結果を使用してもよい。

[0100] このようにして、前記候補配列群における各候補配列から、前記比較元候補配列および前記比較先候補配列を、それぞれ順次セットし、配列間の類似を判断することによって、前記比較元候補配列とそれに類似する比較先候補配列とからなる類似配列群を選択できる。

[0101] 本実施形態における類似選択装置 2 0 において、入力手段 1 1 と類似度計算手段 1 3 1、類似度計算手段 1 3 1 と候補配列選択手段 1 3 2、候補配列選択手段 1 3 2 と類似配列選択手段 1 3 3 とが、それぞれ電氣的に接続されてもよい。また、類似選択装置 2 0 は、例えば、各種記憶部を備えてもよいし、備えていなくてもよい。この場合、例えば、入力手段 1 1 により入力された各配列について、類似度計算手段 1 3 1 により類似度を計算し、計算された類似度について、候補配列選択手段 1 3 2 により候補配列群の選択を行い、さらに、選択された候補配列群について、類似配列選択手段 1 3 3 により類似配列群の選択を行ってもよい。

[0102] [実施形態 3]

実施形態 3 は、実施形態 2 と同様に、本発明の類似選択装置および類似選択方法に関する。本実施形態は、前記実施形態 2 の前記類似配列群の選択において、重複度を用いる一例である。本実施形態は、特に示さない限り、実施形態 1 および 2 の記載を援用できる。

[0103] 本実施形態によれば、配列間の類似度を用いることによって、簡便に、類似配列群を選択できる。

[0104] 図 7 に、本実施形態の類似選択装置の一例を示す。図 7 において、図 4 の類似選択装置 2 0 と同じ箇所には、同じ符号を付している。図 7 に示すように、類似選択装置 3 0 は、類似重複度記憶部 1 2 4 a および類似配列記憶部

1 2 4 b、類似重複度計算手段 1 3 3 a および類似配列選択手段 1 3 3 b を備える。類似重複度計算手段 1 3 3 a および類似配列選択手段 1 3 3 b は、例えば、図 7 に示すように、ハードウェアであるデータ処理手段 1 3 に組み込まれてもよく、ソフトウェアまたは前記ソフトウェアが組み込まれたハードウェアでもよい。類似重複度記憶部 1 2 4 a および類似配列記憶部 1 2 4 b は、例えば、図 7 に示すように、ハードウェアである記憶手段 1 2 に組み込まれてもよい。

[0105] 候補配列記憶部 1 2 3 は、類似重複度計算手段 1 3 3 a と電氣的に接続されており、類似重複度記憶部 1 2 4 a は、類似重複度計算手段 1 3 3 a および類似配列選択手段 1 3 3 b と電氣的に接続されており、類似配列記憶部 1 2 4 b は、類似配列選択手段 1 3 3 b および出力手段 1 4 と、それぞれ電氣的に接続されている。また、候補配列選択手段 1 3 2 は、類似重複度計算手段 1 3 3 a と、類似重複度計算手段 1 3 3 a は、類似配列選択手段 1 3 3 b と、類似配列選択手段 1 3 3 b は、出力手段 1 4 と、それぞれ電氣的に接続されてもよい。

[0106] つぎに、図 8 および図 9 のフローチャートを参照し、本実施形態の類似選択方法を説明する。本実施形態の類似選択方法は、A 1 ステップ（配列入力）、A 2 ステップ（類似度計算）、A 3 ステップ（候補配列選択）および A 4 ステップ（類似配列選択）を含む。本実施形態において、A 4 ステップは、A 4 a ステップ（類似重複度計算）と、A 4 b ステップ（類似重複度の計算結果に基づく類似配列選択）を含む。図 8 および図 9 において、図 5 および図 6 と同じステップには同じ符号を付している。

[0107] 前記 A 1 ステップ、前記 A 2 ステップおよび前記 A 3 ステップは、前記実施形態 2 と同様に行うことができる。本実施形態において、入力する配列の情報項目は、例えば、配列を構成する要素の順序の他に、前記各配列の重複度を含む。

[0108] (A 4) 類似配列選択

前記 A 3 ステップで選択された候補配列群から、新しい比較元候補配列を

セット (A 4 1') し、その重複度が 0 か否かを判断する (A 4 2')。NO の場合、つまり、重複度 0 の場合 (初期重複度が 0 または再設定重複度 0)、再度、新しい比較元候補配列をセットする (A 4 1')。他方、YES の場合、つまり、重複度が 0 でない場合 (初期重複度 ≥ 1)、前記比較元候補配列の重複度をセットする (A 4 3')。そして、新しい比較先候補配列をセット (A 4 4') し、前記比較先候補配列が、前記比較元候補配列に類似するか否かを判断する (A 4 5')。YES の場合、つまり、前記比較先候補配列が前記比較元候補配列に類似する場合、前記比較元候補配列の類似度と前記比較先候補配列の類似度とを合計し、その合計値を類似重複度とする (A 4 6')。この類似重複度は、前記比較元候補配列の類似重複度という。他方、NO の場合、つまり、前記比較先候補配列が、前記比較元候補配列に類似しない場合、未比較の比較先候補配列の有無を確認する (A 4 7')。そして、YES の場合、つまり、未比較の比較先候補配列がある場合、A 4 4' ステップから同様の処理を行う。そして、NO の場合、つまり、未比較の比較先候補配列がない場合、さらに、未比較の比較元候補配列の有無を確認する (A 4 8')。YES の場合、つまり、未比較の比較元候補配列がある場合、A 4 1' ステップから同様の処理を行う。NO の場合、つまり、未比較の比較元候補配列がない場合、最も大きい類似重複度の候補配列以外であって、類似重複度が 0 でない候補配列について、類似重複度をリセット、つまり 0 に再設定する (A 4 9')。さらに、最も大きい類似重複度の候補配列およびそれに類似する候補配列について、重複度を 0 に再設定する (A 4 10')。つぎに、重複度が 0 でない候補配列の有無を確認する (A 4 11')。YES の場合、つまり、重合度が 0 でない候補配列 (初期重複度 ≥ 1) がある場合、これを新しい比較元候補配列とし、A 4 1' ステップから同様の処理を行う。NO の場合、つまり、重複度が 0 でない候補配列が存在しない場合、類似重複度が 0 でない候補配列とそれに類似する候補配列とを類似配列群とし、類似配列群の一覧を出力する (A 4 12')。出力する情報項目は例えば、前記類似配列群に含まれる各配列ならびに類似重複度

等があげられる。

[0109] 前記A4ステップについて、さらなる具体例として、候補配列群に含まれる異なる配列が5種類 (Seq1、Seq2、Seq3、Seq4、Seq5) であり、それぞれの重複度 (すなわち、出現数) が、{5、4、3、2、1} である場合を例にあげて説明する。

[0110] まず、下記表1に、候補配列の種類とその重複度を示す。

[0111] [表1]

