

(19) 日本国特許庁(JP)

再公表特許(A1)

(11) 国際公開番号

W02004/051546

発行日 平成18年4月6日(2006.4.6)

(43) 国際公開日 平成16年6月17日(2004.6.17)

| | | |
|-----------------------------|-----------------|-------------|
| (51) Int. Cl. | F I | テーマコード (参考) |
| G06F 19/00 (2006.01) | G06F 19/00 600 | 5B075 |
| G06F 17/30 (2006.01) | G06F 17/30 170F | |
| | G06F 17/30 210D | |
| | G06F 17/30 350C | |

審査請求 未請求 予備審査請求 未請求 (全 51 頁)

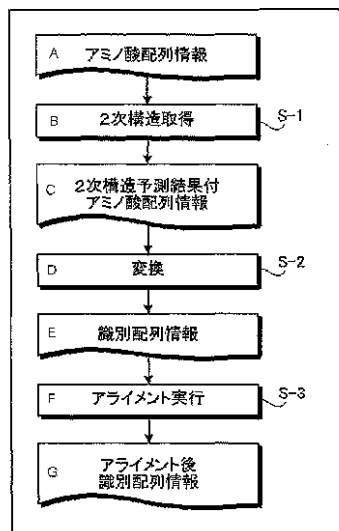
| | |
|---|------------------------------------|
| 出願番号 特願2004-556854 (P2004-556854) | (71) 出願人 502433955 |
| (21) 国際出願番号 PCT/JP2003/015245 | 株式会社インシリコサイエンス |
| (22) 国際出願日 平成15年11月28日(2003.11.28) | 東京都大田区東雪谷二丁目15番9号 |
| (31) 優先権主張番号 特願2002-348678 (P2002-348678) | (74) 代理人 100089118 |
| (32) 優先日 平成14年11月29日(2002.11.29) | 弁理士 酒井 宏明 |
| (33) 優先権主張国 日本国(JP) | (74) 代理人 100113103 |
| | 弁理士 香島 拓也 |
| | (72) 発明者 梅山 秀明 |
| | 千葉県浦安市美浜1丁目7番1002号 |
| | (72) 発明者 岩館 満雄 |
| | 埼玉県羽生市大字尾崎26番地5 |
| | Fターム(参考) 5B075 ND02 NR12 QM08 UU19 |

最終頁に続く

(54) 【発明の名称】 配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体

(57) 【要約】

本発明は、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得し、取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換し、変換された複数の識別配列情報に対してアライメントを実行する。



- A.. AMINO ACID SEQUENCE DATA
- B.. ACQUISITION OF SECONDARY STRUCTURE
- C.. AMINO ACID SEQUENCE DATA WITH SECONDARY STRUCTURE ESTIMATION RESULT
- D.. CONVERSION
- E.. IDENTIFICATION SEQUENCE DATA
- F.. ALIGNMENT PERFORMANCE
- G.. POST-ALIGNMENT IDENTIFICATION SEQUENCE DATA

【特許請求の範囲】**【請求項 1】**

アミノ酸配列情報を取得し、取得された上記アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得手段と、

上記二次構造取得手段にて取得された上記二次構造および上記アミノ酸配列情報を構成する上記アミノ酸情報に対応するアミノ酸の類似性に基づいて、上記アミノ酸配列情報を構成する各アミノ酸情報を、同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の情報として識別するための識別情報に変換することにより、上記アミノ酸配列情報を上記識別情報からなる識別配列情報に変換する変換手段と、

上記変換手段にて変換された複数の上記識別配列情報に対してアライメントを実行するアライメント実行手段と、

を備えたことを特徴とする配列情報処理装置。

10

【請求項 2】

上記二次構造および上記類似性に基づいて同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の上記識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納手段、

を備え、

上記変換手段は、

上記二次構造・類似性分類情報格納手段にて格納された上記二次構造・類似性分類情報に基づいて、上記アミノ酸配列情報を上記識別配列情報に変換する分類情報参照変換手段

20

をさらに備えたことを特徴とする請求の範囲第 1 項に記載の配列情報処理装置。

【請求項 3】

上記分類情報参照変換手段により変換された上記識別配列情報を格納する識別配列情報格納手段、

を備え、

上記変換手段は、

上記識別配列情報格納手段にて格納された上記識別配列情報から、上記アミノ酸配列情報に対応する上記識別配列情報を検索する識別配列情報検索手段、

をさらに備えたことを特徴とする請求の範囲第 1 項または第 2 項に記載の配列情報処理装置。

30

【請求項 4】

上記二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納手段、

および/または、

上記識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納手段、

を備え、

上記アライメント実行手段は、

上記二次構造置換マトリックス格納手段にて格納された上記二次構造置換マトリックスおよび/または上記類似性置換マトリックス格納手段にて格納された上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行手段、

40

をさらに備えたことを特徴とする請求の範囲第 1 項から第 3 項のいずれか一つに記載の配列情報処理装置。

【請求項 5】

上記置換マトリックス基準アライメント実行手段は、予め定めた係数により重み付けされた上記二次構造置換マトリックスおよび/または上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行すること、

を特徴とする請求の範囲第 4 項に記載の配列情報処理装置。

50

【請求項 6】

上記アライメント実行手段にてアライメントが実行された上記識別配列情報を構成する上記識別情報を、対応する上記アミノ酸情報に再変換する再変換手段、

を備えたことを特徴とする請求の範囲第 1 項から第 5 項のいずれか一つに記載の配列情報処理装置。

【請求項 7】

上記二次構造・類似性分類情報は、

上記二次構造を、ヘリックス構造、シート構造、または、その他の構造の 3 つに分類し、上記類似性を、アラニン (Ala) とグリシン (Gly) からなる第 1 の群、アスパラギン酸 (Asp) とグルタミン酸 (Glu) とアスパラギン (Asn) とグルタミン (Gln) からなる第 2 の群、システイン (Cys) からなる第 3 の群、フェニルアラニン (Phe) とヒスチジン (His) とトリプトファン (Trp) とチロシン (Tyr) からなる第 4 の群、イソロイシン (Ile) とロイシン (Leu) とメチオニン (Met) とバリン (Val) からなる第 5 の群、リシン (Lys) とアルギニン (Arg) からなる第 6 の群、プロリン (Pro) とセリン (Ser) とトレオニン (Thr) からなる第 7 の群の 7 つに分類すること、

を特徴とする請求の範囲第 1 項から第 6 項のいずれか一つに記載の配列情報処理装置。

【請求項 8】

上記アミノ酸の上記類似性は、上記アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であること、

を特徴とする請求の範囲第 1 項から第 7 項のいずれか一つに記載の配列情報処理装置。

【請求項 9】

アミノ酸配列情報を取得し、取得された上記アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得ステップと、

上記二次構造取得ステップにて取得された上記二次構造および上記アミノ酸配列情報を構成する上記アミノ酸情報に対応するアミノ酸の類似性に基づいて、上記アミノ酸配列情報を構成する各アミノ酸情報を、同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の情報として識別するための識別情報に変換することにより、上記アミノ酸配列情報を上記識別情報からなる識別配列情報に変換する変換ステップと、

上記変換ステップにて変換された複数の上記識別配列情報に対してアライメントを実行するアライメント実行ステップと、

を含むことを特徴とする配列情報処理方法。

【請求項 10】

上記二次構造および上記類似性に基づいて同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の上記識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納ステップ、

を含み、

上記変換ステップは、

上記二次構造・類似性分類情報格納ステップにて格納された上記二次構造・類似性分類情報に基づいて、上記アミノ酸配列情報を上記識別配列情報に変換する分類情報参照変換ステップ、

をさらに含むことを特徴とする請求の範囲第 9 項に記載の配列情報処理方法。

【請求項 11】

上記分類情報参照変換ステップにより変換された上記識別配列情報を格納する識別配列情報格納ステップ、

を含み、

上記変換ステップは、

上記識別配列情報格納ステップにて格納された上記識別配列情報から、上記アミノ酸配列情報に対応する上記識別配列情報を検索する識別配列情報検索ステップ、

をさらに含むことを特徴とする請求の範囲第 9 項または第 10 項に記載の配列情報処理

10

20

30

40

50

方法。

【請求項 1 2】

上記二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納ステップ、

および/または、

上記識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納ステップ、

を含み、

上記アライメント実行ステップは、

上記二次構造置換マトリックス格納ステップにて格納された上記二次構造置換マトリックスおよび/または上記類似性置換マトリックス格納ステップにて格納された上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行ステップ、

をさらに含むことを特徴とする請求の範囲第 9 項から第 11 項のいずれか一つに記載の配列情報処理方法。

10

【請求項 1 3】

上記置換マトリックス基準アライメント実行ステップは、予め定めた係数により重み付けされた上記二次構造置換マトリックスおよび/または上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行すること、

を特徴とする請求の範囲第 12 項に記載の配列情報処理方法。

20

【請求項 1 4】

上記アライメント実行ステップにてアライメントが実行された上記識別配列情報を構成する上記識別情報を、対応する上記アミノ酸情報に再変換する再変換ステップ、

を含むことを特徴とする請求の範囲第 9 項から第 13 項のいずれか一つに記載の配列情報処理方法。

【請求項 1 5】

上記二次構造・類似性分類情報は、

上記二次構造を、ヘリックス構造、シート構造、または、その他の構造の 3 つに分類し、上記類似性を、アラニン (Ala) とグリシン (Gly) からなる第 1 の群、アスパラギン酸 (Asp) とグルタミン酸 (Glu) とアスパラギン (Asn) とグルタミン (Gln) からなる第 2 の群、システイン (Cys) からなる第 3 の群、フェニルアラニン (Phe) とヒスチジン (His) とトリプトファン (Trp) とチロシン (Tyr) からなる第 4 の群、イソロイシン (Ile) とロイシン (Leu) とメチオニン (Met) とバリン (Val) からなる第 5 の群、リシン (Lys) とアルギニン (Arg) からなる第 6 の群、プロリン (Pro) とセリン (Ser) とトレオニン (Thr) からなる第 7 の群の 7 つに分類すること、

30

を特徴とする請求の範囲第 9 項から第 14 項のいずれか一つに記載の配列情報処理方法。

【請求項 1 6】

上記アミノ酸の上記類似性は、上記アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であること、

40

を特徴とする請求の範囲第 9 項から第 15 項のいずれか一つに記載の配列情報処理方法。

【請求項 1 7】

アミノ酸配列情報を取得し、取得された上記アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得ステップと、

上記二次構造取得ステップにて取得された上記二次構造および上記アミノ酸配列情報を構成する上記アミノ酸情報に対応するアミノ酸の類似性に基づいて、上記アミノ酸配列情報を構成する各アミノ酸情報を、同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の情報として識別するための識別情報に変換することにより、上記

50

アミノ酸配列情報を上記識別情報からなる識別配列情報に変換する変換ステップと、

上記変換ステップにて変換された複数の上記識別配列情報に対してアライメントを実行するアライメント実行ステップと、

を含む配列情報処理方法をコンピュータに実行させることを特徴とするプログラム。

【請求項 18】

上記二次構造および上記類似性に基づいて同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の上記識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納ステップ、

を含み、

上記変換ステップは、

上記二次構造・類似性分類情報格納ステップにて格納された上記二次構造・類似性分類情報に基づいて、上記アミノ酸配列情報を上記識別配列情報に変換する分類情報参照変換ステップ、

をさらに含むことを特徴とする請求の範囲第 17 項に記載のプログラム。

10

【請求項 19】

上記分類情報参照変換ステップにより変換された上記識別配列情報を格納する識別配列情報格納ステップ、

を含み、

上記変換ステップは、

上記識別配列情報格納ステップにて格納された上記識別配列情報から、上記アミノ酸配列情報に対応する上記識別配列情報を検索する識別配列情報検索ステップ、

をさらに含むことを特徴とする請求の範囲第 17 項または第 18 項に記載のプログラム

20

【請求項 20】

上記二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納ステップ、

および/または、

上記識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納ステップ、

を含み、

上記アライメント実行ステップは、

上記二次構造置換マトリックス格納ステップにて格納された上記二次構造置換マトリックスおよび/または上記類似性置換マトリックス格納ステップにて格納された上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行ステップ、

をさらに含むことを特徴とする請求の範囲第 17 項から第 19 項のいずれか一つに記載のプログラム。

30

【請求項 21】

上記置換マトリックス基準アライメント実行ステップは、予め定めた係数により重み付けされた上記二次構造置換マトリックスおよび/または上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行すること、

を特徴とする請求の範囲第 20 項に記載のプログラム。

40

【請求項 22】

上記アライメント実行ステップにてアライメントが実行された上記識別配列情報を構成する上記識別情報を、対応する上記アミノ酸情報に再変換する再変換ステップ、

を含むことを特徴とする請求の範囲第 17 項から第 21 項のいずれか一つに記載のプログラム。

【請求項 23】

上記二次構造・類似性分類情報は、

上記二次構造を、ヘリックス構造、シート構造、または、その他の構造の 3 つに分

50

類し、上記類似性を、アラニン (Ala) とグリシン (Gly) からなる第 1 の群、アスパラギン酸 (Asp) とグルタミン酸 (Glu) とアスパラギン (Asn) とグルタミン (Gln) からなる第 2 の群、システイン (Cys) からなる第 3 の群、フェニルアラニン (Phe) とヒスチジン (His) とトリプトファン (Trp) とチロシン (Tyr) からなる第 4 の群、イソロイシン (Ile) とロイシン (Leu) とメチオニン (Met) とバリン (Val) からなる第 5 の群、リシン (Lys) とアルギニン (Arg) からなる第 6 の群、プロリン (Pro) とセリン (Ser) とトレオニン (Thr) からなる第 7 の群の 7 つに分類すること、

を特徴とする請求の範囲第 17 項から第 22 項のいずれか一つに記載のプログラム。

【請求項 24】

10

上記アミノ酸の上記類似性は、上記アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であること、

を特徴とする請求の範囲第 17 項から第 23 項のいずれか一つに記載のプログラム。

【請求項 25】

上記請求の範囲第 17 項から第 24 項のいずれか一つに記載されたプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

本発明は、配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体に関するものであり、特に、アミノ酸配列の二次構造および/またはアミノ酸の類似性を加味したアライメントを得ることができる、配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体に関する。

20

【背景技術】

現在、タンパク質の二次構造をそのアミノ酸配列から予測できるようになっている。その代表的な二次構造予測ソフトである PSIPRED (protein structure prediction) は高速で精度が高く、PSI-BLAST で作成されるマルチプルアライメントを用いて進化の過程で機能を維持するために保存された構造や、残基置換の制約といったことも考慮して二次構造を予測することができる。

また、ホモロジーモデリング法とは、立体構造既知のタンパク質に関する情報を利用して、目的タンパク質とのアライメントから立体構造を作成する方法である。このホモロジーモデリングを用いて作成されるモデルの精度は近年目覚しく向上しているが、未だ解決すべき問題点も多い。

30

まず、ホモロジーモデリングにおいて、精度の高いモデルを作成するためには、まず信頼性の高いアライメントを得ることが重要である。任意の配列が与えられたとき、PDB (Protein Data Bank) や SCOP (Structure Classification of Protein) のような立体構造データベースから配列類似性の高いタンパク質を単数、または複数選び出し、アライメントを与えるホモロジー検索方法として、FASTA、PSI-BLAST (Position-Specific Iterated BLAST)、RPS-BLAST (Reverse PSI-BLAST)、IMPALA (Integrating Matrix Profiles And Local Alignments) などのアライメントソフトがある。

40

ここで、ホモロジー検索を行なうと、一般的に各アライメントについてホモロジーと e 値を得ることができる。ここで、「ホモロジー」とは残基一致度 (%) のことである。また、「e 値 (Expected Value)」とは、データベースにおいて全く偶然に同じスコアになる配列の数の期待値、すなわちそのアライメントのスコアがどの程度まれであるのかを示す指標であり、小さければ小さいほど似た配列は他に見つかりにくく、偶然には見つかりにくいことを表わす。

また、FASTA は 20 種の天然アミノ酸のアルファベット文字配列マッチングを行なうプログラムであり、高ホモロジー (アミノ酸配列の一致度 30% 以上、e 値 0.01 以下に相当) の参照タンパク質とのアライメントを用いてモデルを作成すると信頼性の高い

50

モデルができるといわれている。一方 P S I - B L A S T では、F A S T A と同じように文字列のマッチングを行なうが、文字が一致しているかどうかの情報ではなく、プロファイルと呼ばれる文字の一致の度合い、置換のしやすさなどを文字配列上の部位ごとに置換配列として算出し、さらに何回も繰り返して算出することによりアライメントを最適化する性質をもっている。なお、P S I - B L A S T ではホモロジーが低い場合であっても e 値が低ければある程度信頼性のあるアライメントを得ることができるといわれている。

現在、タンパク質の二次構造をそのアミノ酸配列から予測できるようになっている。その代表的な二次構造予測ソフトである P S I P R E D (p r o t e i n s t r u c t u r e p r e d i c t i o n) は高速で精度が高く、P S I - B L A S T で作成されるマルチプルアライメントを用いて進化の過程で機能を維持するために保存された構造や、残基の置換の制約といったことも考慮して二次構造を予測することができる。

10

さらに、ホモロジーモデリング法を用いた F A M S (F u l l A u t o m a t i c M o d e l i n g S y s t e m) は、X 線や N M R などにより実験的に構造が決定されたタンパク質の立体構造 (参照タンパク質) をもとにして、立体構造未知のタンパク質 (目的タンパク質) の構造をモデリングすることが可能である。F A M S ではホモロジー検索プログラムなどにより作成されるアライメント (目的タンパク質と参照タンパク質のアミノ酸配列を並置したもの) ファイルを入力としている。

しかしながら、F A S T A、P S I - B L A S T などの従来のホモロジー検索プログラムにおいては、以下の 2 つの問題点が指摘されている。

20

(1) 低ホモロジーで高 e 値の場合、検索が困難であること

P S I - B L A S T のように位置ごとに置換配列を作成 (プロファイル) し、ホモロジー検索を行なうプログラムは P S S M (P o s i t i o n S p e c i f i c S c o r e M a t r i x) と呼ばれ、F A S T A のように文字列の一致度のみを考慮してアライメントを行なう従来のプログラムよりも、より感受性が高く遠縁の配列も検索できるといわれている。しかしながら、低ホモロジーで高 e 値の場合には P S S M によるアライメントであっても信頼性は低下する、または検索が不可能という場合があり限界がある。ここで低ホモロジーの場合、より広く遠縁の配列を探索するホモロジー検索手法の開発が望まれていた。

(2) アライメントの信頼性

従来のアライメントプログラムでは、アミノ酸配列の文字情報のみを用いてアライメントを行なっていた。このアライメントでは、タンパク質の立体構造形成の単位であるヘリックスやシートなどの二次構造を分断してアライメントを行なってしまふことがある。また、このアライメントを用いてモデリングを行なった場合、疎水性アミノ酸が溶媒に露出してしまい、疎水性相互作用の減少から実際の形と比べて大きく異なったモデルを作成してしまふことがある。このように文字のみを考慮したアライメントでは、特に低ホモロジーの場合その信頼性が疑わしいことがある。この問題を解消するため、二次構造情報や残基の疎水性情報をも考慮したアライメントプログラムの開発が望まれていた。

30

また、上記で述べたホモロジー値、e 値はホモロジー検索プログラムが出力するものであり、アミノ酸配列の文字情報のみから算出したものであって、実際に作成されたモデル構造についての評価は行なうことができない。そのため、ホモロジー値、e 値だけではなく作成したモデルの構造をも同時に考慮したモデルの選定方法の開発が望まれていた。

40

さらに、F A M S では、入力ファイルであるアライメント次第でモデリングできる領域、長さが決まってしまう。しかし、異なる参照タンパク質によるアライメントを用いると、より長くモデルを作成できる場合がある。このような場合、よい構造であるが領域が短いモデルに、長い領域をモデリングできたモデルの末端構造を貼り付けることによって、自動的にモデリング領域を伸長させる技術が望まれていた。

本発明は上記問題点に鑑みてなされたもので、上記問題点を解決することのできる、各種の方法等を提供することを目的としている。

【発明の開示】

上述した目的を達成するために、本発明にかかる配列情報処理装置は、アミノ酸配列情

50

報を取得し、取得された上記アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得手段と、上記二次構造取得手段にて取得された上記二次構造および上記アミノ酸配列情報を構成する上記アミノ酸情報に対応するアミノ酸の類似性に基づいて、上記アミノ酸配列情報を構成する各アミノ酸情報を、同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の情報として識別するための識別情報に変換することにより、上記アミノ酸配列情報を上記識別情報からなる識別配列情報に変換する変換手段と、上記変換手段にて変換された複数の上記識別配列情報に対してアライメントを実行するアライメント実行手段とを備えたことを特徴とする。