		比較先					重複度	類似重複度
		Seq1	Seq2	Seq3	Seq4	Seq5		
比較元	Seq1						5	—
	Seq2						4	—
	Seq3						3	—
	Seq4						2	—
	Seq5						1	—

[0112] つぎに、それぞれの配列間における類似を判断する。下記表2において、類似の関係にあるものを、網掛けで示す。

[0113] [表2]

		比較先					重複度	類似重複度
		Seq1	Seq2	Seq3	Seq4	Seq5		
比較元	Seq1						5	—
	Seq2						4	—
	Seq3						3	—
	Seq4						2	—
	Seq5						1	—

[0114] そして、それぞれの比較元候補配列について、前記比較元候補配列の初期重複度とそれに類似する前記比較先候補配列の初期重複度とを合計し、この合計値を比較元候補配列の類似重複度とする。下記表3に、類似重複度を示す。そして、前記比較元候補配列のうち、最も大きい類似重複度を示す比較元候補配列を選択し、前記比較元候補配列とそれに類似する比較先候補配列とを、類似配列群とする。下記表3において、最も大きい類似重複度11を

示すSeq4ならびにそれに類似するSeq1およびSeq2が、同じ類似配列群となる。

[0115] [表3]

		比較先					重複度	類似 重複度
		Seq1	Seq2	Seq3	Seq4	Seq5		
比 較 元	Seq1	5			2	1	5	8
	Seq2		4	3	2		4	9
	Seq3		4	3			3	7
	Seq4	5	4		2		2	11
	Seq5	5				1	1	6

[0116] 続いて、最も大きい類似重複度を示す比較元候補配列以外であって、類似重複度が0ではない候補配列について、類似重複度をリセットし、最も大きい類似重複度を示す比較元候補配列の初期重複度とそれに類似する比較先候補配列の初期重複度とを、0に再設定する（再設定重複度0）。下記表4において、最も大きい類似重複数11を示すSeq4以外の配列について、類似重複度をリセットし、且つ、Seq4と、それに類似するSeq1およびSeq2の初期重複度を、0に再設定する（再設定重複度0）。

[0117] [表4]

		比較先					重複度	類似 重複度
		Seq1	Seq2	Seq3	Seq4	Seq5		
比 較 元	Seq1	0			0	0	5→0	—
	Seq2		0	0	0		4→0	—
	Seq3		0	3			3	—
	Seq4	0	0		0		2→0	11
	Seq5	0				1	1	—

[0118] そして、重複度が0以外（初期重複度 ≥ 1 ）の比較元候補配列について、同様にして、類似重複度の計算、最も大きい類似重複度に基づく類似候補群の選択を行う。類似候補群の選択は、全ての候補配列の初期重複度が0に再設定されるまで、繰り返し行うことが好ましい。下記表5において、重複度が0ではない候補配列のうち、最も大きい類似重複度3を示すSeq3を、

類似配列群とする。

[0119] [表5]

		比較先					重複度	類似 重複度
		Seq1	Seq2	Seq3	Seq4	Seq5		
比 較 元	Seq1	0			0	0	0	—
	Seq2		0	0	0		0	—
	Seq3		0	3			3	3
	Seq4	0	0		0		0	11
	Seq5	0				1	1	1

[0120] なお、配列間の類似について、一方の配列を比較元候補配列とし、他方の配列を比較先候補配列とするのと、前記一方の配列を比較先候補配列とし、前記他方の配列を比較元候補配列とするのは、実質的に同じといえる。そこで、前記類似配列群の選択をより促進できるため、例えば、比較元候補配列と比較先候補配列との組合せから、すでに実行した組合せを除外することが好ましい。この場合、例えば、下記表6のように、異なる配列間の組合せを半分にできる（セル数の半減）。

[0121] [表6]

		比較先					重複度	類似 重複度
		Seq1	Seq2	Seq3	Seq4	Seq5		
比 較 元	Seq1						5	
	Seq2						4	
	Seq3						3	
	Seq4						2	
	Seq5						1	

[0122] これらの処理を繰り返すことによって、候補配列群を類似配列群に分類することができる。

[0123] <目的の類似配列群の濃縮を判定する装置>

本発明の濃縮の判定装置は、前述のように、下記（X）および（Y）手段を備え、前記（X）手段が、前記本発明の類似選択装置であることを特徴とする、目的の類似配列情報群の濃縮の判定装置である。

(X) 配列情報群から、目的配列情報とそれに類似する配列情報とを目的の類似配列情報群として選択する工程を実行する手段

(Y) 前記類似配列情報群における前記目的配列情報と前記類似する配列情報との重複度の合計から、前記類似配列情報群の濃縮を判定する工程を実行する手段

[0124] 本発明の判定装置において、前記(X)手段は、前記本発明の類似選択装置であればよく、前記本発明の類似選択装置の記載を援用できる。

[0125] 本発明の濃縮の判定装置は、前記(X)手段が、比較元となる類似配列情報群と、比較先となる類似配列情報群を、それぞれ選択する工程を実行し、前記(Y)手段が、下記(Y1)および(Y2)工程を実行する手段であることが好ましい。

(Y1) 前記比較元の類似配列情報群における目的の配列情報とそれに類似する配列情報との重複度の合計と、前記比較先の類似配列情報群における目的の配列情報とそれに類似する配列情報との重複度の合計とを、比較する工程

(Y2) 前記比較元の類似配列情報群における前記重複度の合計が、前記比較先の類似配列情報群における前記重複度の合計よりも大きい場合に、前記比較元の類似配列情報群が、前記比較先の配列情報群よりも、濃縮されていると判断する工程

[0126] 本発明において、濃縮の判定は、例えば、同じ配列情報群に含まれる異なる配列情報について、前記配列情報の間における濃縮度合いの違いを比較することにより行ってもよい。この場合、例えば、前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、同じ配列群から選択された類似配列情報群であり、前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の類似配列情報群の前記目的の配列情報とが、異なる配列情報である。これによって、例えば、同じ配列情報群から、相対的に濃縮度合いの高い配列情報およびその類似配列情報を選択することが可能となる。具体例としては、例えば、アプタマーの調製において、特定のラウンドのライブラリー

に含まれる複数の類似配列情報群から、相対的に濃縮度の高い類似配列情報群の選択、つまり濃縮度が高いアプタマー類似配列群の選択を行うことができる。

[0127] また、前記濃縮の判定は、例えば、異なる配列情報群に含まれる同じ配列情報について、前記配列情報群の間における濃縮度合いの違いを比較することにより行ってもよい。この場合、例えば、前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、異なる配列群から選択された類似配列情報群であり、前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の類似配列情報群の前記目的の配列情報とが、同じ配列情報である。これによって、例えば、特定の配列情報の類似配列情報群について、相対的に濃縮度合いの高い配列情報群を選択することができる。具体例としては、例えば、アプタマーの調製において、各ラウンドのライブラリーのうち、特定のアプタマー類似配列群の濃縮度が相対的に高いライブラリーを選択することができる。

[0128] 本発明の濃縮の判定方法は、下記（X）および（Y）工程を含み、前記（X）工程が、前記本発明の類似選択方法を含むことを特徴とする、類似配列情報群の濃縮の判定方法である。本発明の濃縮の判定方法は、特に示さない限り、前記本発明の濃縮の判定装置における記載を援用できる。

（X）配列情報群から、目的の配列情報とそれに類似する配列情報とを判定対象の類似配列情報群として選択する工程

（Y）前記類似配列情報群における前記目的の配列情報と前記類似する配列情報との重複度の合計から、前記類似配列情報群の濃縮を判定する工程

[0129] 本発明の濃縮の判定方法は、前記（X）工程が、比較元となる類似配列情報群と、比較先となる類似配列情報群を、それぞれ選択する工程であり、前記（Y）工程が、下記（Y1）および（Y2）工程を含むことが好ましい。

（Y1）前記比較元の類似配列情報群における目的の配列情報とそれに類似する配列情報との重複度の合計と、前記比較先の類似配列情報群における目

的の配列情報とそれに類似する配列情報との重複度の合計とを、比較する工程

(Y2) 前記比較元の類似配列情報群における前記重複度の合計が、前記比較先の類似配列情報群における前記重複度の合計よりも大きい場合に、前記比較元の類似配列情報群が、前記比較先の配列情報群よりも、濃縮されていると判断する工程

[0130] 本発明の濃縮の判定方法は、前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、同じ配列群から選択された類似配列情報群であり、前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の類似配列情報群の前記目的の配列情報とが、異なる配列情報であってもよい。

[0131] 本発明の濃縮の判定方法は、前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、異なる配列群から選択された類似配列情報群であり、前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の類似配列情報群の前記目的の配列情報とが、同じ配列情報であってもよい。

[0132] 本発明の用途は、特に制限されないが、例えば、アダマーの調製における濃縮の判定に適用することが好ましい。本発明によれば、前述のように、例えば、同じライブラリー内における異なるアダマー類似配列情報群の濃縮度合いの比較、または、異なるライブラリー内における同じアダマー類似配列情報群の濃縮度合いの比較が可能である。

実施例

[0133] つぎに、本発明の実施例について説明する。ただし、本発明は、下記の実施例により制限されない。

[0134] [実施例1]

本実施例では、低分子化合物をターゲットとするライブラリーについて、本発明の類似選択方法により、類似配列群の分類を行った。

[0135] 配列群として、40塩基長の85,800個の核酸配列群を使用した。仮想配列群の条件、許容できるミスマッチの塩基数および許容条件を下記表7に示す。

[0136] [表7]

	仮想配列		許容ミスマッチ数 (M)	許容条件 (N × M)	計算時間 (時間)
	塩基長 (N)	配列数			
比較例	—	—	—	—	83
実施例	1	4	5	5	17
	2	4 ²	5	10	9
	3	4 ³	5	15	1
	4	4 ⁴	5	20	2

[0137] 実施例は、前記条件に従い、前記表6に示すセル数の半減を行って、候補配列群の選択、類似配列群の選択を行った。これらの計算にかかった時間を前記表7にあわせて示す。なお、比較例は、前記配列群の全ての核酸配列について、アラインメントによる類似の判断を行い、類似配列群を選択した。その結果、実施例によれば、比較例よりも格段に短い計算時間で類似配列群の選択を行うことができた。

[0138] 以上、実施形態を参照して本願発明を説明したが、本願発明は、上記実施形態に限定されるものではない。本願発明の構成や詳細には、本願発明のスコープ内で当業者が理解しうる様々な変更をすることができる。

[0139] この出願は、2013年2月15日に出願された日本出願特願2013-027851を基礎とする優先権を主張し、その開示の全てをここに取り込む。

産業上の利用可能性

[0140] 本発明によれば、配列情報間の類似を判断するにあたって、まず、類似を判断するための候補配列群が選択される。このため、例えば、全ての配列情報間の類似を確認する従来の方法とは異なり、簡便に効率よく類似の判断を行うことができる。このため、例えば、アダマーの濃縮の判定等についても、労力、時間およびコストの軽減が可能となる。

符号の説明

[0141] 10 候補選択装置
20、30 類似選択装置

- 1 1 入力手段
- 1 2 記憶手段
 - 1 2 1 配列記憶部
 - 1 2 2 類似度記憶部
 - 1 2 3 候補配列記憶部
 - 1 2 4 類似配列記憶部
 - 1 2 4 a 類似重複度記憶部
 - 1 2 4 b 類似配列記憶部
- 1 3 データ処理手段
 - 1 3 1 類似度計算手段
 - 1 3 2 候補配列選択手段
 - 1 3 3 類似配列選択手段
 - 1 3 3 a 類似重複度計算手段
 - 1 3 3 b 類似配列選択手段
- 1 4 出力手段

請求の範囲

- [請求項1] 下記（a）、（b）、（c）および（d）手段を備えることを特徴とする、配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する候補選択装置。
- （a）配列情報群の各配列情報について、仮想配列情報群の各仮想配列情報の頻度をカウントする工程を実行する手段
- （b）前記配列情報群から、比較元となる配列情報と比較先となる配列情報とを選択する工程を実行する手段
- （c）前記比較元配列情報の前記各仮想配列情報の頻度と、前記比較先配列情報の前記各仮想配列情報の頻度との相違を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程を実行する手段
- （d）前記比較元配列情報に対する前記比較先配列情報の類似度が、前記仮想配列情報群に設定した類似度の許容条件を満たす場合、前記比較元配列情報および前記比較先配列情報を、配列情報間の類似を判断する候補配列情報群として選択する工程を実行する手段
- [請求項2] 前記仮想配列情報群が、配列情報を構成する要素から構築された仮想配列情報の群である、請求項1記載の候補選択装置。
- [請求項3] 前記（c）手段が、下記（c1）および（c2）工程を実行する手段である、請求項1または2記載の候補選択装置。
- （c1）前記仮想配列情報ごとに、前記比較元配列情報における頻度と前記比較先配列情報における頻度との差を求める工程
- （c2）前記各仮想配列情報の頻度の差のうち、正数の差のみの総和の絶対値または負数の差のみの総和の絶対値を求め、前記絶対値を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程
- [請求項4] 前記類似度の許容条件が、2つの配列情報を対比した場合に許容できる mismatches の個数に基づき設定された条件である、請求項1から3