この装置によれば、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得（例えば、アミノ酸配列情報に対応する立体構造が未知の場合には、既存の二次構造予測プログラムなどを用いて予測して取得、また、対応する立体構造が既知の場合には既存の二次構造判定プログラムなどを用いて取得）し、取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性（例えば、各アミノ酸の疎水性情報等のアミノ酸毎の性質に関する類似性等）に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換し、変換された複数の識別配列情報に対してアライメントを実行するので、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効果的に実行することができ、従来のホモロジー検索法と比べ、より遠縁の配列の探索が可能となる。

すなわち、本装置によれば、二次構造情報や疎水性情報なども考慮した探索・アライメントを行なうため、従来のアミノ酸の文字のみを考慮したアライメントと比較すると、立体構造も考慮したアライメントが可能となり、より高精度なアライメント作成が可能となる。

本装置のように二次構造情報や疎水性情報などを考慮して作成したアライメントを用いると、疎水エネルギーが安定化し、また二次構造が分断していないモデルを作成でき、モデルとしても真実構造に近いものができると考えられる。

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記二次構造および上記類似性に基づいて同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の上記識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納手段を備え、上記変換手段は、上記二次構造・類似性分類情報格納手段にて格納された上記二次構造・類似性分類情報に基づいて、上記アミノ酸配列情報を上記識別配列情報に変換する分類情報参照変換手段をさらに備えたことを特徴とする。

これは二次構造・類似性分類情報格納手段および変換手段の一例を一層具体的に示すものである。この装置によれば、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を格納し、変換手段は、二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換するので、二次構造・類似性分類情報を予め作成して格納しておき、変換時に当該情報を参照することにより、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効率よく得ることができる。

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記分類情報参照変換手段により変換された上記識別配列情報を格納する識別配列情報格納手段を備え、上記変換手段は、上記識別配列情報格納手段にて格納された上記識別配列情報から、上記アミノ酸配列情報に対応する上記識別配列情報を検索する識別配列情報検索手段をさらに備えたことを特徴とする。

これは識別配列情報格納手段および変換手段の一例を一層具体的に示すものである。この装置によれば、変換された識別配列情報を格納し、変換手段は、格納された識別配列情

報から、アミノ酸配列情報に対応する識別配列情報を検索するので、変換された識別配列情報をデータベースなどに予め格納しておき、変換時に当該データベースなどを参照することにより、変換処理を効率化、高速化することができるようになる。

また、これにより、例えば公共のデータベース（例えばPDBなど）に登録されたアミノ酸配列情報の二次構造を既存の二次構造判定プログラムなどを用いて判定した後、当該アミノ酸配列情報を本装置により二次構造・類似性分類情報に基づいて変換して作成した識別配列情報を予めデータベースに格納して利用することができるようになり、取得したアミノ酸配列情報に対応する識別配列情報を効率よく検索することができる。

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納手段、および/または、上記識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納手段を備え、上記アライメント実行手段は、上記二次構造置換マトリックス格納手段にて格納された上記二次構造置換マトリックスおよび/または上記類似性置換マトリックス格納手段にて格納された上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行手段をさらに備えたことを特徴とする。

10

これは二次構造置換マトリックス格納手段、類似性置換マトリックス格納手段およびアライメント実行手段の一例を一層具体的に示すものである。この装置によれば、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックス、および/または、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納し、アライメント実行手段は、格納された二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、置換マトリックスを予め作成して格納しておき、アライメント実行時に当該置換マトリックスを参照して各スコア値に置換することにより、アミノ酸配列情報の二次構造および/またはアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したスコア値に置換することにより最適なアライメントを効率よく得ることができる。

20

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記置換マトリックス基準アライメント実行手段は、予め定めた係数により重み付けされた上記二次構造置換マトリックスおよび/または上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行することを特徴とする。

30

これは置換マトリックス基準アライメント実行手段の一例を一層具体的に示すものである。この装置によれば、置換マトリックス基準アライメント実行手段は、予め定めた係数により重み付けされた二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、例えば生物学的な知見などに基づいて二次構造置換マトリックスおよび/または類似性置換マトリックスに対して適切な係数を設定して重み付けをすることによって、生物学的な知見などを反映した最適なアライメントを効率よく得ることができる。

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記アライメント実行手段にてアライメントが実行された上記識別配列情報を構成する上記識別情報を、対応する上記アミノ酸情報に再変換する再変換手段を備えたことを特徴とする。

40

この装置によれば、アライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換するので、二次構造および類似性を加味してアライメントが実行されたアミノ酸配列情報を効率よく得ることができる。

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記二次構造・類似性分類情報は、上記二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、上記類似性を、アラニン（Ala）とグリシン（Gly）からなる第1の群、アスパラギン酸（Asp）とグルタミン酸（Glu）とアスパラギン（Asn）とグルタミン（Gln）からなる第2の群、システイン（Cys）から

50

なる第3の群、フェニルアラニン (Phe) とヒスチジン (His) とトリプトファン (Trp) とチロシン (Tyr) からなる第4の群、イソロイシン (Ile) とロイシン (Leu) とメチオニン (Met) とバリン (Val) からなる第5の群、リシン (Lys) とアルギニン (Arg) からなる第6の群、プロリン (Pro) とセリン (Ser) とトレオニン (Thr) からなる第7の群の7つに分類することを特徴とする。

これは二次構造・類似性分類情報の一例を一層具体的に示すものである。この装置によれば、二次構造・類似性分類情報は、二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、類似性を、アラニン (Ala) とグリシン (Gly) からなる第1の群、アスパラギン酸 (Asp) とグルタミン酸 (Glu) とアスパラギン (Asn) とグルタミン (Gln) からなる第2の群、システイン (Cys) からなる第3の群、フェニルアラニン (Phe) とヒスチジン (His) とトリプトファン (Trp) とチロシン (Tyr) からなる第4の群、イソロイシン (Ile) とロイシン (Leu) とメチオニン (Met) とバリン (Val) からなる第5の群、リシン (Lys) とアルギニン (Arg) からなる第6の群、プロリン (Pro) とセリン (Ser) とトレオニン (Thr) からなる第7の群の7つに分類するので、既知の置換マトリックスである BLOSUM62 に基づいてアミノ酸を7つの群に分類することによりアミノ酸情報の類似性を加味し、かつ、二次構造を3つに分類することによりアミノ酸配列情報の二次構造を加味して、アライメントが実行された識別配列情報を効率よく得ることができる。

つぎの発明にかかる配列情報処理装置は、上記に記載の配列情報処理装置において、上記アミノ酸の上記類似性は、上記アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であることを特徴とする。

これはアミノ酸の類似性の一例を一層具体的に示すものである。この装置によれば、アミノ酸の類似性は、アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であるので、アミノ酸配列情報の立体構造などに影響するアミノ酸の性質 (疎水性、親水性、酸性、塩基性、荷電状態など) の類似性を考慮してアライメントが実行された識別配列情報を効果的に得ることができる。

また、本発明は配列情報処理方法に関するものであり、本発明にかかる配列情報処理方法は、アミノ酸配列情報を取得し、取得された上記アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得ステップと、上記二次構造取得ステップにて取得された上記二次構造および上記アミノ酸配列情報を構成する上記アミノ酸情報に対応するアミノ酸の類似性に基づいて、上記アミノ酸配列情報を構成する各アミノ酸情報を、同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の情報として識別するための識別情報に変換することにより、上記アミノ酸配列情報を上記識別情報からなる識別配列情報に変換する変換ステップと、上記変換ステップにて変換された複数の上記識別配列情報に対してアライメントを実行するアライメント実行ステップとを含むことを特徴とする。

この方法によれば、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得 (例えば、アミノ酸配列情報に対応する立体構造が未知の場合には、既存の二次構造予測プログラムなどを用いて予測して取得、また、対応する立体構造が既知の場合には既存の二次構造判定プログラムなどを用いて取得) し、取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性 (例えば、各アミノ酸の疎水性情報等のアミノ酸毎の性質に関する類似性等) に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換し、変換された複数の識別配列情報に対してアライメントを実行するので、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効果的に実行することができ、従来のホモロジー検索法と比べ、より遠縁の配列の探索が可能となる。

すなわち、本方法によれば、二次構造情報や疎水性情報なども考慮した探索・アライメ

10

20

30

40

50

ントを行なうため、従来のアミノ酸の文字のみを考慮したアライメントと比較すると、立体構造も考慮したアライメントが可能となり、より高精度なアライメント作成が可能となる。

本方法のように二次構造情報や疎水性情報などを考慮して作成したアライメントを用いると、疎水エネルギーが安定化し、また二次構造が分断していないモデルを作成でき、モデルとしても真実構造に近いものができると考えられる。

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記二次構造および上記類似性に基づいて同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の上記識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納ステップを含み、上記変換ステップは、上記二次構造・類似性分類情報格納ステップにて格納された上記二次構造・類似性分類情報に基づいて、上記アミノ酸配列情報を上記識別配列情報に変換する分類情報参照変換ステップをさらに含むことを特徴とする。

10

これは二次構造・類似性分類情報格納ステップおよび変換ステップの一例を一層具体的に示すものである。この方法によれば、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を格納し、変換ステップは、二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換するので、二次構造・類似性分類情報を予め作成して格納しておき、変換時に当該情報を参照することにより、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効率よく得ることができる。

20

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記分類情報参照変換ステップにより変換された上記識別配列情報を格納する識別配列情報格納ステップを含み、上記変換ステップは、上記識別配列情報格納ステップにて格納された上記識別配列情報から、上記アミノ酸配列情報に対応する上記識別配列情報を検索する識別配列情報検索ステップをさらに含むことを特徴とする。

これは識別配列情報格納ステップおよび変換ステップの一例を一層具体的に示すものである。この方法によれば、変換された識別配列情報を格納し、変換ステップは、格納された識別配列情報から、アミノ酸配列情報に対応する識別配列情報を検索するので、変換された識別配列情報をデータベースなどに予め格納しておき、変換時に当該データベースなどを参照することにより、変換処理を効率化、高速化することができるようになる。

30

また、これにより、例えば公共のデータベース（例えばPDBなど）に登録されたアミノ酸配列情報の二次構造を既存の二次構造判定プログラムなどを用いて判定した後、当該アミノ酸配列情報を本方法により二次構造・類似性分類情報に基づいて変換して作成した識別配列情報を予めデータベースに格納して利用することができるようになり、取得したアミノ酸配列情報に対応する識別配列情報を効率よく検索することができる。

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納ステップ、および/または、上記識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納ステップを含み、上記アライメント実行ステップは、上記二次構造置換マトリックス格納ステップにて格納された上記二次構造置換マトリックスおよび/または上記類似性置換マトリックス格納ステップにて格納された上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行ステップをさらに含むことを特徴とする。

40

これは二次構造置換マトリックス格納ステップ、類似性置換マトリックス格納ステップおよびアライメント実行ステップの一例を一層具体的に示すものである。この方法によれば、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックス、および/または、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納し、アライメント実行ステップは、格納された二次構造置

50

換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、置換マトリックスを予め作成して格納しておき、アライメント実行時に当該置換マトリックスを参照して各スコア値に置換することにより、アミノ酸配列情報の二次構造および/またはアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したスコア値に置換することにより最適なアライメントを効率よく得ることができる。

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記置換マトリックス基準アライメント実行ステップは、予め定めた係数により重み付けされた上記二次構造置換マトリックスおよび/または上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行することを特徴とする。

10

これは置換マトリックス基準アライメント実行ステップの一例を一層具体的に示すものである。この方法によれば、置換マトリックス基準アライメント実行ステップは、予め定めた係数により重み付けされた二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、例えば生物学的な知見などに基づいて二次構造置換マトリックスおよび/または類似性置換マトリックスに対して適切な係数を設定して重み付けをすることによって、生物学的な知見などを反映した最適なアライメントを効率よく得ることができる。

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記アライメント実行ステップにてアライメントが実行された上記識別配列情報を構成する上記識別情報を、対応する上記アミノ酸情報に再変換する再変換ステップを含むことを特徴とする。

20

この方法によれば、アライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換するので、二次構造および類似性を加味してアライメントが実行されたアミノ酸配列情報を効率よく得ることができる。

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記二次構造・類似性分類情報は、上記二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、上記類似性を、アラニン(Ala)とグリシン(Gly)からなる第1の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第2の群、システイン(Cys)からなる第3の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第4の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第5の群、リシン(Lys)とアルギニン(Arg)からなる第6の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第7の群の7つに分類することを特徴とする。

30

これは二次構造・類似性分類情報の一例を一層具体的に示すものである。この方法によれば、二次構造・類似性分類情報は、二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、類似性を、アラニン(Ala)とグリシン(Gly)からなる第1の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第2の群、システイン(Cys)からなる第3の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第4の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第5の群、リシン(Lys)とアルギニン(Arg)からなる第6の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第7の群の7つに分類するので、既知の置換マトリックスであるBLOSUM62に基づいてアミノ酸を7つの群に分類することによりアミノ酸情報の類似性を加味し、かつ、二次構造を3つに分類することによりアミノ酸配列情報の二次構造を加味して、アライメントが実行された識別配列情報を効率よく得ることができる。

40

つぎの発明にかかる配列情報処理方法は、上記に記載の配列情報処理方法において、上記アミノ酸の上記類似性は、上記アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であることを特徴とする。

50

これはアミノ酸の類似性の一例を一層具体的に示すものである。この方法によれば、アミノ酸の類似性は、アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であるので、アミノ酸配列情報の立体構造などに影響するアミノ酸の性質（疎水性、親水性、酸性、塩基性、荷電状態など）の類似性を考慮してアライメントが実行された識別配列情報を効果的に得ることができる。

また、本発明はプログラムに関するものであり、本発明にかかる配列情報処理方法をコンピュータに実行させるためのプログラムは、アミノ酸配列情報を取得し、取得された上記アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得ステップと、上記二次構造取得ステップにて取得された上記二次構造および上記アミノ酸配列情報を構成する上記アミノ酸情報に対応するアミノ酸の類似性に基づいて、上記アミノ酸配列情報を構成する各アミノ酸情報を、同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の情報として識別するための識別情報に変換することにより、上記アミノ酸配列情報を上記識別情報からなる識別配列情報に変換する変換ステップと、上記変換ステップにて変換された複数の上記識別配列情報に対してアライメントを実行するアライメント実行ステップとを含むことを特徴とする。

このプログラムによれば、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得（例えば、アミノ酸配列情報に対応する立体構造が未知の場合には、既存の二次構造予測プログラムなどを用いて予測して取得、また、対応する立体構造が既知の場合には既存の二次構造判定プログラムなどを用いて取得）し、取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性（例えば、各アミノ酸の疎水性情報等のアミノ酸毎の性質に関する類似性等）に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換し、変換された複数の識別配列情報に対してアライメントを実行するので、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効果的に実行することができ、従来のホモロジー検索法と比べ、より遠縁の配列の探索が可能となる。

すなわち、本プログラムによれば、二次構造情報や疎水性情報なども考慮した探索・アライメントを行なうため、従来のアミノ酸の文字のみを考慮したアライメントと比較すると、立体構造も考慮したアライメントが可能となり、より高精度なアライメント作成が可能となる。

本プログラムのように二次構造情報や疎水性情報などを考慮して作成したアライメントを用いると、疎水エネルギーが安定化し、また二次構造が分断していないモデルを作成でき、モデルとしても真実構造に近いものができると考えられる。

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記二次構造および上記類似性に基づいて同一の上記二次構造および同一の上記類似性を有する上記アミノ酸情報を同一の上記識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納ステップを含み、上記変換ステップは、上記二次構造・類似性分類情報格納ステップにて格納された上記二次構造・類似性分類情報に基づいて、上記アミノ酸配列情報を上記識別配列情報に変換する分類情報参照変換ステップをさらに含むことを特徴とする。

これは二次構造・類似性分類情報格納ステップおよび変換ステップの一例を一層具体的に示すものである。このプログラムによれば、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を格納し、変換ステップは、二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換するので、二次構造・類似性分類情報を予め作成して格納しておき、変換時に当該情報を参照することにより、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効率よく得ることができる。

10

20

30

40

50

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記分類情報参照変換ステップにより変換された上記識別配列情報を格納する識別配列情報格納ステップを含み、上記変換ステップは、上記識別配列情報格納ステップにて格納された上記識別配列情報から、上記アミノ酸配列情報に対応する上記識別配列情報を検索する識別配列情報検索ステップをさらに含むことを特徴とする。

これは識別配列情報格納ステップおよび変換ステップの一例を一層具体的に示すものである。このプログラムによれば、変換された識別配列情報を格納し、変換ステップは、格納された識別配列情報から、アミノ酸配列情報に対応する識別配列情報を検索するので、変換された識別配列情報をデータベースなどに予め格納しておき、変換時に当該データベースなどを参照することにより、変換処理を効率化、高速化することができるようになる

10

また、これにより、例えば公共のデータベース（例えばPDBなど）に登録されたアミノ酸配列情報の二次構造を既存の二次構造判定プログラムなどを用いて判定した後、当該アミノ酸配列情報を本プログラムにより二次構造・類似性分類情報に基づいて変換して作成した識別配列情報を予めデータベースに格納して利用することができるようになり、取得したアミノ酸配列情報に対応する識別配列情報を効率よく検索することができる。

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納ステップ、および/または、上記識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納ステップを含み、上記アライメント実行ステップは、上記二次構造置換マトリックス格納ステップにて格納された上記二次構造置換マトリックスおよび/または上記類似性置換マトリックス格納ステップにて格納された上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行ステップをさらに含むことを特徴とする。

20

これは二次構造置換マトリックス格納ステップ、類似性置換マトリックス格納ステップおよびアライメント実行ステップの一例を一層具体的に示すものである。このプログラムによれば、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックス、および/または、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納し、アライメント実行ステップは、格納された二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、置換マトリックスを予め作成して格納しておき、アライメント実行時に当該置換マトリックスを参照して各スコア値に置換することにより、アミノ酸配列情報の二次構造および/またはアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したスコア値に置換することにより最適なアライメントを効率よく得ることができる。

30

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記置換マトリックス基準アライメント実行ステップは、予め定めた係数により重み付けされた上記二次構造置換マトリックスおよび/または上記類似性置換マトリックスに基づいて、上記識別配列情報を置換してアライメントを実行することを特徴とする。

40

これは置換マトリックス基準アライメント実行ステップの一例を一層具体的に示すものである。このプログラムによれば、置換マトリックス基準アライメント実行ステップは、予め定めた係数により重み付けされた二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、例えば生物学的な知見などに基づいて二次構造置換マトリックスおよび/または類似性置換マトリックスに対して適切な係数を設定して重み付けをすることによって、生物学的な知見などを反映した最適なアライメントを効率よく得ることができる。

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記アライメント実行ステップにてアライメントが実行された上記識別配列情報を構成する上記識別情報を、対応する上記アミノ酸情報に再変換する再変換ステップを含むことを特徴とする。

50

このプログラムによれば、アライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換するので、二次構造および類似性を加味してアライメントが実行されたアミノ酸配列情報を効率よく得ることができる。

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記二次構造・類似性分類情報は、上記二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、上記類似性を、アラニン(Ala)とグリシン(Gly)からなる第1の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第2の群、システイン(Cys)からなる第3の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第4の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第5の群、リシン(Lys)とアルギニン(Arg)からなる第6の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第7の群の7つに分類することを特徴とする。

これは二次構造・類似性分類情報の一例を一層具体的に示すものである。このプログラムによれば、二次構造・類似性分類情報は、二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、類似性を、アラニン(Ala)とグリシン(Gly)からなる第1の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第2の群、システイン(Cys)からなる第3の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第4の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第5の群、リシン(Lys)とアルギニン(Arg)からなる第6の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第7の群の7つに分類するので、既知の置換マトリックスであるBLOSUM62に基づいてアミノ酸を7つの群に分類することによりアミノ酸情報の類似性を加味し、かつ、二次構造を3つに分類することによりアミノ酸配列情報の二次構造を加味して、アライメントが実行された識別配列情報を効率よく得ることができる。

つぎの発明にかかるプログラムは、上記に記載のプログラムにおいて、上記アミノ酸の上記類似性は、上記アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であることを特徴とする。

これはアミノ酸の類似性の一例を一層具体的に示すものである。このプログラムによれば、アミノ酸の類似性は、アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であるので、アミノ酸配列情報の立体構造などに影響するアミノ酸の性質(疎水性、親水性、酸性、塩基性、荷電状態など)の類似性を考慮してアライメントが実行された識別配列情報を効果的に得ることができる。

また、本発明は記録媒体に関するものであり、本発明にかかるコンピュータ読み取り可能な記録媒体は、上記に記載されたプログラムを記録したことを特徴とする。

この記録媒体によれば、上記に記載されたプログラムをコンピュータを利用して実現することができる。これら各方法と同様の効果を得ることができる。

(関数を用いたモデルの選定)