のいずれか一項に記載の候補選択装置。

- [請求項5] 前記配列情報が、塩基配列であり、前記配列情報を構成する要素が、A、G、C、TおよびUの塩基である、請求項1から4のいずれか一項に記載の候補選択装置。
- [請求項6] 前記仮想配列情報の塩基長が、1～10塩基長である、請求項5記載の候補選択装置。
- [請求項7] 前記仮想配列情報群の各仮想配列情報が、すべて同じ塩基長である、請求項5または6記載の候補選択装置。
- [請求項8] 前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの塩基数に基づき設定された条件である、請求項3から7のいずれか一項に記載の候補選択装置。
- [請求項9] 前記類似度の許容条件が、2つの配列情報を対比した場合に許容できるミスマッチの塩基数(M)に前記仮想配列情報の塩基長(N)を乗じた値である、請求項5から8のいずれか一項に記載の候補選択装置。
- [請求項10] さらに、下記(e)手段を有する、請求項1から9のいずれか一項に記載の候補選択装置。
- (e) 前記(b)、(c)および(d)手段による各工程の反復を実行する手段
- [請求項11] 前記(b)手段は、前記工程の実行ごとに、前記配列情報群から、異なる配列情報を前記比較元配列情報として選択する、請求項10記載の候補選択装置。
- [請求項12] 下記(A)および(B)手段を備え、
前記(A)手段が、請求項1から11のいずれか一項に記載の候補選択装置であることを特徴とする、配列情報群から、相互に類似する類似配列情報群を選択する類似選択装置。
- (A) 配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する工程を実行する手段

(B) 前記候補配列情報群の各候補配列情報を相互に対比し、同一および類似する配列情報を類似配列情報群 (G3) として選択する工程を実行する手段

[請求項13]

前記 (B) 手段が、下記 (B1)、(B2)、(B3)、(B4) および (B5) 工程を実行する手段である、請求項12記載の類似選択装置

(B1) 前記候補配列情報群から、比較元となる候補配列情報と比較先となる候補配列情報とを選択する工程

(B2) 前記比較元候補配列情報に対する前記比較先候補配列情報の類似の有無を決定する工程

(B3) 前記比較元候補配列情報の重複度と、前記比較元候補配列情報に類似する前記比較先候補配列情報の重複度とを合計し、得られた合計値を、前記比較元候補配列情報の類似重複度とする工程

(B4) 前記候補配列情報群から、異なる候補配列情報を、新たな比較元となる候補配列情報として選択し、前記 (B1)、(B2) および (B3) 工程を反復する工程

(B5) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群 (G3) として選択する工程

[請求項14]

前記 (B) 手段が、さらに、下記 (B6)、(B7) および (B8) 工程を実行する手段である、請求項13記載の類似選択装置。

(B6) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報の重複度および前記候補配列情報に類似する候補配列情報の重複度を0に再設定する工程

(B7) 重複度が0以外である他の候補配列情報について、類似重複度を再算出する工程

(B8) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、

類似配列情報群として再選択する工程

[請求項15] 前記（Ｂ）手段が、さらに、下記（Ｂ９）の工程を実行する手段である、請求項１４記載の類似選択装置。

（Ｂ９）前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報の重複度を０に再設定し、前記（Ｂ７）および（Ｂ８）工程を反復する工程

[請求項16] 前記（Ｂ）手段が、前記（Ｂ１）工程における前記比較元補配列情報と前記比較先候補配列情報との組合せとして、すでに実行した組合せの除外を実行する、請求項１３から１５のいずれか一項に記載の類似選択装置。

[請求項17] 下記（Ｘ）および（Ｙ）手段を備え、前記（Ｘ）手段が、請求項１２から１６のいずれか一項に記載の類似選択装置であることを特徴とする、目的の類似配列情報群の濃縮の判定装置。

（Ｘ）配列情報群から、目的配列情報とそれに類似する配列情報とを目的の類似配列情報群として選択する工程を実行する手段

（Ｙ）前記類似配列情報群における前記目的配列情報と前記類似する配列情報との重複度の合計から、前記類似配列情報群の濃縮を判定する工程を実行する手段

[請求項18] 前記（Ｘ）手段が、比較元となる類似配列情報群と、比較先となる類似配列情報群を、それぞれ選択する工程を実行し、

前記（Ｙ）手段が、下記（Ｙ１）および（Ｙ２）工程を実行する手段である、請求項１７記載の判定装置。

（Ｙ１）前記比較元の類似配列情報群における目的の配列情報とそれに類似する配列情報との重複度の合計と、前記比較先の類似配列情報群における目的の配列情報とそれに類似する配列情報との重複度の合計とを、比較する工程

（Ｙ２）前記比較元の類似配列情報群における前記重複度の合計が、

前記比較先の類似配列情報群における前記重複度の合計よりも大きい場合に、前記比較元の類似配列情報群が、前記比較先の配列情報群よりも、濃縮されていると判断する工程

[請求項19] 前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、同じ配列群から選択された類似配列情報群であり、
前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の類似配列情報群の前記目的の配列情報とが、異なる配列情報である、請求項18記載の判定装置。

[請求項20] 前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、異なる配列群から選択された類似配列情報群であり、
前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の類似配列情報群の前記目的の配列情報とが、同じ配列情報である、請求項18記載の判定装置。

[請求項21] 下記(a)、(b)、(c)および(d)工程を含むことを特徴とする、配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する候補選択方法。

(a) 配列情報群の各配列情報について、仮想配列情報群の各仮想配列情報の頻度をカウントする工程

(b) 前記配列情報群から、比較元となる配列情報と比較先となる配列情報とを選択する工程

(c) 前記比較元配列情報の前記各仮想配列情報の頻度と、前記比較先配列情報の前記各仮想配列情報の頻度との相違を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程

(d) 前記比較元配列情報に対する前記比較先配列情報の類似度が、前記仮想配列情報群に設定した類似度の許容条件を満たす場合、前記比較元配列情報および前記比較先配列情報を、配列情報間の類似を判断する候補配列情報群として選択する工程

[請求項22] 前記仮想配列情報群が、配列情報を構成する要素から構築された仮想

配列情報の群である、請求項 2 1 記載の候補選択方法。

[請求項23] 前記 (c) 工程が、下記 (c 1) および (c 2) 工程を含む、請求項 2 1 または 2 2 記載の候補選択方法。

(c 1) 前記仮想配列情報ごとに、前記比較元配列情報における頻度と前記比較先配列情報における頻度との差を求める工程

(c 2) 前記各仮想配列情報の頻度の差のうち、正数の差のみの総和の絶対値または負数の差のみの総和の絶対値を求め、前記絶対値を、前記比較元配列情報に対する前記比較先配列情報の類似度として計算する工程

[請求項24] 前記類似度の許容条件が、2つの配列情報を対比した場合に許容できる mismatches の個数に基づき設定された条件である、請求項 2 1 から 2 3 のいずれか一項に記載の候補選択方法。

[請求項25] 前記配列情報が、塩基配列であり、前記配列情報を構成する要素が、A、G、C、T および U の塩基である、請求項 2 1 から 2 4 のいずれか一項に記載の候補選択方法。