また、本発明は目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法に関するものであり、本発明にかかる方法は、同一の遺伝子のタンパク質について複数の立体構造が算出された場合、アライメントの適合の度合い(e値等)ばかりでなくコンピュータで作成されたモデル構造のタンパク質らしさを表す構造の指標やホモロジーモデリング対象のタンパク質のアミノ酸シーケンスのみから予測する立体構造の指標やモデルを作成するとき参考にする実験で解明された立体構造の指標からなるある目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

また、つぎの発明にかかる方法は、ある目的関数が、アミノ酸シーケンスしかわかっていないタンパク質のアミノ酸残基とそれと文字並べて類似性ありと判定された立体構造が既知の参照タンパク質のアミノ酸残基のアライメントの適合の度合い(e値等)、目的タ

10

20

30

40

50

ンパク質の予測二次構造と参照タンパク質の判定二次構造の適合率、モデルタンパク質の立体構造に発生する疎水エネルギー等の評価基準を組み合わせたものである上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

また、つぎの発明にかかる方法は、二次構造の適合率をクエリー配列（目的配列）から二次構造予測ソフトウェア（例えば P S I - P R E D 等）を用いて予測した結果とモデル構造の二次構造判定ソフトウェア（例えば D S S P、S T R I D E 等）を用いて決定した結果とを見比べてモデルが構築された領域の残基数に対して二次構造が適合した残基数の割合で表した場合の上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

また、つぎの発明にかかる方法は、ホモロジーモデリングされたタンパク質の立体構造に対する疎水エネルギーの指標の評価は疎水性残基の側鎖同士が下記の表 1 の疎水性原子から構成される疎水性残基同士の距離がある閾値（例えば 6 . 6 ）以下であったとき、疎水性相互作用していると見なし溶媒から遮蔽されており接触している状態よりもエネルギー的に安定であると考えて負のエネルギーの指標をカウントしそのエネルギーの指標をモデル立体構造の座標全体における合計で表した時の上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

(表 1)

| 残基名 | 疎水性原子 | | | |
|-----|-------|-----|-----|-----|
| ALA | CB | | | |
| VAL | CG1 | CG2 | | |
| PHE | CB | CD1 | CD2 | CZ |
| PRO | CG | | | |
| MET | CG | CE | | |
| ILE | CG1 | CG2 | CD1 | |
| LEU | CB | CD1 | CD2 | |
| ASP | | | | |
| GLU | CB | | | |
| LYS | CB | CD | | |
| ARG | CG | | | |
| SER | | | | |
| THR | CG2 | | | |
| TYR | CB | CE1 | CE2 | |
| HIS | CB | | | |
| CYS | CB | | | |
| ASN | CB | | | |
| GLN | CB | CG | | |
| TRP | CB | CD2 | CZ2 | CZ3 |
| GLY | | | | |

また、つぎの発明にかかる方法は、アライメントの適合率（e 値等）、二次構造の適合率、疎水エネルギーの指標等の評価基準を組み合わせる際に、複数種のアライメントソフト（FASTA、BLAST、PSI-BLAST、RPS-BLAST、IMPALA 等）で、e 値と二次構造適合率、疎水エネルギーの指標のどれを優先させるべきか決定するのに有用な特別な判定関数をもつ上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

また、つぎの発明にかかる方法は、二次構造適合率、疎水エネルギー値の積の符号を逆にした値が高い構造はタンパク質らしい構造を持つとして、得られた多数のモデル構造の評価関数とする上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルと

して算出する方法である。

また、つぎの発明にかかる方法は、複数アライメントそれぞれの e 値の対数の符号を逆にしたものにある定数を加算した値に複数のアライメント中で最も低い e 値の対数の符号を逆にしたものにある定数を加算した値で除した値と請求項 6 の評価関数との積を評価関数とする上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

また、つぎの発明にかかる方法は、e 値の対数が 0 で 100 の値を持ち、-3 で 50 の値を、-10 で 40 の値を持つようにして、その間では線形の関係があるような値に 1 を加算した値をある定数とする上記に記載の目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法である。

(本発明で得られたアライメントからホモロジーモデリングを用いて立体構造を得た後の埋め込み処理)

また、本発明は遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法に関するものであり、本発明にかかる方法は、同一の遺伝子のタンパク質についてある精密なモデル立体構造(以下局所構造 1)とそのモデルのアミノ酸配列上の領域を含む別の領域をモデリングしているモデル立体構造(局所構造 2)が存在する場合において、局所構造 1 の立体構造を保持しつつ局所構造 2 の領域をも包括し且つ局所構造 1 の精密さを保持した立体構造を算出し、遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法である。

また、つぎの発明にかかる方法は、精密な局所構造 1 を局所構造 2 が完全に包括する場合において、立体構造的にその RMSD (Root Mean Square Distance) を最小化するように重ね合わせて、局所構造 1 の末端部分と局所構造 2 における局所構造 1 の末端部分に相当するアミノ酸残基が距離的に(例えば、8 以上程度)離れていたとき、精密な局所構造の末端部分 1 残基を切り短くして再度重ね合せを行う、という操作を前述の距離がある閾値(例えば、8)よりも小さくなるまで行い、局所構造 1 の精密さを最大限に保持し且つ局所構造 2 のモデル領域までも包括するモデル立体構造を算出する上記に記載の遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法である。

また、つぎの発明にかかる方法は、精密な局所構造 1 とその一部を含む別領域の局所構造 2 が存在しているとき、立体構造的にその RMSD を最小化するように重ね合わせて、重ね合わせた部分の根平均自乗距離(RMSD)が例えば 2 以上のとき、局所構造 2 の重ね合わせの領域中、局所構造 1 と共通領域でもある末端部分のある残基数分(例えば、2 残基)切り、短くして再度重ね合せを行う、という操作を前述の RMSD がある閾値(例えば、2)よりも小さくなるまで重ね合せ領域が短くなる操作を行う。RMSD が閾値より小さくならず局所構造 2 の残基数が 4 残基になった場合は閾値をある値だけ(例えば、1)上げて上記を繰り返し行う。その結果、局所構造 1 の精密さを最大限に保持し、局所構造 2 の領域をモデリングしたグローバル立体構造を複数の立体構造中の 1 つとする上記に記載の遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法である。

また、つぎの発明にかかる方法は、精密な局所構造 1 とその一部または全部を含む別領域の局所構造 2 が存在しているとき、また上記に記載の遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法の両方がうまく行かない場合においても、立体構造的にその RMSD を最小化するように重ね合わせる長さを一定にして移動させ、重ね合わせた部分の根平均自乗距離(RMSD)がある閾値以下(例えば、2)以下で且つ重ねあわせ領域をなるべく大きくするように変化させる、という操作を行い、閾値以下の領域が見つからない場合は閾値を上げて操作を繰り返し行う。閾値(例えば、4)以下にしても領域が見つからない場合は、重ね合わせる長さを二残基短くして、再度閾値を例えば 2 に戻し前述の操作を行う。その結果、局所構造 1 の精密さを最大限に保持し、局所構造 2 の領域を包括したモデル立体構造を得て、それを複数の立体構造中の 1 つとする上記に記載の遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法である。

10

20

30

40

50

また、つぎの発明にかかる方法は、ある遺伝子のある領域をモデリングした精密な立体構造モデル1が存在する場合に置いて、二次構造単位に細切れにした立体構造データベースを予め作成しておき、P S I - B L A S T等のアライメントによって、断片の検索を行い、その中で遺伝子のアミノ酸配列をP S I - P R E D等の二次構造予測によく一致する断片のみを取り出し、F A M S等のホモロジーモデリングによって断片をモデリングし、その断片をモデル1に対して重ね合わせて遺伝子が完全長でない場合は完全長になるまで重ね合わせつづける。その結果、最初のモデル1の精密さを損なわずに且つ遺伝子の完全長のモデルを得て、それを複数の立体構造中の1つとする上記に記載の遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法である。

また、本発明はプログラムに関するものであり、本発明にかかるプログラムは、上記に記載された目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法または遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法をコンピュータに実行させるためのプログラムであることを特徴とする。

このプログラムによれば、当該プログラムをコンピュータに読み取らせて実行することによって、上記に記載された目的関数に従ってもっとも高得点を得た立体構造をモデルとして算出する方法または遺伝子のアミノ酸配列と同じ長さを持つ完全長モデルを得る方法をコンピュータを利用して実現することができ、これら各方法と同様の効果を得ることができる。

また、本発明は記録媒体に関するものであり、本発明にかかる記録媒体は、上記に記載されたプログラムを記録したことを特徴とする。

この記録媒体によれば、当該記録媒体に記録されたプログラムをコンピュータに読み取らせて実行することによって、上記に記載されたプログラムをコンピュータを利用して実現することができ、これら各方法と同様の効果を得ることができる。

【図面の簡単な説明】

第1図は、本発明の基本原理を示す原理構成図であり、第2図は、本発明が適用される本システムの構成の一例を示すブロック図であり、第3図は、本実施形態におけるアミノ酸配列関連情報ファイル106aに格納される情報の一例を示す図であり、第4図は、本実施形態における二次構造・類似性分類情報ファイル106bに格納される情報の一例を示す図であり、第5図は、本実施形態における識別配列情報データベース106cに格納される情報の一例を示す図であり、第6図は、本実施形態における二次構造置換マトリックス106dに格納される情報の一例を示す図であり、第7図は、本実施形態における類似性置換マトリックス106eに格納される情報の一例を示す図であり、第8図は、本実施形態におけるアライメント後識別配列情報ファイル106fに格納される情報の一例を示す図であり、第9図は、本実施形態における再変換アミノ酸配列情報ファイル106gに格納される情報の一例を示す図であり、第10図は、本発明が適用される本システムの変換部102dの構成の一例を示すブロック図であり、第11図は、本発明が適用される本システムのアライメント実行部102gの構成の一例を示すブロック図であり、第12図は、本実施形態における本システムのメイン処理の一例を示すフローチャートであり、第13図は、本発明を有効に適用するための処理を示すフローチャートであり、第14図は、予測結果ファイルの一例を示す図であり、第15図は、疎水性相互作用の指標(enesosui値)の一例を示す図であり、第16図は、SCOPドメインデータベース(95% Non-Redundant: total_seq_dom_all_95(7433ドメイン))中のドメイン構造について、その残基数とenesosui値との関係を調べた結果をまとめた表を示す図であり、第17図は、本発明の入力ファイルの一例を示す図であり、第18図は、本発明の出力ファイルの一例を示す図であり、第19図は、本発明の処理を示すフローチャート(ペアワイズモードの場合)であり、第20図は、本発明の処理を示すフローチャート(ホモロジー検索モードの場合)であり、第21図は、21種(二次構造3パターン×疎水性7グループ)のアミノ酸を分類する表の一例を示す図であり、第22図は、二次構造によるスコアマトリックスの一例を示す図であり、第23図は、アミノ酸類似性によるスコアマトリックスの一例を示す図であり、第24図は、

10

20

30

40

50

高ホモロジーの場合における本発明の入力ファイルの一例を示す図であり、第25図は、PREDFASTAによる1位の参照タンパク質についての出力ファイルを示す図であり、第26図は、低ホモロジーの場合における本発明の入力ファイルの一例を示す図であり、第27図は、PREDFASTAにより1位で検索された参照タンパク質についての出力ファイルを示す図であり、第28図は、タンパク質(T0176:CAFASP8880)のアミノ酸配列についてPSI-BLAST、PREDFASTAにより得られたアライメントを示す図であり、第29図は、PSI-BLASTによるアライメントにPREDFASTAによるアライメントの末端部分を付け加えることにより、長いアライメントを得ることができた場合を示す図であり、第30図は、ペアワイズモードの場合の本発明の入力ファイルの一例を示す図であり、第31図は、本発明の出力ファイルの一例を示す図であり、第32図は、本発明で得られたアライメントからホモロジーモデリングを用いて立体構造を得た後の処理(埋め込み)を示すフローチャートであり、第33図は、つなぎ目の残基の距離を判定して伸長する場合の概念図であり、第34図は、埋め込みの処理によりつなぎ目の残基の距離を判定して伸長する場合の処理の一例を示すフローチャートであり、第35図は、埋め込みの処理により、末端を優先させるフィッティングによる伸長を行う場合の処理の一例を示すフローチャートであり、第36図は、RMSDの閾値、フィッティング残基数の変化を示す図であり、第37図は、埋め込みの処理によりフィッティング領域を移動させてフィッティングさせる場合の処理の一例を示すフローチャートであり、第38図は、RMSDの閾値、フィッティング残基数、フィッティング領域の変化の順序を示す図であり、第39図は、埋め込みの処理による二次構造データベースの作成法を示すフローチャートであり、第40図は、埋め込みの処理による二次構造データベースの作成法を示すフローチャートである。

10

20

【発明を実施するための最良の形態】

以下に、本発明にかかる配列情報処理装置、配列情報処理方法、プログラム、記録媒体等の実施の形態を図面に基づいて詳細に説明する。なお、この実施の形態によりこの発明が限定されるものではない。

(本発明の基本原則)

以下、本発明の概要について説明し、その後、本発明の構成および処理等について詳細に説明する。第1図は、本発明の基本原則を示す原理構成図である。

本発明は、概略的に、以下の基本的特徴を有する。まず、本発明は、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する(ステップS-1)。なお、二次構造の取得には、既存の二次構造予測プログラム(例えば、PSI-PREDなど)などを用いて予測して取得、また、対応する立体構造が既知の場合には既存の二次構造判定プログラム(例えば、DSSP、STRIDEなど)などを用いて取得してもよい。また、ステップS-1において、予めアライメントしたアミノ酸配列情報を取得してもよい。

30

ついで、ステップS-1にて取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性(例えば、各アミノ酸の疎水性情報等のアミノ酸毎の性質に関する類似性等)に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報(例えば、A、B、Cなどの文字など)に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換する(ステップS-2)。

40

ここで、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を予め格納し、ステップS-2において、格納された二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換してもよい。これにより、変換時に当該情報を参照することにより、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効率よく得ることができる。なお、二次構造・類似性分類情報は、例えばアミノ酸間の置換の起こりやすさをスコア値で表したマトリクス形式の既知の情報であるBLOSUM62のスコア値などに基づいてクラスタリング(ク

50

ラスター解析)を行ったことにより得られた複数のグループを示す類似性情報と、複数のグループに分類された二次構造を示す二次構造情報とを相互に関連付けて構成され、類似性情報と二次構造情報との組み合わせ毎に、対応する識別情報が割り当てられた情報でもよい。具体的には、二次構造・類似性分類情報は、例えば、二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、類似性を、アラニン(Ala)とグリシン(Gly)からなる第1の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第2の群、システイン(Cys)からなる第3の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第4の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第5の群、リシン(Lys)とアルギニン(Arg)からなる第6の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第7の群の7つに分類した情報でもよい。

また、二次構造・類似性分類情報に基づいて変換された識別配列情報を予め格納し、ステップS-2において、格納された識別配列情報から、アミノ酸配列情報に対応する識別配列情報を検索してもよい。これにより、変換された識別配列情報をデータベースなどに予め格納しておき、変換時に当該データベースなどを参照することにより、変換処理を効率化、高速化することができるようになる。また、これにより、例えば公共のデータベース(例えばPDBなど)に登録されたアミノ酸配列情報の二次構造を既存の二次構造判定プログラムなどを用いて判定した後、当該アミノ酸配列情報を本装置により二次構造・類似性分類情報に基づいて変換して作成した識別配列情報を予めデータベースに格納して利用することができるようになり、取得したアミノ酸配列情報に対応する識別配列情報を効率よく検索することができる。

ついで、ステップS-2にて変換された複数の識別配列情報に対してアライメントを実行する(ステップS-3)。

ここで、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスおよび/または識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを予め格納し、ステップS-3において、格納された二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行してもよい。これにより、置換マトリックスを予め作成して格納しておき、アライメント実行時に当該置換マトリックスを参照して各スコア値に置換することにより、アミノ酸配列情報の二次構造および/またはアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したスコア値に置換することにより最適なアライメントを効率よく得ることができる。

なお、ステップS-3において、予め定めた係数により重み付けされた二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行してもよい。これにより、例えば生物学的な知見などに基づいて二次構造置換マトリックスおよび/または類似性置換マトリックスに対して適切な係数を設定して重み付けをすることによって、生物学的な知見などを反映した最適なアライメントを効率よく得ることができる。

ここで、ステップS-3にてアライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換してもよい。これにより、二次構造および/またはアミノ酸類似性を加味してアライメントが行われたアミノ酸配列情報を得ることができる。

なお、上述におけるアミノ酸の類似性は、アミノ酸の疎水性のほか、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性を用いてもよい。これにより、アミノ酸配列情報の立体構造などに影響するアミノ酸の性質(疎水性、親水性、酸性、塩基性、荷電状態など)の類似性を考慮してアライメントが実行された識別配列情報を効果的に得ることができる。

(システム構成)

まず、本システムの構成について説明する。第2図は、本発明が適用される本システムの構成の一例を示すブロック図であり、該構成のうち本発明に関係する部分のみを概念的に示している。本システムは、概略的に、第2図に示すように、配列情報処理装置100と、塩基配列情報やアミノ酸配列情報などに関する外部データベース、配列情報に対するホモロジー検索やアミノ酸配列情報に対するアライメントの実行やタンパク質の立体構造（例えば二次構造、三次構造など）の解析などを行うための外部プログラム等を提供する外部システム200とを、ネットワーク300を介して通信可能に接続して構成されている。

第2図において、ネットワーク300は、配列情報処理装置100を相互に接続する機能を有し、例えば、インターネット等である。

第2図において、配列情報処理装置100は、概略的に、配列情報処理装置100の全体を統括的に制御するCPU等の制御部102、通信回線等に接続されるルータ等の通信装置（図示せず）に接続される通信制御インタフェース部104、入力装置112や出力装置114に接続される入出力制御インタフェース部108、および、各種のデータベースやテーブルなどを格納する記憶部106を備えて構成されており、これら各部は任意の通信路を介して通信可能に接続されている。さらに、この配列情報処理装置100は、ルータ等の通信装置および専用線等の有線または無線の通信回線を介して、ネットワーク300に通信可能に接続されている。

記憶部106に格納される各種のデータベースやテーブル（アミノ酸配列関連情報ファイル106a～再変換アミノ酸配列情報ファイル106g）は、固定ディスク装置等のストレージ手段であり、各種処理に用いる各種のプログラムやテーブルやファイルやデータベースやウェブページ用ファイル等を格納する。

これら記憶部106の各構成要素のうち、アミノ酸配列関連情報ファイル106aは、取得されたアミノ酸配列情報と、当該アミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造、類似性および識別情報と、を格納するアミノ酸配列関連情報格納手段である。ここで、アミノ酸配列関連情報ファイル106aに格納される情報について第3図を参照して説明する。

第3図は、本実施形態におけるアミノ酸配列関連情報ファイル106aに格納される情報の一例を示す図である。第3図に示すように、アミノ酸配列関連情報ファイル106aに格納される情報は、アミノ酸配列情報を一意に識別するためのアミノ酸配列識別情報（第3図に示す「A001」）と、アミノ酸配列情報を構成する各アミノ酸情報（第3図に示す「Ala」、「Phe」、「Trp」、・・・、「Lys」）と、各アミノ酸情報に対応する二次構造情報（第3図に示す「1」、「4」、「4」、・・・、「6」）と、アミノ酸の類似性に基づいて分類されたグループ（群）を示す類似性情報（第3図に示す「1」、「4」、「4」、・・・、「6」）と、変換後の識別配列情報を構成する識別情報（第3図に示す「A」、「K」、「M」、・・・、「S」）と、を相互に関連付けて構成されている。

また、二次構造・類似性分類情報ファイル106bは、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納手段である。ここで、二次構造・類似性分類情報ファイル106bに格納される情報について第4図を参照して説明する。

第4図は、本実施形態における二次構造・類似性分類情報ファイル106bに格納される情報の一例を示す図である。第4図に示すように、二次構造・類似性分類情報ファイル106bに格納される情報は、アミノ酸の類似性に基づいて分類されたグループ（群）を示す類似性情報（第4図に示す「1」、「2」、「3」、・・・、「7」）と、類似性情報により分類されるグループに属するアミノ酸情報（第4図は、例えば、類似性情報「1」に属するアミノ酸情報が「Ala」および「Gly」であることを示している。）と、二次構造情報（第4図に示す「1」、「2」および「その他」）とを相互に関連付けて構成され、類似性情報および二次構造情報の組み合わせ毎に、対応する識別情報が

10

20

30

40

50

格納されている。例えば、第4図に示す例では、類似性情報が「1」、二次構造情報が「」の組み合わせの場合には、識別情報が「A」であることを示している。

また、識別配列情報データベース106cは、二次構造・類似性分類情報に基づいてアミノ酸配列情報を構成するアミノ酸情報を識別情報に変換した識別配列情報を格納する識別配列情報格納手段である。ここで、識別配列情報データベース106cに格納される情報について第5図を参照して説明する。