[請求項26] 前記仮想配列情報の塩基長が、1～10塩基長である、請求項 2 5 記載の候補選択方法。

[請求項27] 前記仮想配列情報群の各仮想配列情報が、すべて同じ塩基長である、請求項 2 5 または 2 6 記載の候補選択方法。

[請求項28] 前記類似度の許容条件が、2つの配列情報を対比した場合に許容できる mismatches の塩基数に基づき設定された条件である、請求項 2 3 から 2 7 のいずれか一項に記載の候補選択方法。

[請求項29] 前記類似度の許容条件が、2つの配列情報を対比した場合に許容できる mismatches の塩基数 (M) に前記仮想配列情報の塩基長 (N) を乗じた値である、請求項 2 5 から 2 8 のいずれか一項に記載の候補選択方法。

[請求項30] さらに、下記 (e) 工程を含む、請求項 2 1 から 2 9 のいずれか一項に記載の候補選択方法。

(e) 前記 (b)、(c) および (d) 工程を反復する工程

[請求項31] 前記 (b) 工程において、前記工程の実行ごとに、前記配列情報群から、異なる配列情報を前記比較元配列情報として選択する、請求項30記載の候補選択方法。

[請求項32] 下記 (A) および (B) 工程を含み、
前記 (A) 工程が、請求項21から31のいずれか一項に記載の候補選択方法を含むことを特徴とする、配列情報群から、相互に類似する類似配列情報群を選択する類似選択方法。

(A) 配列情報群から、配列情報間の類似の判断候補となる候補配列情報群を選択する工程

(B) 前記候補配列情報群の各候補配列情報を相互に対比し、同一および類似する配列情報を類似配列情報群 (G3) として選択する工程

[請求項33] 前記 (B) 工程が、下記 (B1)、(B2)、(B3)、(B4) および (B5) 工程を含む、請求項32記載の類似選択方法

(B1) 前記候補配列情報群から、比較元となる候補配列情報と比較先となる候補配列情報とを選択する工程

(B2) 前記比較元候補配列情報に対する前記比較先候補配列情報の類似の有無を決定する工程

(B3) 前記比較元候補配列情報の重複度と、前記比較元候補配列情報に類似する前記比較先候補配列情報の重複度とを合計し、得られた合計値を、前記比較元候補配列情報の類似重複度とする工程

(B4) 前記候補配列情報群から、異なる候補配列情報を、新たな比較元となる候補配列情報として選択し、前記 (B1)、(B2) および (B3) 工程を反復する工程

(B5) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群 (G3) として選択する工程

[請求項34] 前記 (B) 工程が、さらに、下記 (B6)、(B7) および (B8)

工程を含む、請求項 3 3 記載の類似選択方法。

(B 6) 前記候補配列情報のうち、最も大きな類似重複度を示した候補配列情報の重複度および前記候補配列情報に類似する候補配列情報の重複度を 0 に再設定する工程

(B 7) 重複度が 0 以外である他の候補配列情報について、類似重複度を再算出する工程

(B 8) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報を、類似配列情報群として再選択する工程

[請求項35] 前記 (B) 工程が、さらに、下記 (B 9) 工程を含む、請求項 3 4 記載の類似選択方法。

(B 9) 前記他の候補配列情報のうち、最も大きな類似重複度を示した候補配列情報および前記候補配列情報に類似する候補配列情報の重複度を 0 に再設定し、前記 (B 7) および (B 8) 工程を反復する工程

[請求項36] 前記 (B) 工程において、前記 (B 1) 工程における前記比較元補配列情報と前記比較先候補配列情報との組合せとして、すでに実行した組合せを除外する、請求項 3 3 から 3 5 のいずれか一項に記載の類似選択方法。

[請求項37] 下記 (X) および (Y) 工程を含み、前記 (X) 工程が、請求項 3 2 から 3 6 のいずれか一項に記載の類似選択方法を含むことを特徴とする、類似配列情報群の濃縮の判定方法。

(X) 配列情報群から、目的の配列情報とそれに類似する配列情報とを判定対象の類似配列情報群として選択する工程

(Y) 前記類似配列情報群における前記目的の配列情報と前記類似する配列情報との重複度の合計から、前記類似配列情報群の濃縮を判定する工程

[請求項38] 前記 (X) 工程が、比較元となる類似配列情報群と、比較先となる類

似配列情報群を、それぞれ選択する工程であり、
前記（Ｙ）工程が、下記（Ｙ１）および（Ｙ２）工程を含む、請求項
３７記載の判定方法。

（Ｙ１）前記比較元の類似配列情報群における目的の配列情報とそれ
に類似する配列情報との重複度の合計と、前記比較先の類似配列情報
群における目的の配列情報とそれに類似する配列情報との重複度の合
計とを、比較する工程

（Ｙ２）前記比較元の類似配列情報群における前記重複度の合計が、
前記比較先の類似配列情報群における前記重複度の合計よりも大きい
場合に、前記比較元の類似配列情報群が、前記比較先の配列情報群よ
りも、濃縮されていると判断する工程

[請求項39] 前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、同
じ配列群から選択された類似配列情報群であり、
前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の
類似配列情報群の前記目的の配列情報とが、異なる配列情報である、
請求項３８記載の判定方法。

[請求項40] 前記比較元の類似配列情報群と前記比較先の類似配列情報群とが、異
なる配列群から選択された類似配列情報群であり、
前記比較元の類似配列情報群の前記目的の配列情報と、前記比較先の
類似配列情報群の前記目的の配列情報とが、同じ配列情報である、請
求項３８記載の判定方法。

[請求項41] 請求項２１から３１のいずれか一項に記載の候補選択方法を、コンピ
ュータ上で実行可能なことを特徴とするプログラム。

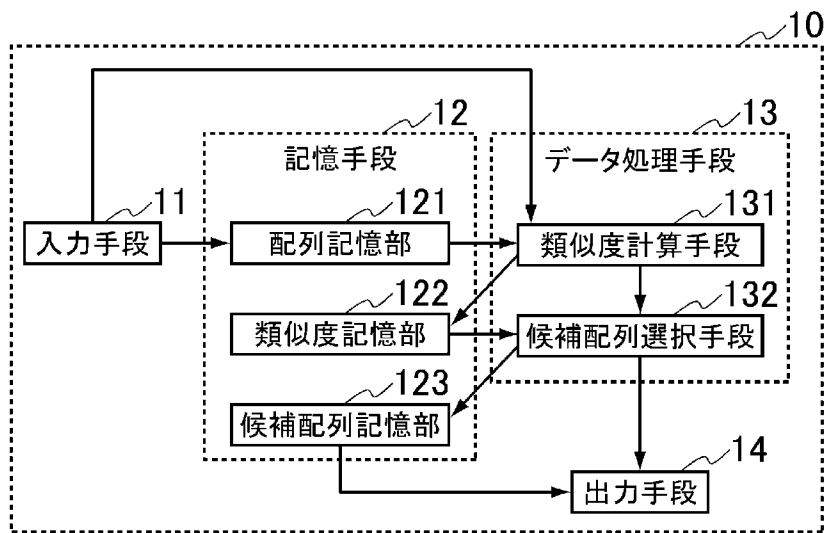
[請求項42] 請求項３２から３６のいずれか一項に記載の類似選択方法を、コンピ
ュータ上で実行可能なことを特徴とするプログラム。