第5図は、本実施形態における識別配列情報データベース106cに格納される情報の一例を示す図である。第5図に示すように、識別配列情報データベース106cに格納される情報は、アミノ酸配列情報を一意に識別するためのアミノ酸配列識別情報(第5図に示す「K__001」と、アミノ酸配列情報を構成する各アミノ酸情報(第5図に示す「Ala」、「Phe」、「Trp」、・・・、「Lys」)と、アミノ酸配列情報を構成する各アミノ酸情報が二次構造・類似性分類情報に基づいて変換された各識別情報(第5図に示す「A」、「K」、「M」、・・・、「S」)と、を相互に関連付けて構成されている。

10

また、二次構造置換マトリックス106dは、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納手段である。ここで、二次構造置換マトリックス106dに格納される情報について第6図を参照して説明する。

第6図は、本実施形態における二次構造置換マトリックス106dに格納される情報の一例を示す図である。第6図に示すように、二次構造置換マトリックス106dに格納される情報は、二次構造情報(第6図に示す「」、「」および「その他」)の組み合わせに対応する構造スコア値(第6図は、例えば、二次構造情報が「」と「」の組み合わせの場合、構造スコア値が「-6」であることを示している。)で構成されている。

20

また、類似性置換マトリックス106eは、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納手段である。ここで、類似性置換マトリックス106eに格納される情報について第7図を参照して説明する。

第7図は、本実施形態における類似性置換マトリックス106eに格納される情報の一例を示す図である。第7図に示すように、類似性置換マトリックス106eに格納される情報は、識別情報(第7図に示す「A」、「R」、・・・、「X」)の組み合わせに応じて割り当てられた類似性スコア値(第7図は、例えば、識別情報が「A」と「N」の組み合わせの場合、類似性スコア値が「-1」であることを示している。)で構成されている。

30

また、アライメント後識別配列情報ファイル106fは、後述するアライメント実行部102gによりアライメントが実行された識別配列情報を格納するアライメント後識別配列情報格納手段である。ここで、アライメント後識別配列情報ファイル106fに格納される情報について第8図を参照して説明する。

第8図は、本実施形態におけるアライメント後識別配列情報ファイル106fに格納される情報の一例を示す図である。第8図に示すように、アライメント後識別配列情報ファイル106fに格納される情報は、識別配列情報に対応するアミノ酸配列識別情報(第8図に示す「A001」と、アライメントが実行された識別配列情報(第8図に示す「A」、「B」、「C」、・・・、「P」)と、を相互に関連付けて構成されている。

40

また、再変換アミノ酸配列情報ファイル106gは、後述する再変換部102hにより再変換されたアミノ酸配列情報を格納する再変換アミノ酸配列情報格納手段である。ここで、再変換アミノ酸配列情報ファイル106gに格納される情報について第9図を参照して説明する。

第9図は、本実施形態における再変換アミノ酸配列情報ファイル106gに格納される情報の一例を示す図である。第9図に示すように、再変換アミノ酸配列情報ファイル106gに格納される情報は、アミノ酸配列識別情報(第9図に示す「A001」と、再変

50

換されたアミノ酸配列情報（第9図に示す「Ala」、「Phe」、「Trp」、・・・、「Lys」）と、を相互に関連付けて構成されている。

また、第2図において、通信制御インタフェース部104は、配列情報処理装置100とネットワーク300（またはルータ等の通信装置）との間における通信制御を行う。すなわち、通信制御インタフェース部104は、他の端末と通信回線を介してデータを通信する機能を有する。

また、第2図において、入出力制御インタフェース部108は、入力装置112や出力装置114の制御を行う。ここで、出力装置114としては、モニタ（家庭用テレビを含む）の他、スピーカを用いることができる（なお、以下においては出力装置114をモニタとして記載する場合がある）。また、入力装置112としては、キーボード、マウス、および、マイク等を用いることができる。また、モニタも、マウスと協働してポインティングデバイス機能を実現する。

また、第2図において、制御部102は、OS（Operating System）等の制御プログラム、各種の処理手順等を規定したプログラム、および所要データを格納するための内部メモリを有し、これらのプログラム等により、種々の処理を実行するための情報処理を行う。制御部102は、機能概念的に、二次構造取得部102a、二次構造・類似性分類情報格納部102b、識別配列情報格納部102c、変換部102d、二次構造置換マトリックス格納部102e、類似性置換マトリックス格納部102f、アライメント実行部102gおよび再変換部102hを備えて構成される。

ここで、二次構造取得部102aは、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得する二次構造取得手段である。

また、二次構造・類似性分類情報格納部102bは、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を格納する二次構造・類似性分類情報格納手段である。

また、識別配列情報格納部102cは、上記の二次構造・類似性分類情報に基づいてアミノ酸配列情報を構成するアミノ酸情報を識別情報に変換した識別配列情報を格納する識別配列情報格納手段である。

また、変換部102dは、二次構造取得部102aにて取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換する変換手段である。ここで、変換部102dは、第10図に示すように、分類情報参照変換部102iおよび識別配列情報検索部102jをさらに含んで構成される。

第10図は、本発明が適用される本システムの変換部102dの構成の一例を示すブロック図であり、該構成のうち本発明に関する部分のみを概念的に示している。

第10図において、分類情報参照変換部102iは、二次構造・類似性分類情報格納部102bにて格納された二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換する分類情報参照変換手段である。

また、識別配列情報検索部102jは、識別配列情報格納部102cにて格納された識別配列情報から、アミノ酸配列情報に対応する識別配列情報を検索する識別配列情報検索手段である。

再び第2図に戻り、二次構造置換マトリックス格納部102eは、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを格納する二次構造置換マトリックス格納手段である。

また、類似性置換マトリックス格納部102fは、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納する類似性置換マトリックス格納手段である。

また、アライメント実行部102gは、変換部102dにて変換された複数の識別配列

10

20

30

40

50

情報に対してアライメントを実行するアライメント実行手段である。ここで、アライメント実行部 102g は、第 11 図に示すように、置換マトリックス基準アライメント実行部 102k をさらに含んで構成されている。

第 11 図は、本発明が適用される本システムのアライメント実行部 102g の構成の一例を示すブロック図であり、該構成のうち本発明に関する部分のみを概念的に示している。

第 11 図において、置換マトリックス基準アライメント実行部 102k は、二次構造置換マトリックス格納部 102e にて格納された二次構造置換マトリックスおよび/または類似性置換マトリックス格納部 102f にて格納された類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行する置換マトリックス基準アライメント実行手段である。

10

再び第 2 図に戻り、再変換部 102h は、アライメント実行部 102g にてアライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換する再変換手段である。

なお、これら各部によって行なわれる処理の詳細については、後述する。

(システムの処理)

次に、このように構成された本実施の形態における本システムの処理の一例について、以下に第 12 図等を参照して詳細に説明する。第 12 図は、本実施形態における本システムのメイン処理の一例を示すフローチャートである。

まず、配列情報処理装置 100 は、二次構造取得部 102a の処理により、アミノ酸配列情報を取得してアミノ酸配列関連情報ファイル 106a の所定の記憶領域(第 3 図に示す「アミノ酸情報」の項)に格納し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得し、取得された二次構造に関する情報をアミノ酸配列関連情報ファイル 106a の所定の記憶領域(第 3 図に示す「二次構造情報」の項)に格納する(ステップ SA-1)。なお、二次構造の取得には、既存の二次構造予測プログラム(例えば、PSI-PRED など)などを用いて予測して取得、また、対応する立体構造が既知の場合には既存の二次構造判定プログラム(例えば、DSSP、STRIDE など)などを用いて取得してもよい。また、ステップ SA-1 において、予めアライメントしたアミノ酸配列情報を取得してもよい。

20

ついで、配列情報処理装置 100 は、制御部 102 の処理により、アミノ酸配列情報を構成するアミノ酸情報を、同じ類似性を有するアミノ酸を同じグループに属させることにより 20 種類の天然アミノ酸を複数のグループに分類した情報である予め定めたアミノ酸類似性分類情報に基づいて分類し、分類されたグループを示す類似性情報をアミノ酸配列関連情報ファイル 106a の所定の記憶領域(第 3 図に示す「類似性情報」の項)に格納する。なお、アミノ酸類似性分類情報は、例えば、アミノ酸間の置換の起こりやすさをスコア値で表したマトリックス形式の既知の情報である BLOSUM62 のスコア値などに基づいてクラスタリング(クラスター解析)を行ったことにより得られた複数のグループからなる情報でもよい。具体的には、アミノ酸類似性分類情報は、例えば、下記の表 2 に示すように、類似性を、アラニン(Ala)とグリシン(Gly)からなる第 1 の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第 2 の群、システイン(Cys)からなる第 3 の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第 4 の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第 5 の群、リシン(Lys)とアルギニン(Arg)からなる第 6 の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第 7 の群の 7 つに分類した情報でもよい。

30

40

(表2)

| アミノ酸 (カッコは1文字コード) | | | | |
|-------------------|-----------|-----------|-----------|-----------|
| 1 | A l a (A) | G l y (G) | | |
| 2 | A s p (D) | G l u (E) | A s n (N) | G l n (Q) |
| 3 | C y s (C) | | | |
| 4 | P h e (F) | H i s (H) | T r p (W) | T y r (Y) |
| 5 | I l e (I) | L e u (L) | M e t (M) | V a l (V) |
| 6 | L y s (K) | A r g (R) | | |
| 7 | P r o (P) | S e r (S) | T h r (T) | |

10

20

30

40

50

ついで、配列情報処理装置100は、変換部102dの処理により、ステップSA-1にて予測された二次構造(第3図に示す「二次構造情報」)およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性(第3図に示す「類似性情報」)に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造(第3図に示す「二次構造情報」)および同一の類似性(第3図に示す「類似性情報」)を有するアミノ酸情報を同一の情報として識別するための識別情報(例えば、A、B、Cなどの文字など)に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換し、アミノ酸配列関連情報ファイル106aの所定の記憶領域(第3図に示す「識別情報」の項)に格納する(ステップSA-2)。

ここで、配列情報処理装置100は、二次構造・類似性分類情報格納部102bの処理により、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を予め二次構造・類似性分類情報ファイル106bの所定の記憶領域(第4図参照)に格納し、ステップSA-2において、変換部102dは、分類情報参照変換部102iの処理により、二次構造・類似性分類情報ファイル106bに格納された二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換し、アミノ酸配列関連情報ファイル106aの所定の記憶領域(第3図に示す「識別情報」の項)に格納してもよい。なお、二次構造・類似性分類情報は、例えば、アミノ酸間の置換の起こりやすさをスコア値で表したマトリックス形式の既知の情報であるBLOSUM62のスコア値などに基づいてクラスタリング(クラスター解析)を行ったことにより得られた複数のグループを示す類似性情報と、複数に分類された二次構造を示す二次構造情報とを相互に関連付けて構成され、類似性情報と二次構造情報との組み合わせ毎に、対応する識別情報が割り当てられた情報でもよい。具体的には、二次構造・類似性分類情報は、例えば、二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、類似性を、アラニン(A l a)とグリシン(G l y)からなる第1の群、アスパラギン酸(A s p)とグルタミン酸(G l u)とアスパラギン(A s n)とグルタミン(G l n)からなる第2の群、システイン(C y s)からなる第3の群、フェニルアラニン(P h e)とヒスチジン(H i s)とトリプトファン(T r p)とチロシン(T y r)からなる第4の群、イソロイシン(I l e)とロイシン(L e u)とメチオニン(M e t)とバリン(V a l)からなる第5の群、リシン(L y s)とアルギニン(A r g)からなる第6の群、プロリン(P r o)とセリン(S e r)とトレオニン(T h r)からなる第7の群の7つに分類した情報でもよい。

また、配列情報処理装置100は、二次構造取得部102aの処理により、例えば外部システム200の外部データベースに格納されたアミノ酸配列情報(例えば、公共の立体構造データベースであるPDBに登録されているアミノ酸配列情報など)をネットワーク300を介して予め取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得し、分類情報参照変換部102iの処理により、取得したアミノ酸配列情報を構成するアミノ酸情報を二次構造・類似性分類情報に基づいて識別情報に予め

変換し、識別配列情報格納部 102c の処理により、変換した識別情報からなる識別配列情報を識別配列情報データベース 106c の所定の記憶領域に予め格納し、ステップ SA-2 において、変換部 102d は、識別配列情報検索部 102j の処理により、識別配列情報データベース 106c に格納された識別配列情報から、アミノ酸配列情報に対応する識別配列情報を検索し、検索した識別配列情報をアミノ酸配列関連情報ファイル 106a の所定の記憶領域（第 3 図に示す「識別情報」の項）に格納してもよい。

ついで、配列情報処理装置 100 は、アライメント実行部 102g の処理により、ステップ SA-2 にて変換された複数の識別配列情報に対してアライメントを実行し、アライメント後識別配列情報ファイル 106f の所定の記憶領域に格納する（ステップ SA-3）。

ここで、配列情報処理装置 100 は、二次構造置換マトリックス格納部 102e の処理により、二次構造（例えば、ヘリックス、シートおよびその他）の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックスを二次構造置換マトリックス 106d の所定の記憶領域に予め格納し、および/または、類似性置換マトリックス格納部 102f の処理により、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを類似性置換マトリックス 106e の所定の記憶領域に予め格納し、ステップ SA-3 において、アライメント実行部 102g は、置換マトリックス基準アライメント実行部 102k の処理により、二次構造置換マトリックス 106d に格納された二次構造置換マトリックスおよび/または類似性置換マトリックス 106e に格納された類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行し、アライメント後識別配列情報ファイル 106f の所定の記憶領域に格納してもよい。

なお、ステップ SA-3 において、アライメント実行部 102g は、置換マトリックス基準アライメント実行部 102k の処理により、予め定めた係数により重み付けされた二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行し、アライメント後識別配列情報ファイル 106f の所定の記憶領域に格納してもよい。例えば、二次構造置換マトリックスを重んじる場合は、例えば二次構造置換マトリックスの構造スコア値に例えば 1 以上の数値を掛ける、および/または、例えば類似性置換マトリックスの類似性スコア値に例えば 1 未満の数値を掛ける、などすることにより、類似性スコア値に対して構造スコア値が大きくなるように重み付けを行ってもよい。

ここで、配列情報処理装置 100 は、再変換部 102h の処理により、ステップ SA-3 にてアライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換し、再変換アミノ酸配列情報ファイル 106g の所定の記憶領域に格納してもよい。具体的には、配列情報処理装置 100 は、再変換部 102h の処理により、ステップ SA-3 にてアライメントが実行された識別配列情報に対応するアミノ酸配列識別情報と同一のアミノ酸配列識別情報のアミノ酸配列情報を取得し、当該識別配列情報を構成する各識別情報を、取得したアミノ酸配列情報のアミノ酸情報と対応付けることにより再変換し、再変換アミノ酸配列情報ファイル 106g の所定の記憶領域に格納してもよい。

なお、本実施形態におけるアミノ酸の類似性は、アミノ酸の疎水性のほか、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性を用いてもよい。これにより、アミノ酸配列情報の立体構造などに影響するアミノ酸の性質（疎水性、親水性、酸性、塩基性、荷電状態など）の類似性を考慮してアライメントが実行された識別配列情報を効果的に得ることができる。

これにて、メイン処理が終了する。

【実施例】

（関数を用いたモデルの選定）

第 13 図は、本発明を有効に適用するための処理を示すフローチャートである。以下に、第 13 図に示す本フローチャートを詳細に説明する。

（1）ホモロジー検索

目的タンパク質のアミノ酸配列を用いてホモロジー検索を行ない、アライメントを得る

10

20

30

40

50

。このとき本発明の `pred_fasta` のホモロジー検索モードの実行をする。さらに加えて複数のホモロジー検索プログラム（例えば、`BLAST`、`FASTA` など）を用いて探索するのが望ましい。これは同じ参照タンパク質を探索してもアライメント方法が異なる場合や、あるホモロジー検索プログラムでは探索できなかった参照タンパク質を他のホモロジー検索プログラムが探索できる場合など、モデル作成時に様々なアライメントや参照タンパク質を用いたほうがより精度の高いモデルを作成できる可能性があるためである。

(2) モデリング

(1) で得られたアライメントを用いて、`FAMS` などのホモロジーモデリングプログラムにより目的タンパク質のモデリングを行なう。ただし、(1) で得られた全てのアライメントについてモデルを作成するのは非常に非能率的であるので、`e` 値やホモロジーの基準でモデリングを行なうアライメントを絞った。

10

(3) 目的タンパク質の二次構造を予測する

目的タンパク質のアミノ酸配列から、二次構造予測プログラム `PSI-PRED` を用いて目的タンパク質の二次構造を予測する。`PSI-PRED` はアミノ酸配列から二次構造を予測するプログラムであり、かなり高い信頼度（予測率の指標である Q_3 値は約 78% とされている）で二次構造の予測を行なうことが可能である。予測結果は第 14 図に示す形式のファイルに出力される。

第 14 図における `MA-1` ~ `MA-6` について以下に説明する。

`MA-1` : 予測する目的タンパク質アミノ酸配列の残基番号（1 から始まる）

20

`MA-2` : アミノ酸の種類（一文字表記）

`MA-3` : この残基の二次構造判定結果（`MA-4` ~ `MA-6` の値のうち一番高いものの二次構造を表記）

以下の 3 文字で表現される。

H : ヘリックス

E : シート

C : その他（コイル）

`MA-4` : この残基がコイルである確率

`MA-5` : この残基がヘリックスである確率

`MA-6` : この残基がシートである確率

30

(4) 作成したモデルについて、二次構造を判定する

(2) で作成したモデルについて、その 3 次元座標から二次構造判定プログラム `STRIDE` により二次構造を判定した。これによりモデルの各残基について、どのような二次構造をとっているかを判定することができる。

(5) 二次構造適合率の算出

モデルの二次構造判定結果を二次構造予測結果と比較することにより二次構造適合率を求める。この割合の算出方法は以下の通りである。

モデル中のある残基についての二次構造適合率を求めたい場合、(4) において `STRIDE` により求めたモデルの残基ごとの二次構造判定結果から、その残基についての二次構造（ヘリックス、シート、その他（コイル））を求める。次に(3)において `PSI-PRED` により求めた残基ごとの二次構造確率から当該の残基について、求めた二次構造についての確率を求める。この確率がその残基についての二次構造適合率であり、全ての残基について二次構造適合率を求めて和を算出する。この和（二次構造適合残基数）は二次構造の確率を足し合わせたもので、二次構造が適合した残基数を意味している。この二次構造適合残基数をモデル残基数で割って求めた平均を、このタンパク質全体についての二次構造適合率とした。

40

`PSI-PRED` は高い予測率のため、ある程度信頼度の高い二次構造適合率を求めることが可能である、と考えられる。

(6) ホモロジー検索時の `e` 値の算出

モデル作成に用いたアライメントに対応するホモロジー検索結果から、そのアライメン

50

トについての e 値を求めた。ただし、ホモロジー検索時に用いる立体構造データベースは異なる場合があり一概に e 値を比較することはできない。e 値は以下の式で求められるため、異なるデータベースを使用した場合にはデータベース間の大きさの比をかけることにより e 値を補正した。

DB - size exp - value (= E - value)

(7) 疎水性相互作用の指標 (enesosui 値) の算出

タンパク質が立体構造を構成する因子として、疎水性相互作用がある。この疎水性相互作用を評価することにより、予測したモデルの立体構造を評価する指標として利用することが可能である。この指標を求めるために、「enesosui 値」を用いた。enesosui 値は疎水性相互作用の指標であり、第 15 図に示す表で示したアミノ酸残基の側鎖中の疎水性を示す代表的な原子同士の距離が、ある閾値 (例えば、6.6) より小さい場合、疎水性相互作用をしていると見なし、その残基に enesosui 値として「-1」を与える。この場合、溶媒から遮蔽されており溶媒と接触している状態よりもエネルギー的に安定であると考えられる。enesosui 値の合計はタンパク質全体における疎水性相互作用の指標とすることができる。

また、enesosui により SCOP ドメインデータベース (95% Non-Redundant 化したもの) 中のドメイン構造について、その残基数と enesosui 値との関係を調べた (第 16 図参照)。ここで、第 16 図は、SCOP ドメインデータベース (95% Non-Redundant: total_seq_dom_all_95 (7433 ドメイン)) 中のドメイン構造について、その残基数と enesosui 値との関係を調べた結果をまとめた表を示す図である。その結果、「 $y = 7.1118x - 17.34$ 」という検量線 (第 16 図の表中の直線を参照) を得た。構築したモデルについて enesosui 値を計算したとき、この検量線付近の値をとることが望ましいといえる。この検量線を考慮した評価を行なうには、enesosui 適合率を算出する必要がある。この enesosui 適合率は以下の数式 1 により求められる。

$$\text{enesosui 適合率} = 1 - C \left(\frac{E_0 - E}{E_0} \right)^2 \dots \text{(数式 1)}$$

E = モデルの enesosui 値

E_0 = モデル残基数のときに検量線により求められる enesosui 値

C = 定数

(8) 二次構造率を求める

二次構造率とは、モデル構造について二次構造判定プログラム STRIDE により二次構造判定を行ない、二次構造 (ヘリックス、シート) であると判定された残基の数を数え、その値をモデル残基数で割ったものである。二次構造が長くできているモデルほど、この値は大きくなる。

(9) スコア値の算出

二次構造、疎水性エネルギー、e 値、最低の e 値をもとにモデルに対するスコアを計算する (二次構造率、enesosui 適合率がなくてもスコアによる評価は可能である)。具体的には、以下の数式 2 に示すように、二次構造適合率と疎水性相互作用の指標 (enesosui 値) の積に e 値と最低の e 値の対数比による重みをかけることによりスコアを計算する。このスコアはアライメント時に得られる e 値だけではなくモデルの立体構造状態も考慮したものである。さらに (任意的に) 二次構造率、enesosui 適合率をかけることにより、モデルの構造をより厳しく判定することができる。

$$\text{スコア} = - \left(\text{二次構造適合率} \times \text{enesosui 値} \right) \times \frac{-\log_{10}(e) + \alpha}{-\log_{10}(e_{\min}) + \alpha} \times$$

(二次構造率 \times enesosui 適合率)

... (数式 2)

10

20

30

40

50

$m_{i n}$ はモデル作成時に用いたアライメント中で、最も低い e 値をもつものの e 値である。は定数であり、 $m_{i n}$ 値の対数が 0 で 100 の値を持ち、-3 で 50 の値を、-10 で 40 の値を持つようにして、その間では線形の関係があるような値に 1 を加算した値である。

ここで 1 を加算する理由を述べる。通常の FASTA、PSI-BLAST などのホモロジー検索プログラムでは、 e 値の上限が 10 で区切られて検索結果が出力される。このため、 $-\log_{10}(\quad)$ の最低値は -1 であり、最高値は + である。このため、 e 値、最低 e 値の組み合わせによっては $-\log_{10}(\quad)$ と $-\log_{10}(m_{i n})$ の符号が異なる場合がある。そこで、符号をそろえるために分子・分母の値に 1 を加えた。この値は、 $\log_{10}(\quad)$ (ホモロジー検索プログラムでの 値の上限値) で表わされる。

10

この定数により、 e 値が低い場合には二次構造適合率、疎水エネルギー値の重みが大きくなり、逆に e 値が高い場合には e 値の重みが大きくなる。これは e 値が大きい場合、目的タンパク質とその参照タンパク質とは類似性が薄く、構造もあまり類似していないと思われるため、二次構造や疎水エネルギーといったモデル構造に依存する値には重みを与えないようにしたためである。

このスコアをモデルごとに算出する。

(10) スコア値の比較

(9) によりモデルごとにスコアが得られるが、これらのスコアを比較することによりモデルの精度の指標とした。スコア値が大きいほどアライメント的にも、立体構造的に見ても精度が高いモデルである、ということができると考えられる。

20

結果ファイルはスコア値順に出力される。

次に、上記プロセスの入出力画面について、第 17 図および第 18 図を参照して説明する。

第 17 図は本発明の入力ファイルの一例を示す図である。入力ファイルとして目的タンパク質のアミノ酸配列を記したファイルを用いる。このファイルは FASTA フォーマット (1 行目に ">タンパク質名"、2 行目にアミノ酸配列) でなければならない。

また、第 18 図は、本発明の出力ファイルの一例を示す図である。このファイルはスコア値順に出力される。以下に第 18 図に示す MB-1 ~ MB-15 について詳細に説明する。

MB-1 : 探索された参照タンパク質名

30

MB-2 : 参照タンパク質のアミノ酸配列の全長

MB-3 : スコア値

MB-4 : e 値

MB-5 : ホモロジー値

MB-6 : 使用したホモロジー検索プログラムによる結果における順位

MB-7 : 目的タンパク質のモデリング領域 (残基番号)

MB-8 : 参照タンパク質の対応領域 (残基番号)

MB-9 : 使用したホモロジー検索プログラムの略名

MB-10 : 使用したデータベース名

PDB ... 立体構造をもつ鎖単位のタンパク質をアミノ酸配列に変換して作成したデータベース

40

SCOP ... PDB を SCOP ドメインにより分類したアミノ酸配列データベース

PDB95 ... PDB データベースをホモロジー 95% 以上によりクラスタリングし、配列冗長性を排除したデータベース

SCOP95 ... SCOP ドメインデータベースをホモロジー 95% 以上によりクラスタリングし、配列冗長性を排除したデータベース

MB-11 : モデルの $e_{n e s o s u i}$ 値

MB-12 : モデルと参照タンパク質の $e_{n e s o s u i}$ 値の比 (モデルの $e_{n e s o s u i}$ 値 / 参照タンパク質の対応する領域の $e_{n e s o s u i}$ 値)

MB-13 : MB-12 の値をモデルの残基数で割ったもの

50

MB - 14 : 目的タンパク質全体の二次構造適合残基数

MB - 15 : 二次構造適合率 (MB - 14 の値をモデルの残基数で割ったもの)

(本発明の P R E D _ F A S T A)

第 19 図および第 20 図は、本発明の処理を示すフローチャートである。以下に、第 19 図 (ペアワイズモードの場合) および第 20 図 (ホモロジー検索モードの場合) に示す本フローチャートを詳細に説明する。

(1) アミノ酸類似性データベース (上述の実施形態におけるアミノ酸類似性分類情報に対応) の作成

通常、アライメントを行なう場合にはアミノ酸同士の相同性スコアを計算する。このときに使用される置換マトリックスとして代表的なものに B L O S U M 6 2 がある。このファイルには例えば次のようなことが書かれている。 F (P h e) と Y (T y r) との間のスコア値 (置換のしやすさを示す) は 3 であり、 F (P h e) と P (P r o) との間のスコア値は - 4 である。このことから F と P をアライメントするよりも F と Y をアライメントしたときに、より高いスコアが付くようにしている。この置換マトリックスに記されている、あるアミノ酸同士の置換のしやすさはそれらのアミノ酸の間の類似性を示している。このため、 B L O S U M 6 2 に記されているスコア値をもとにクラスタリングを行ない、 20 種の天然アミノ酸をその特質により 7 つのグループに分類した。

(2) 二次構造・アミノ酸類似性データベース (第 5 図に示す識別配列情報データベース 106c に対応) の準備

あらかじめ P R E D _ F A S T A で使用する二次構造・アミノ酸類似性データベース (第 5 図に示す識別配列情報データベース 106c に対応) を準備しておく。

そして、タンパク質の立体構造から二次構造とその位置を判定するプログラム「 S T R I D E 」により、通常ホモロジー検索に用いる立体構造データベースに収載されているタンパク質について二次構造判定をする。これは通常ホモロジー検索を行なったときに、その検索されたタンパク質についての二次構造情報を必要とする場合があるためである。今回は立体構造データベース中のタンパク質配列について、ホモロジー 95% 以上かつ e 値 0.01 以下の基準でクラスタリングしたものをを用いた。

ここで、残基ごとに二次構造の違い (ヘリックス、 シート、 その他) 3 種類に分類できるが、この情報のほかに (1) で作成したアミノ酸類似性分類情報 (20 種のアミノ酸を 7 つのグループに分類) を組み合わせて表現したアミノ酸を作成する。すなわち、第 21 図に示すように、天然のアミノ酸をその二次構造・アミノ酸類似性の基準から便宜上二次構造 3 パターン × 疎水性 7 グループ = 21 種のアミノ酸として新たな配列に置き換えるという作業を行なった (第 4 図に示す二次構造・類似性分類情報ファイル 106b に対応) 。なお、第 21 図に示す例は一例であり、これ以外のグループに分類してもよい。

(3) アライメント、参照タンパク質の準備

まず、目的・参照タンパク質のアミノ酸配列間の関係を示したアライメントを作成する。アライメント用ソフトウェアとして、 F A S T A 、 P S I - B L A S T などがある。これらのプログラムを用いたときに、ある一定の閾値以下の条件で得られたアライメントは信頼性が高いといえる。今回の閾値として、 P S I - B L A S T を用いた際の e 値が 0.01 以下とした。このようにして得られたアライメント領域については、その対応するアミノ酸には大きな誤りはないといえる。しかし、このような通常のアライメントソフトはアミノ酸配列の文字情報のみからアライメントを得ているため、参照タンパク質側の二次構造を分断したアライメントを作成してしまう恐れがある。このため、 F A S T A を用いて二次構造・アミノ酸類似性を崩さないアライメントを以下の手順により作成した。

(4) 目的タンパク質の配列から二次構造を予測する (第 2 図に示す二次構造取得部 102a の処理に対応)

与えられた目的タンパク質の配列を用いて、 P S I P R E D (p r o t e i n s t r u c t u r e p r e d i c t i o n) を実行する。 P S I P R E D は高速で精度の高い二次構造予測プログラムであり、 P S I - B L A S T で作成されるマルチプルアライメントを単純なニューラルネットで評価することにより二次構造を予測する。 P S I - P R E

10

20

30

40

50

D中でマルチプルアライメントの結果を使用する理由は、進化の過程で機能を維持する上で構造の保存性は高く、構造を変えないように残基の置換も制約されるといったことも考慮して二次構造を予測することができるためである。

(5)アライメントを二次構造・アミノ酸類似性により置換した配列に変換(第2図に示す変換部102d、分類情報参照変換部102iおよび識別配列情報検索部102jの処理に対応)

(3)で得られたアライメントの配列を(2)で定めた二次構造とアミノ酸の類似性の基準(上述の実施形態における二次構造・類似性分類情報に対応)によりアミノ酸一文字表記を変換し、新たなアミノ酸配列(上述の実施形態における識別配列情報に対応)を作成する。

10

目的タンパク質の配列に関しては、(4)で行なった二次構造予測結果とアミノ酸類似性(上述の実施形態におけるアミノ酸類似性分類情報に対応)に基づいて変換する。参照タンパク質の配列に関しては(2)で作成した二次構造・アミノ酸類似性データベース(第5図に示す識別配列情報データベース106cに対応)から当該参照タンパク質についての配列を抜き出した。このようにして目的・参照タンパク質共に二次構造・アミノ酸類似性により変換した配列を用意した。この新たに用意したアミノ酸配列は、(3)で得られたアライメント領域について作成した。

(6)FASTAによる再アライメント(第2図に示すアライメント実行部102gおよび置換マトリックス基準アライメント実行部102kの処理に対応)

(5)で用意した目的・参照タンパク質配列(上述の実施形態における識別配列情報に対応)をFASTAにより再アライメントした。ここで用いた文字比較のための置換マトリックスはアミノ酸一文字表記の比較のためのBLOSUM62ではなく、二次構造とアミノ酸類似性を表わす文字(疎水性が同程度のアミノ酸残基であり、且つ同じ二次構造をもつ場合同じ文字として扱われる)の2つのマトリックス(上述の実施形態における二次構造置換マトリックスおよび類似性置換マトリックスに対応)から作成されたものである(第22図および第23図を参照)。デフォルトは1:1の均等の重み付けであるが、オプションによりその重み付けを変更することも可能である。ここで作成したマトリックスを用いることにより、二次構造とアミノ酸類似性情報を加味したアライメントを得ることができる。

20

ここで、第22図は二次構造によるスコアマトリックスの一例を示す図である。第22図に示すスコア(上述の実施形態における構造スコア値に対応)は、二次構造が一致したときに高くなるよう設定し、逆に相反する二次構造の場合にはアライメントしにくいように低く設定した。

30

また、第23図はアミノ酸類似性によるスコアマトリックスの一例を示す図である。第23図のスコアマトリックスは、(1)で定義したアミノ酸類似性により分類した識別情報のグループをもとに作成した。例えば、ある識別情報のグループ内のアミノ酸同士のスコア(上述の実施形態における類似性スコア値に対応)は、そのグループに属する全てのアミノ酸の組み合わせによるBLOSUM62スコアの平均(小数点以下は切り捨て、整数で表記)を割り当てた。同様に、あるグループに属するアミノ酸とあるグループに属するアミノ酸との間のスコアは、その2つのグループに属する全てのアミノ酸の組み合わせによるBLOSUM62スコアの平均を割り当てた。

40

(7)二次構造・アミノ酸類似性による配列をもとのアミノ酸配列に戻す(第2図に示す再変換部102hの処理に対応)

(6)で得られたアライメントは、二次構造とアミノ酸類似性により変換された配列で出力される。このため、もとのアミノ酸配列に再変換したアライメントを作成した。つまり、二次構造、アミノ酸類似性により変換された配列を20種類のアミノ酸で表現される通常の配列に再変換する。

(8)モデリング

(7)で得られた修正アライメントを元に、適当なモデリングソフト、例えばFAMSを用いて立体構造を構築し、出来上がったモデルの二次構造が分断されていないかどうか

50

を検査する。

上記(1)～(8)に示した手順は、第19図に示すペアワイズモード(1対1でアライメントすること)の場合の手順であり、これらの手順は、第20図に示すホモロジー検索モードの場合には、目的タンパク質配列について、データベースを(2)で作成した二次構造・アミノ酸類似性データベース(第2図における識別配列情報データベース106cに対応)全体に対して検索することによって、通常行なわれるホモロジー検索としての利用も同様に可能である。

次に、本発明の入出力画面について、第24図から第31図を参照して説明する。

(1) ホモロジー検索モード

ホモロジー検索モードを使用する場合、目的タンパク質のアミノ酸配列を記したFASTA形式のファイルを入力ファイルとする。FASTA形式とは、ファイルの一行目に「>目的タンパク質名」、2行目に「目的タンパク質のアミノ酸配列」を記したファイルである。出力は通常のFASTAの出力と同じ形式である。

(例1；高ホモロジーの場合)

第24図は、高ホモロジーの場合における本発明の入力ファイルの一例を示す図である。第24図は、PDBに登録しており、かつPREDFASTAで用いられる立体構造データベースにも登録されている1A4FのA鎖(1A4FA)の配列を目的タンパク質test1としてPREDFASTAを実行した場合の入力ファイルである。

また、第25図は、PREDFASTAによる1位の参照タンパク質についての出力ファイルを示す図である。これにより自身の配列である1A4F_Aを探索できたことがわかる。第25図のアミノ酸配列を見るとホモロジーが91.489%となっているが、これは二次構造・アミノ酸類似性によりアミノ酸配列を変換してアライメントさせたときのホモロジーであるためである。

(例2；低ホモロジーの場合)

第26図は、低ホモロジーの場合における本発明の入力ファイルの一例を示す図である。第26図は、立体構造予測コンテストCASP5において出題された低ホモロジー配列(T0194：CAFASP9058)について、PREDFASTAによりホモロジー検索を行なった際の入力ファイルである。

また、第27図は、PREDFASTAにより1位で検索された参照タンパク質についての出力ファイルを示す図である。第27図において、ホモロジー18.797%と出力されているが、これは二次構造・アミノ酸類似性によりアライメントさせたときのホモロジーであり、実際のホモロジーは7%と非常に低い。このように非常にホモロジーが低く、他のホモロジー検索プログラムでは探索できない場合でもPREDFASTAにより参照タンパク質を探索することが可能である。

(例3；アライメントを伸長させる)

PREDFASTAにより参照タンパク質を探索してアライメントを得ると、通常のホモロジー検索プログラムによるアライメントよりも長いアライメントが得られることがある。このような場合、通常のホモロジー検索プログラムによるアライメントとPREDFASTAによるアライメントを組み合わせることにより、より長いアライメントを得ることができる。この例を以下に述べる。

第28図は、立体構造予測コンテストCASP5において出題されたタンパク質(T0176：CAFASP8880)のアミノ酸配列についてPSI-BLAST、PREDFASTAにより得られたアライメントを示す図である。ともに同じ参照タンパク質を探索しているがPREDFASTAの方がアライメント領域を長くアライメントをしている。

また、第29図は、PSI-BLASTによるアライメントにPREDFASTAによるアライメントの末端部分を付け加えることにより、長いアライメントを得ることができた場合を示す図である。

(2) ペアワイズモード

ペアワイズモードで使用する場合、目的タンパク質と参照タンパク質とのアライメント

ファイルを入力とする。第30図は、ペアワイズモードの場合の本発明の入力ファイルの一例を示す図である。この入力ファイルは1行目に「>目的タンパク質名」、2行目に「アライメントされた目的タンパク質のアミノ酸配列」、3行目に「>参照タンパク質名」、4行目に「アライメントされた参照タンパク質のアミノ酸配列」を記したものである。

ここで、第30図に示す例では、PDBに登録してある1A4FのA鎖(1A4F__A)の一部を目的タンパク質test3とした。このtest3と1A4F__Aとのアライメント(ホモロジーは100%)について、1A4F__Aの二次構造が壊れるようにわざとギャップを入れたアライメントファイルを入力としてPREDFASTAによるアライメント修正を実行した。

また、第31図は、本発明の出力ファイルの一例を示す図である。第31図に示す出力ファイルの例では、二次構造上に入ったギャップは除かれてきちんとしたアライメントになっていることがわかる。

(埋め込み)

第32図は、本発明でアライメントを行い立体構造モデルが得られた後、より現実に近い立体構造モデルを得るための処理を示すフローチャートである。以下に、第32図に示す本フローチャートを詳細に説明する。

あるタンパク質Aを参照タンパク質として作成したモデルA(第32図におけるモデル1)より、あるタンパク質Bを参照タンパク質として作成したモデルB(第32図におけるモデル2)の方がより長くモデリングできたとする。しかし領域の短いモデルAの方が精度よい構造である場合(すなわち高優先順位である場合)、モデルAの構造にモデルBの末端の立体構造を貼り付けてモデリング領域を伸長させることができる場合がある。

このような場合に局所構造(短いモデル:ベースとなる)の末端にグローバル構造(長いモデル)の構造を貼り付けることによりモデルの伸長を行なう方法論を開発した。これらの伸長方法には以下の4つがあり、以下にその手順を述べる。

手順(1): つなぎ目の残基の距離を判定して伸長(第33図、第34図参照)

手順(2): 末端を優先させるフィッティングによる伸長(第35図参照)

手順(3): 領域を移動させるフィッティングによる伸長(第37図参照)

手順(4): 二次構造等を主に考慮して断片を再検索して伸長(第39図、第40図参照)

以下に、これらの各手順の方法論などについて詳細に説明する。

手順(1): つなぎ目の残基の距離を判定して伸長

第34図は、本発明によりつなぎ目の残基の距離を判定して伸長する場合の処理の一例を示すフローチャートである。

局所構造(短いモデル:ベースとなる)の末端にグローバル構造(長いモデル)の構造を貼り付けることによりモデルの伸長を行なう方法論を開発した。この方法論では構造のつなぎ目の残基間の距離を考慮して、つなぎ目ができるべく滑らかになるように局所構造を削っていく。以下にその手順を述べる。

(1) アライメントにより共通領域を求める

グローバル構造(長いモデル)のアミノ酸配列と、局所構造(短いモデル:ベースとなる)のアミノ酸配列をアライメントし、共通領域を求めた。つまり、局所構造モデルとグローバル構造モデルのアミノ酸配列を用いてアライメントして対応する残基、共通領域を求めた。局所構造とグローバル構造が全く同一の領域をモデリングしている場合、つまりモデルの長さや領域が同じ場合はモデルの伸長させることができないので、ここで終了する。

(2) フィッティングを行なう

(1)で求めた共通領域を用いて最小二乗法などによりフィッティングを行なう。これは局所構造モデルとグローバル構造モデルの座標系を揃えるために行なう必要がある。局所構造とグローバル構造をフィッティングさせたときのRMSDが2より大きい場合は、この埋め込みは行なわない。これはあまりにも構造が異なる場合、埋め込みを行なってもよいモデルを作成することができないためである。

10

20

30

40

50

(3) 局所構造モデルをグローバル構造モデルに埋め込む

座標系を一致させた局所構造モデルの座標を、グローバル構造モデル座標上の対応する領域にはめ込む。

(4) 境界の残基間距離を計算する

(3)で埋め込んだ局所構造モデルとグローバル構造モデルとの境界残基同士のC 原子間距離を計算する。具体的には、例えば、局所構造モデルとグローバル構造モデルの境界部分の隣接する残基のC 原子間距離を求め、ある閾値(例えば、8.0)以下かどうかを調べる。この距離は次の(5)において、つなぎ目の構造を評価するときに必要な。