[請求項43] 請求項３７から４０のいずれか一項に記載の判定方法を、コンピ
ュータ上で実行可能なことを特徴とするプログラム。

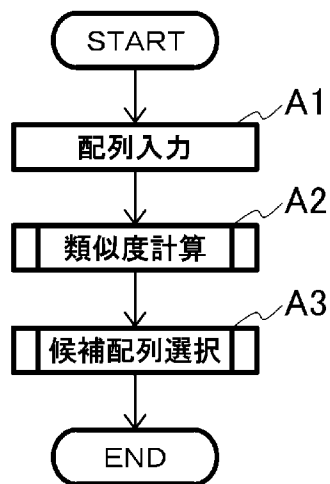
[請求項44] 請求項４１から４３のいずれか一項に記載のプログラムを記録してい

ることを特徴とする記録媒体。

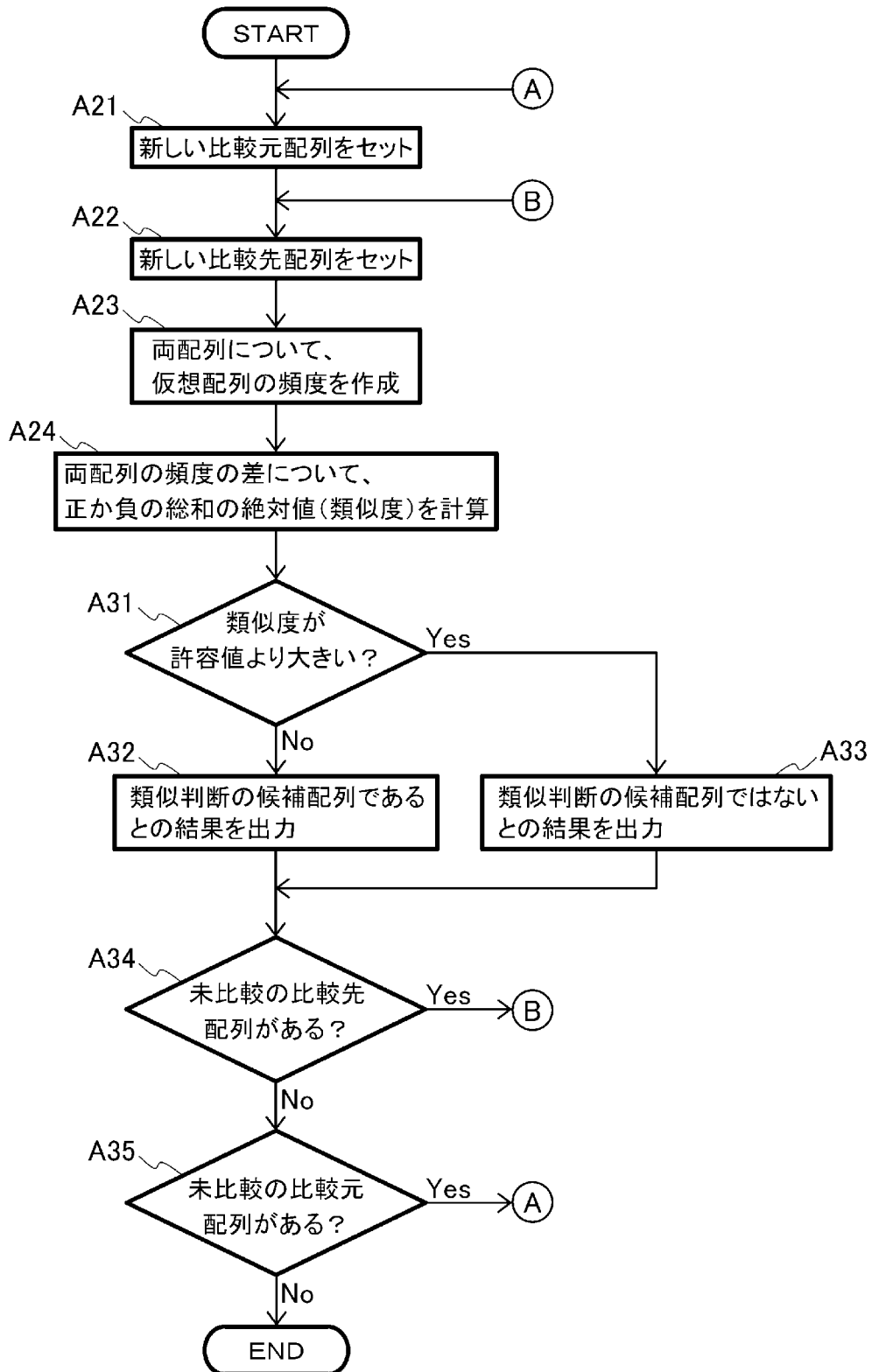
[図1]



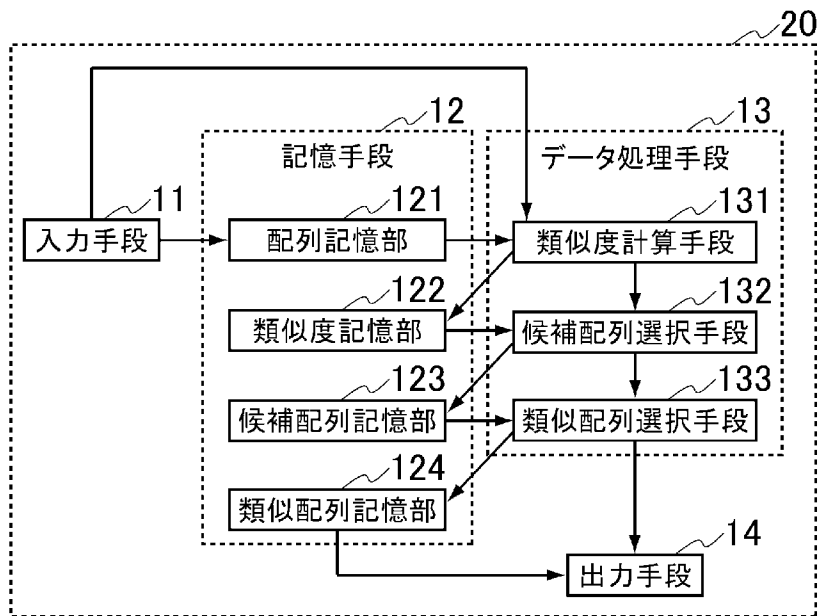
[図2]