(5) つなぎ目の判定

(4)で計算したつなぎ目残基のC 原子間距離がある閾値(例えば、8.0 : F A M Sにより処理可能な限界値など)以上離れている場合、局所構造モデルの末端を一残基カットし、再び(2)のフィッティングに戻る。つまり、閾値よりもC 原子間距離が大きい場合は、局所構造モデルを末端から一残基カットして距離を調べる、という手順を繰り返す。例えば、グローバル構造モデルの残基2(第34図の残基番号2の残基)と局所構造モデルの残基1'(第34図の残基番号1'の残基)のC 原子間距離が閾値よりも大きかった場合、局所構造モデルの残基1'(第34図の残基番号1'の残基)をカットし、その残基については代わりにグローバル構造モデルの残基3(第34図の残基番号3の残基)を用いる。もう一方の末端側について距離が閾値よりも小さい場合はこのカットは行なわない。この手順を距離が閾値以下になり、つなぎ目が滑らかになるまで繰り返して埋め込みモデルを作成する。残基を削っていても最後までこの閾値を満たすことができない場合は埋め込みは不可能として終了する。

この作業を行なう理由は、あまりにも接続残基部分の距離が離れているとF A M Sでの再モデリング時に参照タンパク質座標異常として、参照タンパク質(局所構造モデルをグローバル構造モデルに埋め込んだもの)の座標を参照せずにF A M S内部のデータベースを使用してしまうためである。

このようにしてなるべく局所構造モデル領域を保存し、末端のみグローバル構造領域を使用するようにして座標ファイルをつなげる。

(6) F A M Sを用いて再モデリング(モデルを伸長させることができた場合のみ)

(5)までの作業で埋め込んで作成したP D B形式の座標ファイルを参照タンパク質として、F A M Sにより同一領域の再モデリング(ホモロジー100%のモデリング)を行なった。つまり、埋め込みを行った構造を参照タンパク質にして、F A M Sにより再モデリングを行った。再モデリングすることにより、つなぎ目部分の不自然な構造や、側鎖のぶつかりなどを解消して全体の形を整えることができる。

手順(2): 末端を優先させるフィッティングによる伸長

第35図は、本発明により、末端を優先させるフィッティングによる伸長を行う場合の処理の一例を示すフローチャートである。本発明者は、本方法により、局所構造(短いモデル: ベースとなる)の末端にグローバル構造(長いモデル)の構造を貼り付けることによりモデルの伸長を行なう方法論を開発した。この方法では伸長させようとする方の末端の構造を主に考慮してフィッティングさせる。この方法論について以下に述べる。なお、第35図では、局所構造の末端にグローバル構造の一部を貼り付けることによりモデルを伸長させる。以下に局所構造モデルの1~10残基目がグローバル構造モデルの6~15残基目と対応しているときのフィッティングの手順を一例に解説する。

(1) R M S Dによる閾値の設定

フィッティングの良し悪しの状態を評価するためにR M S Dを用いた。R M S D(根平均二乗距離: R o o t M e a n S q u a r e D e v i a t i o n)とは二つの構造を重ね合わせ、対応する残基間(C 間)の距離の離れ具合を表わす指標である。R M S Dが小さいほどよくフィッティングしているといえる。

フィッティング時のR M S Dによる閾値の初期値を2(デフォルト)とした。このR M S Dの閾値は、閾値以下のフィッティングができなかった場合に1ごとに高くしてい

10

20

30

40

50

く。条件を徐々に甘くすることによって最適なフィッティングを探索する。このRMSDの閾値を上限(デフォルト4)まで変化させることにより繰り返してフィッティングを行なう。これらの基準値は変更することも可能である。

(2) 共通領域全体についてのフィッティング

はじめに共通領域全体についてフィッティングを行なう。RMSDの閾値以下であればそのフィッティングを採用して終了する。共通領域全体においてよくフィッティングすることが理想であるが、互いの構造があまり似ていないことも多い。このようにRMSDが閾値よりも大きい場合は(3)に進んだ。

(3) 残基を削って繰り返し再フィッティング

先のフィッティングに用いた領域から2残基削って再びフィッティングを行なった。ここで削る2残基は、伸長させたい方ではない末端からの2残基とした。RMSDが閾値以下であればそのフィッティングを採用して終了するが、RMSDが閾値よりも大きい場合は再び2残基削ってフィッティングを行なう、ということを繰り返す。

(4) RMSDの閾値の変更

(3)により繰り返してフィッティングを行なうと、フィッティング領域が削られていく。フィッティング領域が3残基以下になった場合、短い残基でフィッティングさせても意味がない。このような場合フィッティングをせずにRMSDの閾値を1増やし、基準を甘くして(2)~(3)を繰り返した。このRMSDの閾値の変更とフィッティングは、RMSDの閾値が(1)で設定した上限(デフォルト4)に達し、かつフィッティング領域が3残基以下になるまで行なわれる。

最後まで閾値以下のフィッティングができなかった場合は、フィッティング不可能とみなして終了する。

また、(3)、(4)によるRMSDの閾値、フィッティング残基数の変化は、第36図に示す図を参照する。

(5) 構造の貼り付け

(4)まででフィッティングがうまく行なうことができた場合、フィッティングで用いた領域をのりしろとして、グローバル構造の座標を局所構造モデルの末端に貼り付けることによりモデルを伸長させる。この貼り付けのときに貼り付けられるグローバル構造部分が局所構造にぶつからないかを確認する。ここでぶつかってしまうようであれば、ぶつからなくなるまで貼り付けられるグローバル構造部分をカットする。

(6) FAMSを用いて再モデリング(モデルを伸長させることができた場合のみ)

(5)までの作業で埋め込んで作成したPDB形式の座標ファイルを参照タンパク質として、FAMSにより同一領域を再モデリング(ホモロジー100%のモデリング)を行なった。再モデリングすることにより、つなぎ目部分の不自然な構造や、側鎖のぶつかりなどを解消して全体の形を整えることができる。

手順(3): 領域を移動させるフィッティングによる伸長

第37図は、本発明によりフィッティング領域を移動させてフィッティングさせる場合の処理の一例を示すフローチャートである。本発明者は、本方法により、局所構造(短いモデル: ベースとなる)の末端にグローバル構造(長いモデル)の構造を貼り付けることによりモデルの伸長を行なう方法論を開発した。この方法では貼り付けようとする方の末端から反対の末端に向けてフィッティングさせる領域を移動させていく。これにより伸長させようとする末端の構造が少しくらい異なっても、フィッティング領域を内部に移動させることによりフィッティングを行なうことができる。この方法論について以下に述べる。なお、第37図では、局所構造の末端にグローバル構造の一部を貼り付けることによりモデルを伸長させるものであり、局所構造モデルの1~24残基目がグローバル構造モデルの6~29残基目と対応しているときのフィッティングの手順を一例に解説する。

(1) RMSDによる閾値、フィッティング残基数の設定

フィッティング時の基準であるRMSDによる閾値(デフォルト2)を設定する。この閾値はフィッティングができない場合に上限(デフォルト4)まで1ずつ高くしていく。

10

20

30

40

50

次にフィッティング残基数（フィッティングさせる領域の残基数）の初期値（デフォルト20残基）を設定する。このフィッティング残基数はフィッティングができない場合に下限（デフォルト12残基）まで2残基ずつ削って繰り返しフィッティングさせる。

このRMSDの閾値、フィッティング残基数の二つの条件を変化させることにより、フィッティング時の基準を徐々に甘くして繰り返しフィッティングを行なう。これらの基準値は変更することも可能である。

(2) 共通領域全体についてのフィッティング

はじめに共通領域全体についてフィッティングを行なった。RMSDの閾値以下であればそのフィッティングを採用して(6)へ移る。RMSDの閾値は2（初期値）から4（上限値）まで1ずつ変化し、フィッティングの具合を調べる。共通領域全体においてよくフィッティングすることが理想であるが、互いの構造があまり似ていないことも多い。このようにRMSDが閾値よりも大きい場合は(3)に進んだ。

(3) フィッティング領域を移動させて繰り返しフィッティング

伸長させたい末端側からフィッティング残基数分の領域についてフィッティングを行なった。RMSDが閾値以下であればそのフィッティングを採用するが、RMSDが閾値よりも大きい場合はこのフィッティング領域を伸長させたい方ではない末端側に2残基ずらし、再びフィッティングを行なう。これをフィッティング領域が反対側の末端まで移動するまで繰り返した。

(4) RMSDの閾値の変更

(3)により閾値以下のフィッティングができない場合、RMSDの閾値を1ずつ増やし条件を甘くして(3)を行なう。これはRMSDの閾値が2（初期値）から1ずつ変化させ4（上限値）になるまで(4を含む)行なわれる。

(5) フィッティングさせる残基数を削る

(1)～(4)を行なっても(RMSDの閾値が上限に達しても)フィッティングがうまくいかない場合、フィッティング残基数を2残基減らして(3)に戻ってフィッティングを繰り返す。この場合RMSDの閾値も(1)で設定した初期値（デフォルト2）に戻り、その上限（デフォルト4）になるまで変化する。

(5)により繰り返してフィッティングを行なうと、フィッティング領域が削られていく。削ることによってフィッティング残基数が下限以下になった場合は、フィッティングが不可能とみなして終了する。

(2)から(5)の手順において、RMSDの閾値、フィッティング残基数、フィッティング領域が複雑に変化する。この変化の順序について第38図に示す表にまとめた。

(6) 構造の貼り付け

(5)まででフィッティングがうまく行なうことができた場合、フィッティングで用いた領域をのりしるとして、グローバル構造の座標を局所構造モデルの末端に貼り付けることによりモデルを伸長させる。この貼り付けのときに貼り付けられるグローバル構造部分が局所構造にぶつからないかを確認する。ここでぶつかってしまうようであれば、ぶつからなくなるまで貼り付けられるグローバル構造部分をカットする。

(7) FAMSを用いて再モデリング（モデルを伸長させることができた場合のみ）

(6)までの作業で埋め込んで作成したPDB形式の座標ファイルを参照タンパク質として、FAMSにより同一領域の再モデリング（ホモロジー100%のモデリング）を行なった。再モデリングすることにより、つなぎ目部分の不自然な構造や、側鎖のぶつかりなどを解消して全体の形を整えることができる。

手順(4)：二次構造等を主に考慮して断片を再検索して伸長

このプログラムでは、手順(1)、手順(2)、手順(3)の結果、遺伝子が完全長ではない場合にその座標をさらに伸長する方法を提案する。本方法は、以下の3段階に大きく分かれる。

- 1) 二次構造データベース上で検索
- 2) 二次構造のマッチの度合いで1)の結果をふるいにかける
- 3) 二次構造データベースから選ばれた断片を用いて、伸長を行う

10

20

30

40

50

以下に、1)～3)について詳細に説明する。

1) 二次構造データベースの検索はプログラム P S I - B L A S T を用いている。二次構造データベースは、2、3、4、5個の二次構造単位を持ったデータベースがそれぞれ用意される。そのそれぞれに対して検索が行われる。検索が行われる順番は、5、4、3、2個の二次構造単位を持ったデータベースの順である。二次構造データベースの作成方法について後述する「二次構造データベースの作成法」を参照(第39図参照)。

2) P S I - B L A S T で得られた多数のアライメントを、p s i - p r e d による二次構造予測とのアライメントで得られたシーケンスの二次構造情報の一致度に基づいて振るいにかける。

3) 2) で残った複数のアライメントを用いて、F A M S によるモデリングを行う。できた複数モデル(フラグメントモデル)を使って、伸長を行う。 10

(i) 手順(1)、手順(2)、手順(3)で得られた座標データ(ベースモデル)とアミノ酸残基が一致する部分(のりしろ)を二残基以上持つフラグメントモデルを選び出す。

(i i) のりしろ部分で、ベースモデルの末端から二残基を最小自乗フィットさせる。

(i i i) この段階で全長でない場合は、他のフラグメントモデルから更に伸長する。即ち(i)(i i)を繰り返す。

ここで、第39図および第40図は、本発明による二次構造データベースの作成法を示すフローチャートである。二次構造データベースとは、クラスタリングされたP D B データベースを加工して作られたものである。具体的な方法として、P D B シーケンスとその二次構造の情報をもとに、2、3、4、5個の二次構造単位でシーケンスをフラグメント化したものである。このときの二次構造単位とは、ループ構造を含めないヘリックス、シート等を指す。 20

次に、本発明の入出力画面について説明する。

プログラムの入力として、引数に局所構造モデルの座標ファイル名とグローバル構造モデルの座標ファイル名を入力する。

手順(2)、手順(3)の方法では、これらの引数の他にどちらの末端がベースとなる局所構造であるのかを入力する。伸長させた構造は標準出力により出力される。

また手順(1)では目的タンパク質名を引数で入力することにより、伸長させた構造がこの名前のファイルに出力される。 30

手順(4)の方法では、手順(1)、手順(2)、手順(3)の方法で伸長されたP D B ファイル名を指定する。出力は単数および複数の座標ファイルとして出力される。

これらの出力は全て座標ファイル(P D B 形式)である。

(他の実施の形態)

さて、これまで本発明の実施の形態について説明したが、本発明は、上述した実施の形態以外にも、請求の範囲に記載した技術的思想の範囲内において種々の異なる実施の形態にて実施されてよいものである。

例えば、装置の分散・統合の具体的な形態は明細書および図面に示すものに限られず、その全部または一部を、各種の負荷等に応じた任意の単位で、機能的または物理的に分散・統合して構成することができる(例えば、グリッド・コンピューティングなど)。 40

例えば、配列情報処理装置100は、配列情報処理装置100とは別筐体で構成されるクライアント端末からの要求に応じて処理を行い、その処理結果を当該クライアント端末に返却するように構成してもよい。

また、実施形態において説明した各処理のうち、自動的に行なわれるものとして説明した処理の全部または一部を手動的に行うこともでき、あるいは、手動的に行なわれるものとして説明した処理の全部または一部を公知の方法で自動的に行うこともできる。

この他、上記文書中や図面中で示した処理手順、制御手順、具体的名称、各種の登録データや検索条件等のパラメータを含む情報、画面例、データベース構成については、特記する場合を除いて任意に変更することができる。

また、配列情報処理装置100に関して、図示の各構成要素は機能概念的なものであり 50

、必ずしも物理的に図示の如く構成されていることを要しない。

例えば、配列情報処理装置100の各部または各装置が備える処理機能、特に制御部102にて行なわれる各処理機能については、その全部または任意の一部を、CPU(Central Processing Unit)および当該CPUにて解釈実行されるプログラムにて実現することができ、あるいは、ワイヤードロジックによるハードウェアとして実現することも可能である。なお、プログラムは、後述する記録媒体に記録されており、必要に応じて配列情報処理装置100に機械的に読み取られる。

すなわち、ROMまたはHDなどの記憶部106などには、OS(Operating System)と協働してCPUに命令を与え、各種処理を行うためのコンピュータプログラムが記録されている。このコンピュータプログラムは、RAM等にロードされること
10
によって実行され、CPUと協働して制御部102を構成する。また、このコンピュータプログラムは、配列情報処理装置100に対して任意のネットワーク300を介して接続されたアプリケーションプログラムサーバに記録されてもよく、必要に応じてその全部または一部をダウンロードすることも可能である。

また、本発明にかかるプログラムを、コンピュータ読み取り可能な記録媒体に格納することもできる。ここで、この「記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM、EPROM、EEPROM、CD-ROM、MO、DVD等の任意の「可搬用の物理媒体」や、各種コンピュータシステムに内蔵されるROM、RAM、HD等の任意の「固定用の物理媒体」、あるいは、LAN、WAN、インターネットに代表されるネットワークを介してプログラムを送信する場合の通信回線や搬送波のように、短期にプログラ
20
ムを保持する「通信媒体」を含むものとする。

また、「プログラム」とは、任意の言語や記述方法にて記述されたデータ処理方法であり、ソースコードやバイナリコード等の形式を問わない。なお、「プログラム」は必ずしも単一的に構成されるものに限られず、複数のモジュールやライブラリとして分散構成されるものや、OS(Operating System)に代表される別個のプログラムと協働してその機能を達成するものをも含む。なお、実施の形態に示した各装置において記録媒体を読み取るための具体的な構成、読み取り手順、あるいは、読み取り後のインストール手順等については、周知の構成や手順を用いることができる。

記憶部106に格納される各種のデータベース等(アミノ酸配列関連情報ファイル106a~再変換アミノ酸配列情報ファイル106g)は、RAM、ROM等のメモリ装置、
30
ハードディスク等の固定ディスク装置、フレキシブルディスク、光ディスク等のストレージ手段であり、各種処理やウェブサイト提供に用いる各種のプログラムやテーブルやファイルやデータベースやウェブページ用ファイル等を格納する。

また、配列情報処理装置100は、既知のパーソナルコンピュータ、ワークステーション等の情報処理端末等の情報処理装置にプリンタやモニターやイメージスキャナ等の周辺装置を接続し、該情報処理装置に本発明の方法を実現させるソフトウェア(プログラム、データ等を含む)を実装することにより実現してもよい。

さらに、配列情報処理装置100の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷等に応じた任意の単位で、機能的または物理的に分散・統合して構成することができる。例えば、各データベースを独立したデータベース装置
40
として独立に構成してもよく、また、処理の一部をCGI(Common Gateway Interface)を用いて実現してもよい。

また、ネットワーク300は、複数の配列情報処理装置100を相互に接続する機能を有し、例えば、インターネットや、イントラネットや、LAN(有線/無線の双方を含む)や、VANや、パソコン通信網や、公衆電話網(アナログ/デジタルの双方を含む)や、専用回線網(アナログ/デジタルの双方を含む)や、CATV網や、IMT2000方式、GSM方式またはPDC/PDC-P方式等の携帯回線交換網/携帯パケット交換網や、無線呼出網や、Bluetooth等の局所無線網や、PHS網や、CS、BSまたはISDB等の衛星通信網等のうちいずれかを含んでもよい。すなわち、本システムは、有線・無線を問わず任意のネットワークを介して、各種データを送受信することができる
50

。以上詳細に説明したように、本発明によれば、アミノ酸配列情報を取得し、取得されたアミノ酸配列情報を構成する各アミノ酸情報に対応する二次構造を取得（例えば、アミノ酸配列情報に対応する立体構造が未知の場合には、既存の二次構造予測プログラムなどを用いて予測して取得、また、対応する立体構造が既知の場合には既存の二次構造判定プログラムなどを用いて取得）し、取得された二次構造およびアミノ酸配列情報を構成するアミノ酸情報に対応するアミノ酸の類似性（例えば、各アミノ酸の疎水性情報等のアミノ酸毎の性質に関する類似性等）に基づいて、アミノ酸配列情報を構成する各アミノ酸情報を、同一の二次構造および同一の類似性を有するアミノ酸情報を同一の情報として識別するための識別情報に変換することにより、アミノ酸配列情報を識別情報からなる識別配列情報に変換し、変換された複数の識別配列情報に対してアライメントを実行するので、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効果的に実行することができ、従来のホモロジー検索法と比べ、より遠縁の配列の探索が可能となる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明によれば、二次構造情報や疎水性情報なども考慮した探索・アライメントを行なうため、従来のアミノ酸の文字のみを考慮したアライメントと比較すると、立体構造も考慮したアライメントが可能となり、より高精度なアライメント作成が可能となる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明のように二次構造情報や疎水性情報などを考慮して作成したアライメントを用いると、疎水エネルギーが安定化し、また二次構造が分断していないモデルを作成でき、モデルとしても真実構造に近いものができると考えられる。

また、本発明によれば、二次構造および類似性に基づいて同一の二次構造および同一の類似性を有するアミノ酸情報を同一の識別情報に分類するための二次構造・類似性分類情報を格納し、変換手段は、二次構造・類似性分類情報に基づいて、アミノ酸配列情報を識別配列情報に変換するので、二次構造・類似性分類情報を予め作成して格納しておき、変換時に当該情報を参照することにより、アミノ酸配列情報の二次構造およびアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したアライメントを効率よく得ることができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明によれば、変換された識別配列情報を格納し、変換手段は、格納された識別配列情報から、アミノ酸配列情報に対応する識別配列情報を検索するので、変換された識別配列情報をデータベースなどに予め格納しておき、変換時に当該データベースなどを参照することにより、変換処理を効率化、高速化することができるようになる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、これにより、例えば公共のデータベース（例えばPDBなど）に登録されたアミノ酸配列情報の二次構造を既存の二次構造判定プログラムなどを用いて判定した後、当該アミノ酸配列情報を本発明により二次構造・類似性分類情報に基づいて変換して作成した識別配列情報を予めデータベースに格納して利用することができるようになり、取得したアミノ酸配列情報に対応する識別配列情報を効率よく検索することができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明によれば、二次構造の組み合わせに応じた構造スコア値が割り当てられた二次構造置換マトリックス、および/または、識別情報の組み合わせに応じた類似性スコア値が割り当てられた類似性置換マトリックスを格納し、アライメント実行手段は、格納された二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、置換マトリックスを予め作成して格納しておき、アライメント実行時に当該置換マトリックスを参照して各スコア値に置換することにより、アミノ酸配列情報の二次構造および/またはアミノ酸配列情報を構成する各アミノ酸情報の類似性を加味したスコア値に置換することにより最適なアライメントを

効率よく得ることができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明によれば、置換マトリックス基準アライメント実行手段は、予め定めた係数により重み付けされた二次構造置換マトリックスおよび/または類似性置換マトリックスに基づいて、識別配列情報を置換してアライメントを実行するので、例えば生物学的な知見などに基づいて二次構造置換マトリックスおよび/または類似性置換マトリックスに対して適切な係数を設定して重み付けをすることによって、生物学的な知見などを反映した最適なアライメントを効率よく得ることができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明によれば、アライメントが実行された識別配列情報を構成する識別情報を、対応するアミノ酸情報に再変換するので、二次構造および類似性を加味してアライメントが実行されたアミノ酸配列情報を効率よく得ることができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明によれば、二次構造・類似性分類情報は、二次構造を、ヘリックス構造、シート構造、または、その他の構造の3つに分類し、類似性を、アラニン(Ala)とグリシン(Gly)からなる第1の群、アスパラギン酸(Asp)とグルタミン酸(Glu)とアスパラギン(Asn)とグルタミン(Gln)からなる第2の群、システイン(Cys)からなる第3の群、フェニルアラニン(Phe)とヒスチジン(His)とトリプトファン(Trp)とチロシン(Tyr)からなる第4の群、イソロイシン(Ile)とロイシン(Leu)とメチオニン(Met)とバリン(Val)からなる第5の群、リシン(Lys)とアルギニン(Arg)からなる第6の群、プロリン(Pro)とセリン(Ser)とトレオニン(Thr)からなる第7の群の7つに分類するので、既知の置換マトリックスであるBLOSUM62に基づいてアミノ酸を7つの群に分類することによりアミノ酸情報の類似性を加味し、かつ、二次構造を3つに分類することによりアミノ酸配列情報の二次構造を加味して、アライメントが実行された識別配列情報を効率よく得ることができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

さらに、本発明によれば、アミノ酸の類似性は、アミノ酸の疎水性、親水性、酸性、塩基性、荷電状態のうち少なくとも一つに基づく情報に関する類似性であるので、アミノ酸配列情報の立体構造などに影響するアミノ酸の性質(疎水性、親水性、酸性、塩基性、荷電状態など)の類似性を考慮してアライメントが実行された識別配列情報を効果的に得ることができる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体を提供することができる。

また、本発明をより効果的に用いるためにホモロジー検索において、e値だけでなく、二次構造や疎水性相互作用によりモデルの立体構造をも考慮したモデルの順位付けが重要と考えられる。この全体の流れの中でpred_fastaを用いることにより、アライメント(e値)、立体構造(二次構造・疎水性相互作用)を考慮したスコア付けを行うことが出来、このスコアが大きいほどアライメント的にもモデル構造的にも精度の高いモデルの構築が期待できるといえる。

本発明によれば、PRED_FASTAは従来のホモロジー検索法と比べ、より遠縁の配列の探索が可能である。また、二次構造情報や疎水性情報なども考慮した探索・アライメントを行なうため、従来のアミノ酸の文字のみを考慮したアライメントと比較すると、より高精度なアライメント作成が可能であり、また立体構造も考慮したアライメントが可能になったといえる。

このように二次構造情報や疎水性情報を考慮して作成したアライメントを用いると、疎水エネルギーが安定化し、また二次構造が分断していないモデルを作成でき、モデルとしても真実構造に近いものができると考えられる。

PRED_FASTAは、従来のホモロジー検索プログラムの短所を克服したホモロジー検索、アライメントプログラムであり、画期的なホモロジー検索の最終ツールとして生命科学分野の発展に大きく貢献するものと期待される。

10

20

30

40

50

また、本発明の適用後に、構造がよいと思われるが領域の短いモデル（局所構造）を領域の長いモデルの構造を用いて、その局所構造を生かしつつ伸長させることができ、特に上述した手順（４）においては、遺伝子の完全長モデルを構築することはより効果的な適用法であると考えられる。

【産業上の利用可能性】

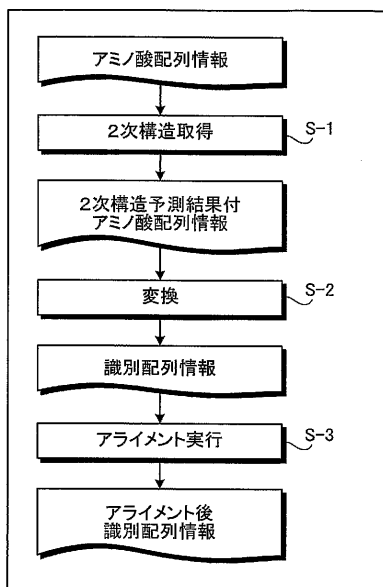
以上のように、本発明にかかる配列情報処理装置、配列情報処理方法、プログラムおよび記録媒体は、アミノ酸配列の二次構造および／またはアミノ酸の類似性を加味したアライメントを得ることができる。

これにより、本発明は、配列情報に対するホモロジー検索やアミノ酸配列のアライメントやタンパク質の立体構造などの解析を行うバイオインフォマティクス分野や生命科学分野において好適に利用することができる。

本発明は、産業上多くの分野、特に、医療、化学、製薬、食品等の分野で広く実施することができ、極めて有用である。

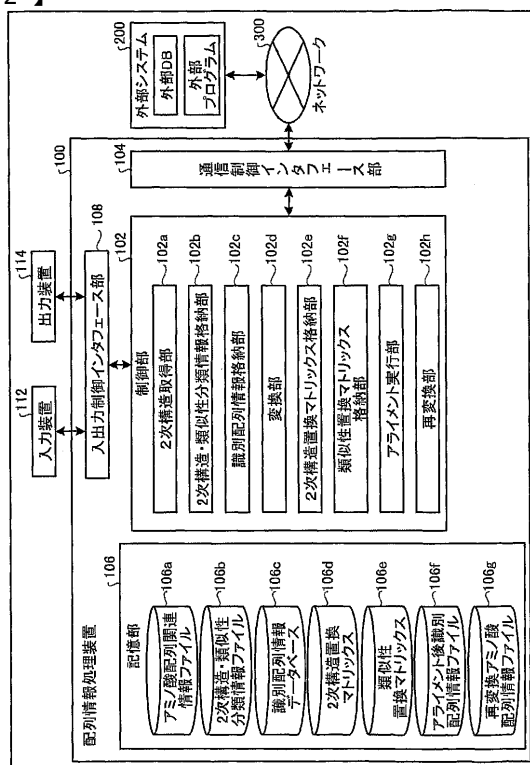
【図 1】

第1図



【図 2】

第2図



【 図 3 】

第3図

アミノ酸配列関連情報ファイル
106a

| アミノ酸配列識別情報 | 情報 | | | | | |
|------------|--------|----------|----------|-----|-----|---------|
| | アミノ酸情報 | Ala | Phe | Trp | ... | Lys |
| A001 | 2次構造情報 | α | α | その他 | ... | β |
| | 類似性情報 | 1 | 4 | 4 | ... | 6 |
| | 識別情報 | A | K | M | ... | S |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

【 図 5 】

第5図

識別配列情報データベース
106c

| アミノ酸配列識別情報 | 情報 | | | | | |
|------------|--------|-----|-----|-----|-----|-----|
| | アミノ酸情報 | Ala | Phe | Trp | ... | Lys |
| K.001 | 識別情報 | A | K | M | ... | S |
| | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

【 図 4 】

第4図

2次構造・類似性分類情報ファイル
106b

| 類似性情報 | アミノ酸情報 | 2次構造情報 | | |
|-------|-----------------|----------|---------|-----|
| | | α | β | その他 |
| 1 | Ala Gly | A | B | C |
| 2 | Asp Glu Asn Gln | D | E | F |
| 3 | Cys | G | H | I |
| 4 | Phe His Trp Tyr | K | L | M |
| 5 | Ile Leu Met Val | N | P | Q |
| 6 | Lys Arg | R | S | T |
| 7 | Pro Ser Thr | V | W | X |

【 図 6 】

第6図

2次構造置換マトリックス
106d

| | α | β | その他 |
|----------|----------|---------|-----|
| α | 6 | -6 | -1 |
| β | -6 | 6 | -1 |
| その他 | -1 | -1 | 6 |

【 図 7 】

第7図

類似性置換マトリックス
106e

| | A | R | N | D | C | Q | E | G | H | I | L | K | M | F | P | S | T | W | Y | V | B | Z | X |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|----|----|----|----|
| A | 1 | 0 | -1 | 0 | 1 | -1 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| R | 0 | 1 | -1 | 0 | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -1 | -1 | 0 | -1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| N | -1 | 1 | 1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | 1 | -1 | -1 | 0 | 0 | 0 | -1 | 0 |
| D | 0 | 0 | -1 | 1 | 0 | -1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C | 1 | 0 | -1 | 0 | 1 | -1 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Q | -1 | 1 | 1 | -1 | 1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | 1 | -1 | -1 | 0 | 0 | 0 | -1 | 0 |
| E | 0 | 0 | -1 | 1 | 0 | -1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | -1 | 0 | -1 | 0 | 0 | -1 | 4 | 4 | 4 | -1 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 |
| H | 0 | -1 | 0 | -1 | 0 | 0 | -1 | 4 | 4 | 4 | -1 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | -1 | 0 | -1 | 0 | 0 | -1 | 4 | 4 | 4 | -1 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 |
| L | -1 | 1 | 0 | -1 | -1 | 0 | -1 | -1 | -1 | -1 | 1 | 1 | 1 | -1 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | -1 |
| K | -1 | -1 | 0 | -1 | -1 | 0 | -1 | -1 | -1 | -1 | 1 | 1 | 1 | -1 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | -1 |
| M | -1 | 1 | 0 | -1 | -1 | 0 | -1 | -1 | -1 | -1 | 1 | 1 | 1 | -1 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | -1 |
| F | 0 | 0 | -1 | 1 | 0 | -1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| P | -1 | -1 | 1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | 1 | -1 | -1 | 0 | 0 | 0 | -1 | 0 |
| S | 0 | 1 | -1 | 0 | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -1 | -1 | 0 | -1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| T | 0 | 1 | -1 | 0 | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -1 | -1 | 0 | -1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Y | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| V | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| B | 1 | 0 | -1 | 0 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Z | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| X | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

【 図 8 】

第8図

アライメント後識別配列情報ファイル
106f

| アミノ酸配列識別情報 | 識別配列情報 | | | | |
|------------|--------|---|---|-----|---|
| A001 | A | B | C | ... | P |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

【 図 9 】

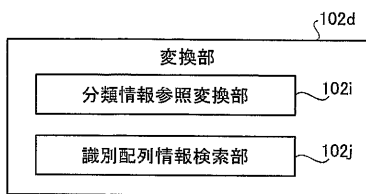
第9図

再変換アミノ酸配列情報ファイル
106g

| アミノ酸配列識別情報 | アミノ酸配列情報 | | | | |
|------------|----------|-----|-----|-----|-----|
| A001 | Ala | Phe | Trp | ... | Lys |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

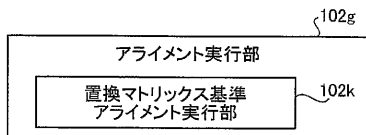
【 図 1 0 】

第10図



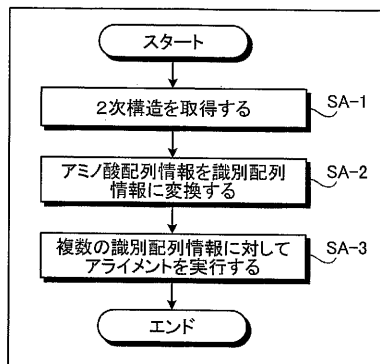
【 図 1 1 】

第11図



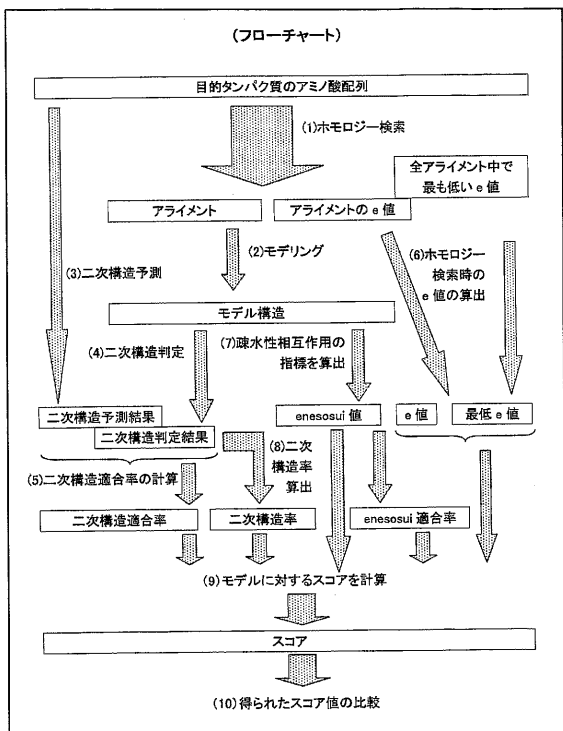
【 図 1 2 】

第12図



【 図 1 3 】

第13図



【 図 1 4 】

第14図

| MA-1 | MA-2 | MA-3 | MA-4 | MA-5 | MA-6 |
|------|------|------|-------|-------|-------|
| 1 | M | C | 0.986 | 0.009 | 0.005 |
| 2 | N | C | 0.921 | 0.043 | 0.021 |
| 3 | E | C | 0.721 | 0.282 | 0.014 |
| 4 | S | H | 0.470 | 0.545 | 0.011 |
| 5 | E | H | 0.065 | 0.928 | 0.008 |

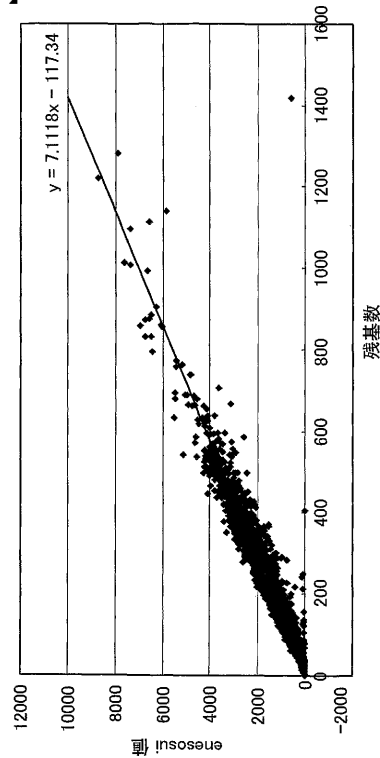
【 図 1 5 】

第15図

| 残基名 | 原子 | | | |
|-----|-----|-----|-----|-----|
| ALA | CB | | | |
| VAL | CG1 | CG2 | | |
| PHE | CB | CD1 | CD2 | CZ |
| PRO | CG | | | |
| MET | CG | CE | | |
| ILE | CG1 | CG2 | CD1 | |
| LEU | CB | CD1 | CD2 | |
| ASP | | | | |
| GLU | CB | | | |
| LYS | CB | CD | | |
| ARG | CG | | | |
| SER | | | | |
| THR | CG2 | | | |
| TYR | CB | CE1 | CE2 | |
| HIS | CB | | | |
| CYS | CB | | | |
| ASN | CB | | | |
| GLN | CB | CG | | |
| TRP | CB | CD2 | CZ2 | CZ3 |
| GLY | | | | |

【 図 1 6 】

第16図



【 図 1 7 】

第17図

```
>test_sequence
MDIRITITSSDYEMVTVSLNEWGGGRLKEKLPRLFEEHFQDTSFITSEHNSMTGFL
IGFQSQSDPETAYIHFSGVHPDFRKMQIGKQLYDVFIETVKQRGCTRVKCVTSPVNK
VSLAYHTKLGFDIEKGTKTVNGISVFANYDGGPQDRVLFVKNI
```

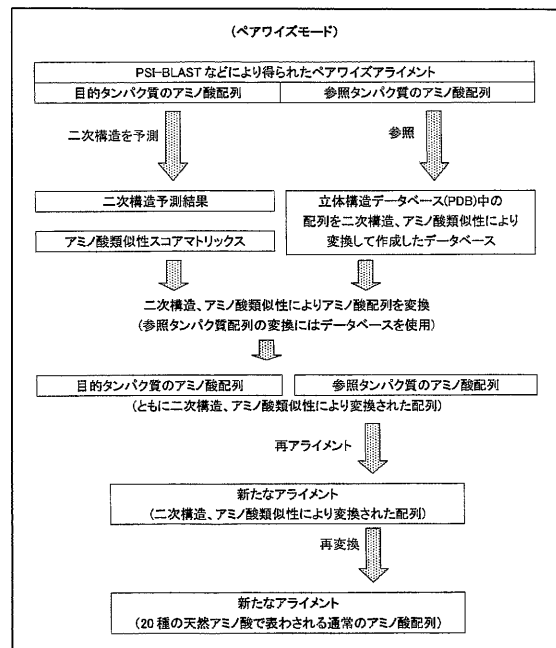
【 図 1 8 】

第18図

| MB-1 | MB-2 | MB-3 | MB-4 | MB-5 | MB-6 | |
|-------------------|--------------|-----------------|--------|---------|--------------|----------|
| 5GCN_A length=166 | score=721.23 | Value=4e-41 | Hom=15 | round=2 | rank=1 | region=1 |
| MB-7 MB-8 | MB-9 | MB-10 MB-11 | MB-12 | MB-13 | MB-14 MB-15 | |
| 13-136 24-162 | PSI | PDB95 -981.000 | 0.88 | -6.67 | 108.131 0.75 | |
| 5GCN_A length=166 | score=719.11 | Value=4e-41 | Hom=15 | round=1 | rank=1 | region=1 |
| 1-156 5-158 | RPS | SCOP95 -897.000 | 0.77 | -6.39 | 112.536 0.72 | |

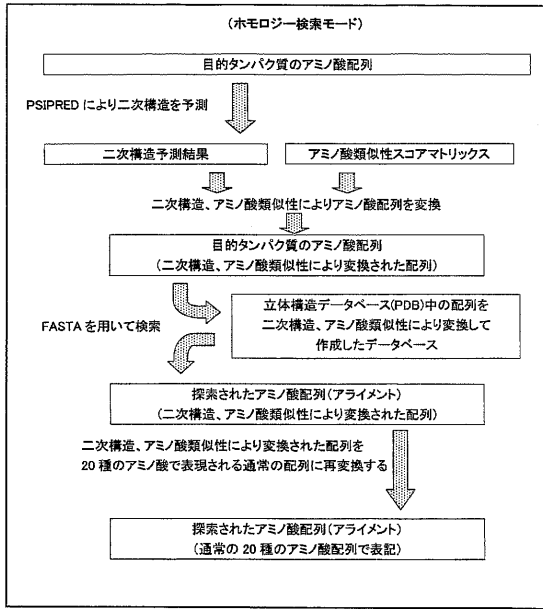
【 図 1 9 】

第19図



【 図 2 0 】

第20図



【 図 2 1 】

第21図

| | グループに所属するアミノ酸 (アミノ酸類似性による分類) | 便宜上付けられるアミノ酸一字表記 | | |
|---|---------------------------------|------------------|-------------|-----|
| | | α ヘリックス | β シート | その他 |
| 1 | Ala Gly | A | B | C |
| 2 | Asp Glu Asn Gln | D | E | F |
| 3 | Cys | G | H | I |
| 4 | Phe His Trp Tyr | K | L | M |
| 5 | Ile Leu Met Val | N | P | Q |
| 6 | Lys Arg | R | S | T |
| 7 | Pro Ser Thr | V | W | X |

【 図 2 2 】

第22図

| | α ヘリックス | β シート | その他 |
|----------------|----------------|-------------|-----|
| α ヘリックス | 6 | -6 | -1 |
| β シート | -6 | 6 | -1 |
| その他 | -1 | -1 | 6 |

【 図 2 3 】

第23図

| | A | R | N | D | C | Q | E | G | H | I | L | K | M | F | P | S | T | W | Y | V | B | Z | X | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|----|----|---|----|
| A | 1 | 0 | -1 | 0 | 1 | -1 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| R | 0 | 1 | -1 | 0 | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -1 | -1 | 0 | -1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| N | -1 | -1 | 1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | 1 | -1 | -1 | 0 | 0 | 0 | -1 | 0 | 0 |
| D | 0 | 0 | -1 | 1 | 0 | -1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C | 1 | 0 | -1 | 0 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Q | -1 | -1 | 1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | 1 | -1 | -1 | 0 | 0 | 0 | -1 | 0 | 0 |
| E | 0 | 0 | -1 | 1 | 0 | -1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | -1 | 0 | -1 | 0 | 0 | -1 | 4 | 4 | 4 | -1 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H | 0 | -1 | 0 | -1 | 0 | 0 | -1 | 4 | 4 | 4 | -1 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | -1 | 0 | -1 | 0 | 0 | -1 | 4 | 4 | 4 | -1 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| L | -1 | -1 | 0 | -1 | -1 | 0 | -1 | -1 | -1 | 1 | 1 | 1 | 1 | 1 | -1 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | -1 |
| K | -1 | -1 | 0 | -1 | -1 | 0 | -1 | -1 | -1 | 1 | 1 | 1 | 1 | 1 | -1 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | -1 |
| M | -1 | -1 | 0 | -1 | -1 | 0 | -1 | -1 | -1 | 1 | 1 | 1 | 1 | 1 | -1 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | 0 | -1 |
| F | 0 | 0 | -1 | 1 | 0 | -1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| P | -1 | -1 | 1 | -1 | -1 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | 1 | -1 | -1 | 0 | 0 | 0 | -1 | 0 | 0 |
| S | 0 | 1 | -1 | 0 | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 0 | -1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| T | 0 | 1 | -1 | 0 | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -1 | -1 | 0 | -1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Y | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| V | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| B | 1 | 0 | -1 | 0 | 1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Z | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| X | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