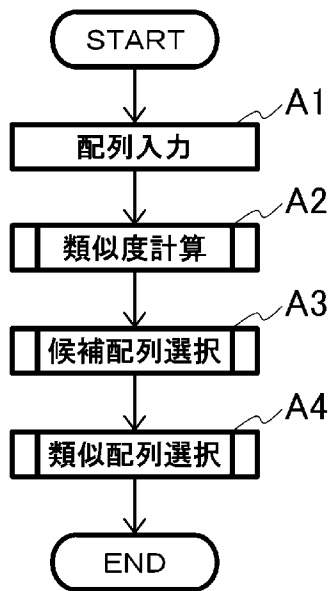
[図3]



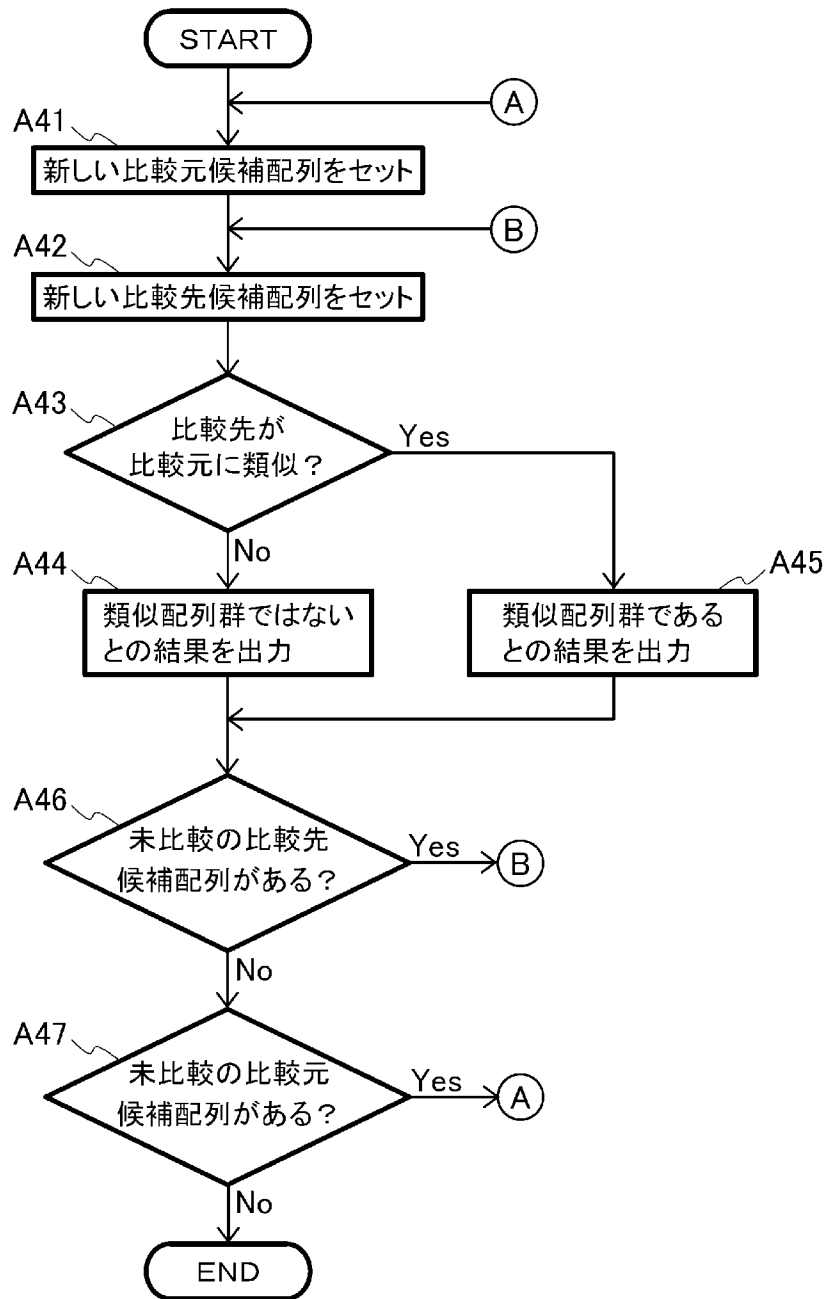
[図4]



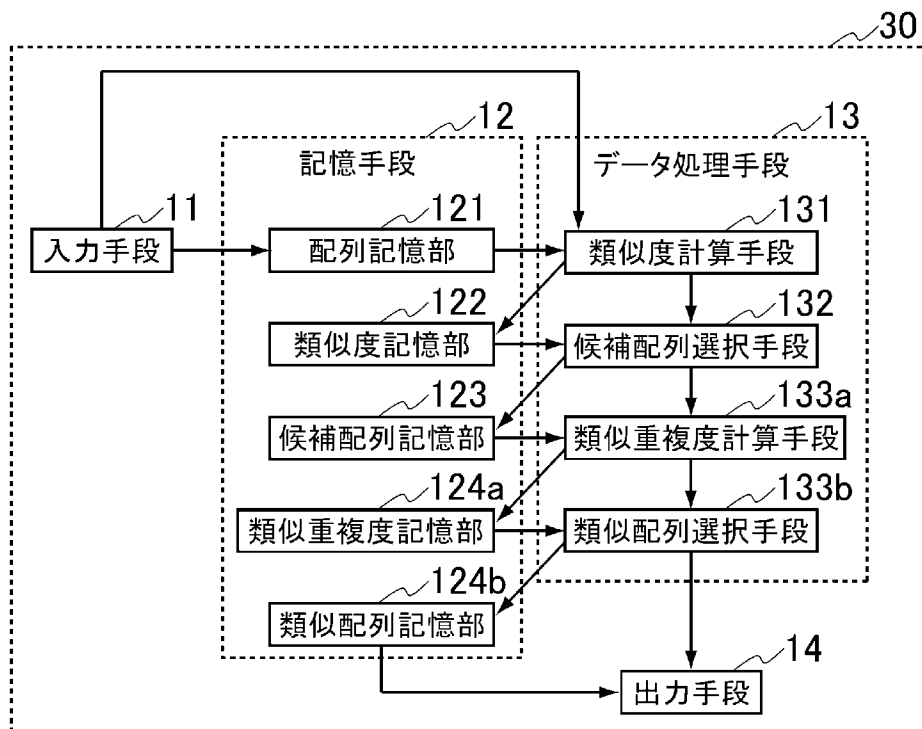
[図5]



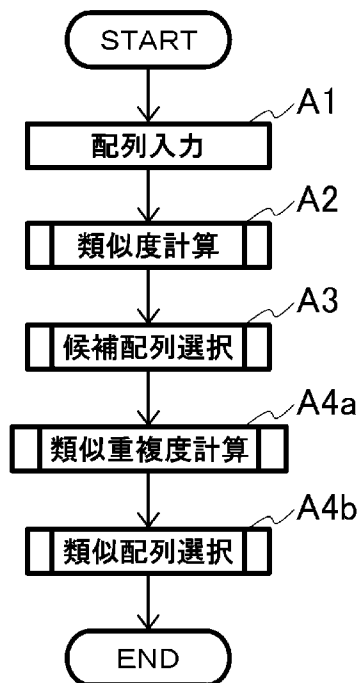
[図6]