【 図 2 4 】

第24図

```
>test1
VLSAADKTNVKGVFSGHAEYGAETLERMFTAYPQTCTYFPFDLQHGSAQIK
AHGKKVVAALVEAVNHDDIAGALSKLSDLHAQKLRVDPVNFKFLGHGFLVVAI
HHPSALTAEVHASLDFLCVAVGTVLTAKYR
```

【 図 2 5 】

第25図

```
>>1A4F_A A (141 aa)
initn: 892 initl: 892 opt: 892 Z-score: 134.4 bits: 30.7 E0: 0.065
Smith-Waterman score: 892; 91.489% identity in 141 aa overlap (1-141:1-141)

10 20 30 40 50 60
test1 VLSAADKTNVKGVFSGHAEYGAETLERMFTAYPQTCTYFPFDLQHGSAQIKAHGK
1A4F_A VLSAADKTNVKGVFSGHAEYGAETLERMFTAYPQTCTYFPFDLQHGSAQIKAHGK
10 20 30 40 50 60

70 80 90 100 110 120
test1 KVVAALVEAVNHIDDIAGALSKLSDLHAQKLRVDPVNFKFLGHGFLVVAIHHPSALTAE
1A4F_A KVVAALVEAVNHIDDIAGALSKLSDLHAQKLRVDPVNFKFLGHGFLVVAIHHPSALTAE
70 80 90 100 110 120

130 140
test1 VHASLDFLCVAVGTVLTAKYR
1A4F_A VHASLDFLCVAVGTVLTAKYR
130 140
```

【 図 2 6 】

第26図

```
>CAFASP9058
MANAPKGVKFTSQTTEIVRAKVSSELVEQYDLIGIPIDKEISASYNWKLVEKDY
ASKFHVLDARAVEMIIDWLNIDIKGWLNTNPSREGLDRFVEQYRKKTAALVFTD
TKPLVAGGRLSNLKSFIKSLKFLHLSFLGENGLQFFDKARSADHAHLKLVKFL
KKQLKHSAKLRGAFLLTVFDEQNTNVLVGEFDKLLDHSINIEGSLSPVICHTT
GLNSYVIVLQGE
```

【 図 2 7 】

第27図

```

>>1E0J_F F (288 aa)
initn: 344 initl: 205 opt: 727 Z-score: 128.0 bits: 31.3 E0: 0.14
Smith-Waterman score: 737: 18.79% identity in 266 aa overlap (1-236:19-274)

      10      20      30      40
CAFASP  MANAPKGVKPFSTSQTTTEIIVRAKVSELVEQYDLIIGIPIDKE
1E0J_F  PDGVVSALSLRERIREHLSSEESVGLLPSGCTGINDKTLGA----RGGEVIMVTSBSGM
      10      20      30      40      50

      50      60      70      80      90
CAFASP  ISASYANWKLVEKD--YASKFPHVLDARAVEMIIDWINDIKGWLNTNPYSREGLD----
1E0J_F  GKSIFVRAQALQWGTAMGKKVGLAMLEESVEETAEDLIGLHNVRRLRQSDSKKREITENG
      60      70      80      90     100     110

      100     110     120     130     140     150
CAFASP  RFVQYRKKTAAVLFVDTKPLVAGGRLSNLSKFSIIRKSLKPHLLISFLGKNGKIQPFDKA
1E0J_F  KFDQWFDLFGNDTIFHLVDSFAEATDRLLAKLAYMRSGLGCDVILLDHISIVVSASGES
      120     130     140     150     160     170

      160     170     180     190     200
CAFASP  RSAHDAHKLAVKFLKKQLLKHSAKLRKGAFLITVFDQATNTNLVQE-----F
1E0J_F  DERKMDINDLMTKLGFAKSTGV----VLVVICHLKKNPKGKAHEEGRPYSITDLRGSQA
      180     190     200     210     220     230

      210     220     230
CAFASP  DRLLDHSINIEQSL---LSPVICTH---TGLNSVVIIVLQGE
1E0J_F  LRQLSDTIIALERNNQQGDMPNLVLVRLKCRFTGDTGLAGYMEYNKETGWLPSSSYSG
      240     250     260     270     280

```

【 図 2 8 】

第28図

```

PSI-BALSTIによるアライメント
>CAFASP8880
MSAVTVNDDGLVLRLYQPKASRDSIVGLHG--DEVKVAITAPPVDDGOANSHLVKFLGKQFRVAKSOVVEKGGELGRHKQIKI
>1JRM.A
MDCLREVGDLLVNIIEVSPASGKFGIFPSYNEWKRKRIEVKIHSPPOKGANREIIEFSETFG----RDVEVSGQKSRQKTIPI

(長くアライメントが得られた部分)
>CAFASP8880
TVNDDGLVLRLYQPKASRDSIVGLHGDEVKVAITAPPVDDGOANSHLVKFLGKQFRVAKSOVVEKGGELGRHKQIKI--INPQQIPPEVAALIN
>1JRM.A
LREVGDLLVNIIEVSPASGKFGIFPSYNEWKRKRIEVKIHSPPOKGANREIIEFSETFGGQKSRQKTIPIQGMGRDLFLKLVSEKFG

```

【 図 2 9 】

第29図

```

>CAFASP8880
MSAVTVNDDGLVLRLYQPKASRDSIVGLHG--DEVKVAITAPPVDDGOANSHLVKFLGKQFRVAKSOVVEKGGELGRHKQIKI--INPQQIPPEVAALIN
>1JRM.A
MDCLREVGDLLVNIIEVSPASGKFGIFPSYNEWKRKRIEVKIHSPPOKGANREIIEFSETFG----RDVEVSGQKSRQKTIPIQGMGRDLFLKLVSEKFG

```

【 図 3 0 】

第30図

```

>test3
VLSAADK---TNVKGVFSKISGHAEYGAET
>1A4F.A
VLSAADKTNVKG---VFSKISGHAEYGAET

```

【 図 3 1 】

第31図

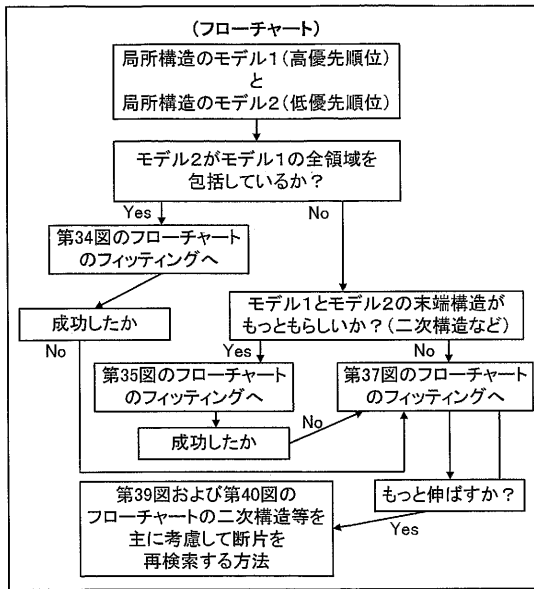
```

>test3
VLSAADKTNVKGVFSKISGHAEYGAET
>1A4F.A
VLSAADKTNVKGVFSKISGHAEYGAET

```

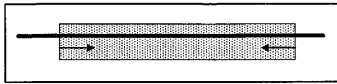
【 図 3 2 】

第32図



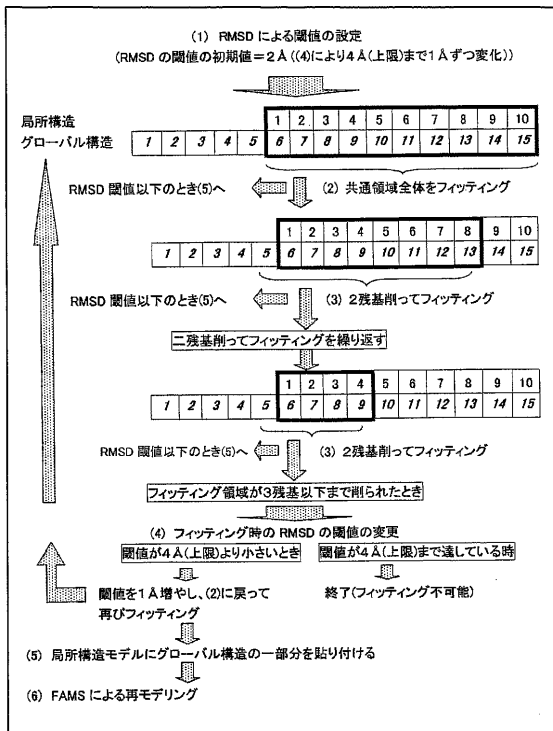
【 図 3 3 】

第33図

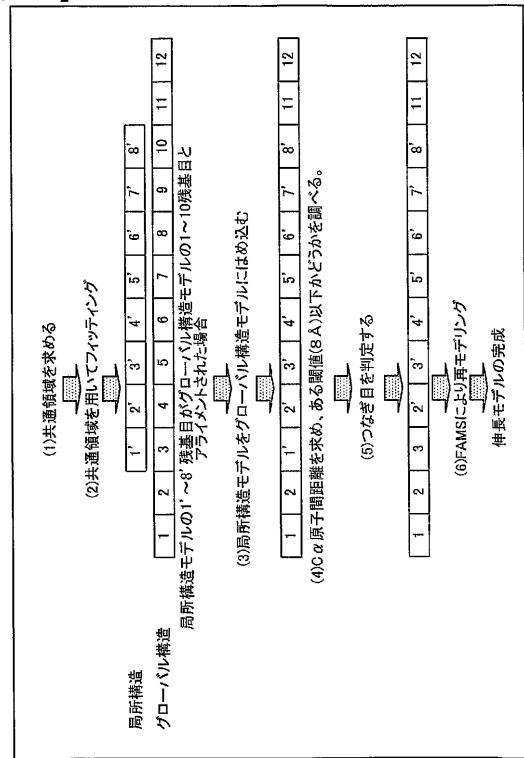


【 図 3 5 】

第35図

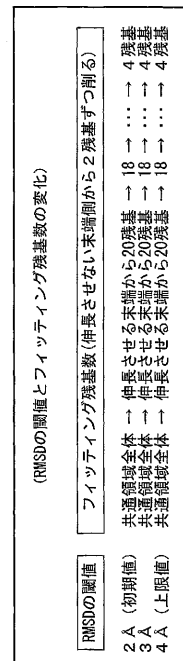


【 図 3 4 】



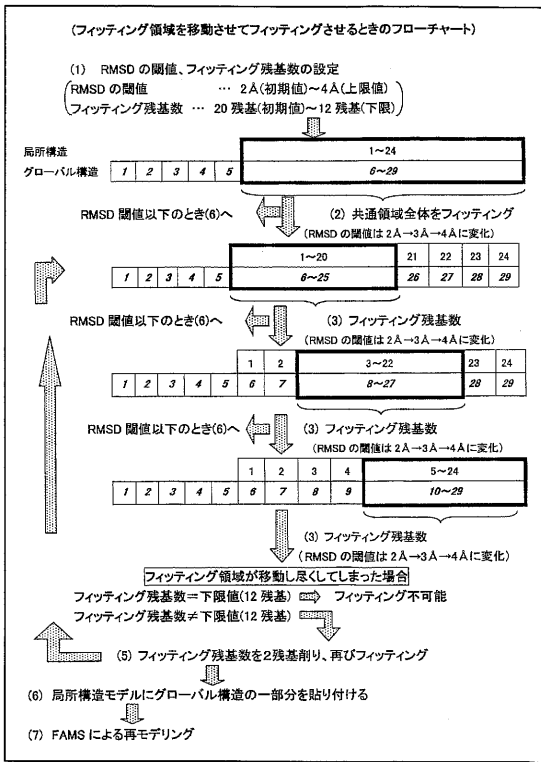
【 図 3 6 】

第36図



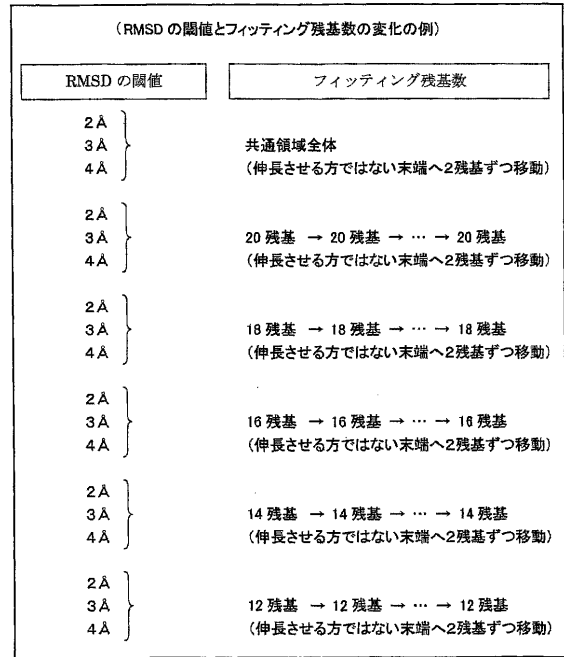
【 図 3 7 】

第37図



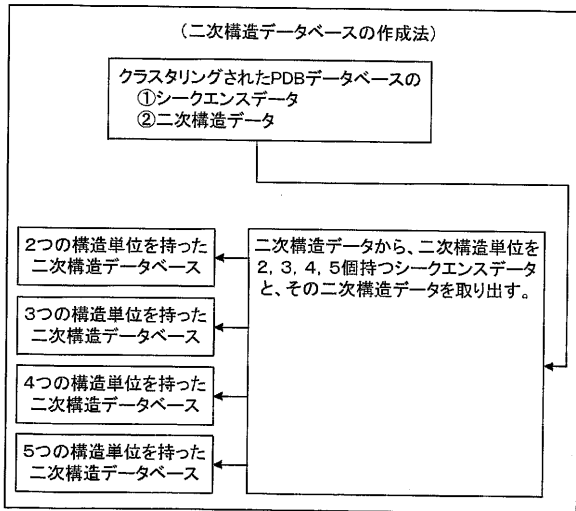
【 図 3 8 】

第38図



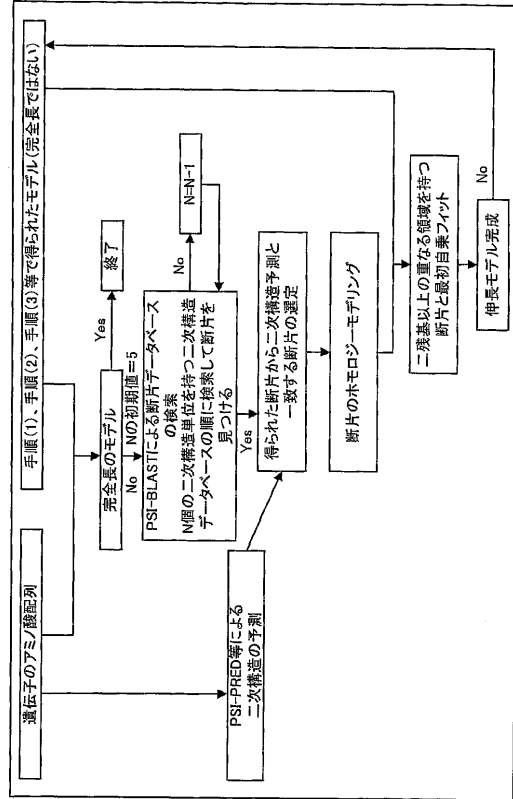
【 図 3 9 】

第39図



【 図 4 0 】

第40図



【 国際調査報告 】

| INTERNATIONAL SEARCH REPORT | | International application No. PCT/JP03/15245 |
|--|---|--|
| A. CLASSIFICATION OF SUBJECT MATTER Int.Cl ⁷ G06F19/00, G06F17/30 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) Int.Cl ⁷ G06F19/00, G06F17/30 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Jitsuyo Shinan Koho 1922-1996 Toroku Jitsuyo Shinan Koho 1994-2004 Kokai Jitsuyo Shinan Koho 1971-2004 Jitsuyo Shinan Toroku Koho 1996-2004 Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) JICST FILE (JOIS), WPI, INSPEC (DIALOG) | | |
| C. DOCUMENTS CONSIDERED TO BE RELEVANT | | |
| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
| X | FISCHEL-GHODSIAN, F. et al., 'Alignment of prote in sequences using secondary structure: a modifi ed dynamic programming method', PROTEIN ENGINEER ING, July 1990, Vol.3, No.7, pages 577 to 581, especially, pages 577 to 578, Methods | 1-25 |
| <input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> 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 document 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 03 February, 2004 (03.02.04) | | Date of mailing of the international search report 17 February, 2004 (17.02.04) |
| Name and mailing address of the ISA/ Japanese Patent Office | | Authorized officer |
| Facsimile No. | | Telephone No. |

| | | | | | | | | | | |
|---|--|---|-----------|------------|-------------|------------|-------------|------------|-------------|------------|
| 国際調査報告 | | 国際出願番号 PCT/JPO3/15245 | | | | | | | | |
| A. 発明の属する分野の分類 (国際特許分類 (IPC)) | | | | | | | | | | |
| Int. Cl ⁷ G06F19/00, G06F17/30 | | | | | | | | | | |
| 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 ⁷ G06F19/00, G06F17/30 | | | | | | | | | | |
| 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 | | | | | | | | | | |
| <table border="0"> <tr> <td>日本国実用新案公報</td> <td>1922-1996年</td> </tr> <tr> <td>日本国公開実用新案公報</td> <td>1971-2004年</td> </tr> <tr> <td>日本国登録実用新案公報</td> <td>1994-2004年</td> </tr> <tr> <td>日本国実用新案登録公報</td> <td>1996-2004年</td> </tr> </table> | | | 日本国実用新案公報 | 1922-1996年 | 日本国公開実用新案公報 | 1971-2004年 | 日本国登録実用新案公報 | 1994-2004年 | 日本国実用新案登録公報 | 1996-2004年 |
| 日本国実用新案公報 | 1922-1996年 | | | | | | | | | |
| 日本国公開実用新案公報 | 1971-2004年 | | | | | | | | | |
| 日本国登録実用新案公報 | 1994-2004年 | | | | | | | | | |
| 日本国実用新案登録公報 | 1996-2004年 | | | | | | | | | |
| 国際調査で使用了電子データベース (データベースの名称、調査に使用了用語) | | | | | | | | | | |
| JICSTファイル (JOIS), WPI, INSPEC (DIALOG) | | | | | | | | | | |
| 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 請求の範囲の番号 | | | | | | | | |
| X | FISCHEL-GHODSIAN, F. et al. 'Alignment of protein sequences using secondary structure: a modified dynamic programming method' PROTEIN ENGINEERING, July 1990, Vol. 3, No7, p. 577-581 especially p. 577-578, Methods | 1-25 | | | | | | | | |
| <input 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 | | の日後に公表された文献 | | | | | | | | |
| 「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 | | 「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 | | | | | | | | |
| 「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) | | 「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 | | | | | | | | |
| 「O」 口頭による開示、使用、展示等に言及する文献 | | 「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 | | | | | | | | |
| 「P」 国際出願日前で、かつ優先権の主張の基礎となる出願 | | 「&」 同一パテントファミリー文献 | | | | | | | | |
| 国際調査を完了した日 | 03.02.2004 | 国際調査報告の発送日 | | | | | | | | |
| | | 17.2.2004 | | | | | | | | |
| 国際調査機関の名称及びあて先 日本国特許庁 (ISA/JP) 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) 高瀬 勲 | 5M 9069 | | | | | | | | |
| | 電話番号 03-3581-1101 内線 3597 | | | | | | | | | |

フロントページの続き

(81) 指定国 AP(BW, GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW

(特許庁注：以下のものは登録商標)

B l u e t o o t h

(注) この公表は、国際事務局(WIPO)により国際公開された公報を基に作成したものである。なおこの公表に係る日本語特許出願(日本語実用新案登録出願)の国際公開の効果は、特許法第184条の10第1項(実用新案法第48条の13第2項)により生ずるものであり、本掲載とは関係ありません。