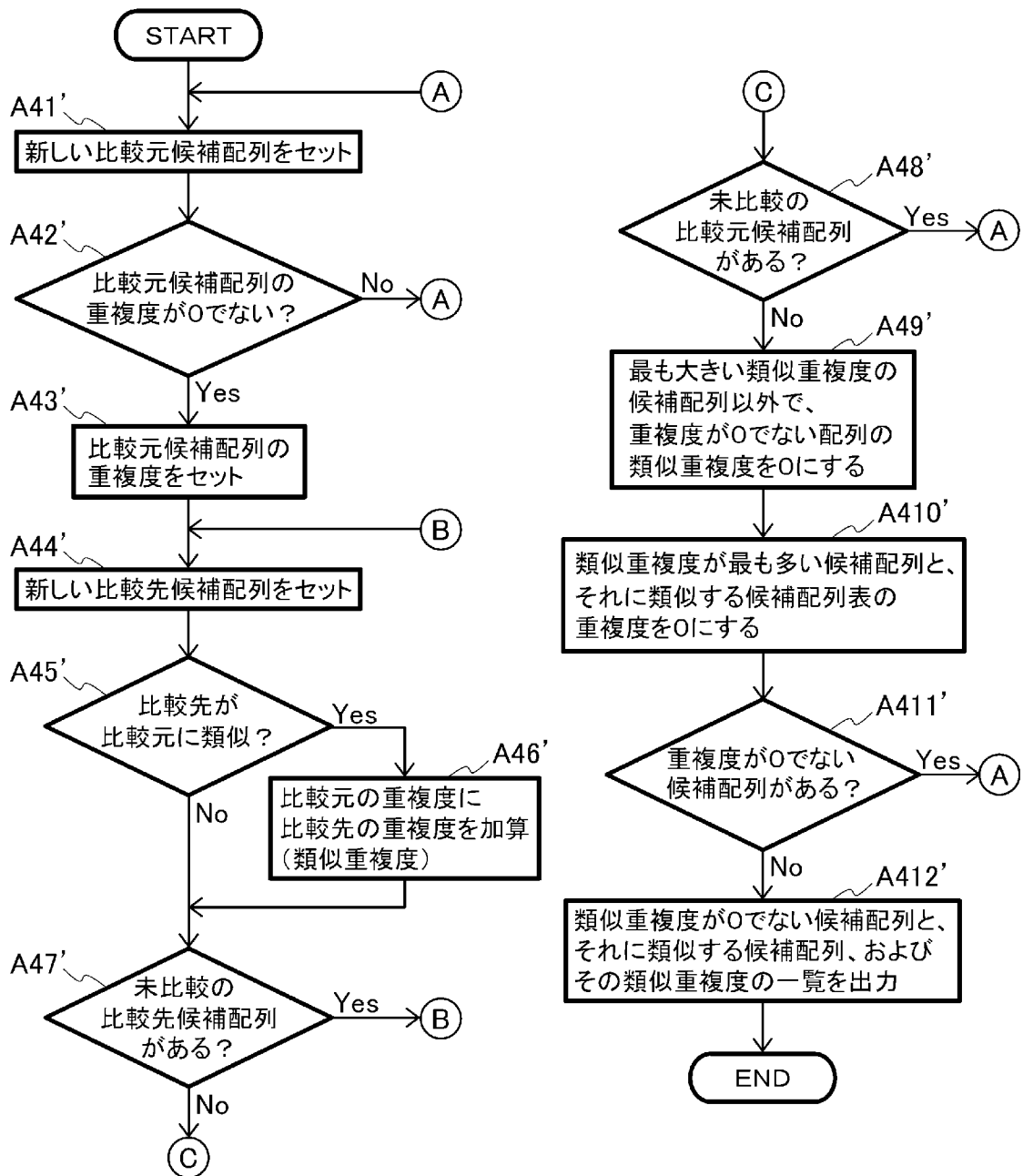
[図7]



[図8]



[図9]



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2014/053516

A. CLASSIFICATION OF SUBJECT MATTER

G06F19/22(2011.01)i, C12N15/115(2010.01)n

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, C12N15/115

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 1996-2014
Kokai Jitsuyo Shinan Koho 1971-2014 Toroku Jitsuyo Shinan Koho 1994-2014

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.
X	JP 2012-146067 A (Japan Software Management Co., Ltd.),	1-11, 21-31, 41, 44
Y	02 August 2012 (02.08.2012), claims 1 to 10; paragraphs [0010] to [0172] & WO 2012/096016 A1 & EP 2665009 A1 & US 2014/019062 A1 & CN 103339632 A	12-20, 32-40, 42, 43
Y	JP 2012-146066 A (Japan Software Management Co., Ltd.), 02 August 2012 (02.08.2012), claims 1 to 6; paragraphs [0010] to [0165] & WO 2012/096015 A1 & EP 2665010 A1 & US 2014/058682 A1 & CN 103348350 A	12-20, 32-40, 42, 43

Further documents are listed in the continuation of Box C. See patent family annex.

* Special categories of cited documents:
 "A" document defining the general state of the art which is not considered to be of particular relevance
 "E" earlier application or patent but published on or after the international filing date
 "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)
 "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
 "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
 "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
 "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
 "&" document member of the same patent family

Date of the actual completion of the international search
09 April, 2014 (09.04.14)

Date of mailing of the international search report
22 April, 2014 (22.04.14)

Name and mailing address of the ISA/
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2014/053516

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	JP 2005-102695 A (National Institute of Advanced Industrial Science and Technology), 21 April 2005 (21.04.2005), paragraphs [0034] to [0042] (Family: none)	17-20, 37-40, 43

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, C12N15/115(2010.01)n		
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, C12N15/115		
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 1971-2014年 日本国実用新案登録公報 1996-2014年 日本国登録実用新案公報 1994-2014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
X	JP 2012-146067 A (日本ソフトウェアマネジメント株式会社) 2012.08.02, 【請求項1】 - 【請求項10】, 【0010】 - 【0172】 & WO	1-11, 21-31, 41, 44
Y	2012/096016 A1 & EP 2665009 A1 & US 2014/019062 A1 & CN 103339632 A	12-20, 32-40, 42, 43
Y	JP 2012-146066 A (日本ソフトウェアマネジメント株式会社) 2012.08.02, 【請求項1】 - 【請求項6】, 【0010】 - 【0165】 & WO	12-20, 32-40, 42, 43
	2012/096015 A1 & EP 2665010 A1 & US 2014/058682 A1 & CN 103348350 A	
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> 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	09.04.2014	国際調査報告の発送日
		22.04.2014
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) 山内 裕史 電話番号 03-3581-1101 内線 3562	5 L 4064

C (続き) . 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
Y	JP 2005-102695 A (独立行政法人産業技術総合研究所) 2005.04.21, 【0034】 - 【0042】 (ファミリーなし)	17-20, 37-40, 43