

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6373977号

(P6373977)

(45) 発行日 平成30年8月15日(2018.8.15)

(24) 登録日 平成30年7月27日(2018.7.27)

(51) Int.Cl. F I  
**G 0 6 F 19/22 (2011.01)** G O 6 F 19/22  
**C 1 2 N 15/00 (2006.01)** C 1 2 N 15/00  
**C 1 2 Q 1/68 (2018.01)** C 1 2 Q 1/68

請求項の数 15 (全 16 頁)

(21) 出願番号	特願2016-514498 (P2016-514498)	(73) 特許権者	590000248
(86) (22) 出願日	平成26年4月30日(2014.4.30)		コーニンクレッカ フィリップス エヌ ヴェ
(65) 公表番号	特表2016-524749 (P2016-524749A)		KONINKLIJKE PHILIPS N. V.
(43) 公表日	平成28年8月18日(2016.8.18)		オランダ国 5656 アーエー アイン ドーフエン ハイテック キャンパス 5
(86) 国際出願番号	PCT/IB2014/061098		High Tech Campus 5, NL-5656 AE Eindhoven
(87) 国際公開番号	W02014/188290		
(87) 国際公開日	平成26年11月27日(2014.11.27)	(74) 代理人	100122769
審査請求日	平成29年4月24日(2017.4.24)		弁理士 笛田 秀仙
(31) 優先権主張番号	61/826,619	(74) 代理人	100163809
(32) 優先日	平成25年5月23日(2013.5.23)		弁理士 五十嵐 貴裕
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 DNA配列の高速かつ安全な検索

(57) 【特許請求の範囲】

【請求項 1】

データベースに記憶されたデオキシリボ核酸(DNA)又はリボ核酸(RNA)配列に対する配列モデルを有する配列指標を生成するステップであって、当該生成するステップは、有限記憶木ソースモデル及び前記有限記憶木ソースモデルに対するパラメータとして前記データベースに記憶された各DNA又はRNA配列に対する前記配列モデルを計算するステップを含み、前記配列モデルが、文脈木重み付け(CTW)を使用して計算される、ステップと、

クエリDNA又はRNA配列に前記配列モデルを適用すること、並びにどれだけ良好に各配列モデルが前記クエリDNA又はRNA配列にフィットするかを決定することに基づいて前記クエリDNA又はRNA配列に最も類似しているものとして前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップと、  
 を含む方法を実行するように電子データ処理装置により実行可能な命令を記憶する非一時的記憶媒体。

【請求項 2】

前記識別するステップが、

有限記憶木ソースモデル及び前記有限記憶木ソースモデルに対するパラメータとして前記クエリDNA又はRNA配列に対するクエリモデルを計算するステップであって、前記クエリモデルが、文脈木重み付け(CTW)を使用して計算される、ステップと、

前記クエリモデルを使用して達成可能な前記クエリDNA又はRNA配列の圧縮の量を

10

20

測定する圧縮計量の基準値を計算するステップと、  
を含み、

前記クエリDNA又はRNA配列に前記配列モデルを適用することが、前記圧縮計量の  
前記基準値と、前記配列モデルを使用して前記クエリDNA又はRNA配列の圧縮率を測  
定する前記圧縮計量の値との間の差に基づいて各配列モデルに対する情報利得を推定する  
ことを含む、

請求項1に記載の非一時的記憶媒体。

【請求項3】

前記識別するステップが、前記配列モデルを使用し、前記データベースに記憶された前  
記DNA又はRNA配列を使用しない、請求項1乃至2のいずれか一項に記載の非一時的  
記憶媒体。

10

【請求項4】

前記クエリDNA又はRNA配列に前記配列モデルを適用することが、  
各配列モデルに対して、前記配列モデルを使用して前記クエリDNA又はRNA配列に  
対する符号語長を計算する、  
ことを含む、請求項1に記載の非一時的記憶媒体。

【請求項5】

前記識別するステップが、  
CTWを使用して有限記憶木ソースモデル及び前記有限記憶木ソースモデルに対するパ  
ラメータとして前記クエリDNA又はRNA配列に対するクエリモデルを計算するステッ  
プと、

20

前記クエリモデルを使用して前記クエリDNA又はRNA配列に対する基準符号語長を  
計算するステップと、  
を含み、

前記クエリDNA又はRNA配列に前記配列モデルを適用することが、前記基準符号語  
長と、前記配列モデルを使用して前記クエリDNA又はRNA配列に対して計算された符  
号語長との間の差に基づいて各配列モデルに対する情報利得を推定することを含む、  
請求項1に記載の非一時的記憶媒体。

【請求項6】

前記データベースに記憶された前記DNA又はRNA配列が、DNA染色体配列であり

30

、  
前記クエリDNA又はRNA配列が、染色体より小さいクエリDNA配列フラグメント  
である、

請求項1乃至5のいずれか一項に記載の非一時的記憶媒体。

【請求項7】

データベースに記憶されたデオキシリボ核酸(DNA)又はリボ核酸(RNA)配列に  
対する文脈木重み付け(CTW)モデル $\{S_x, s_x\}$ を有する配列指標を生成するステップ  
であって、 $S_x$ が前記DNA又はRNA配列 $x$ に対する前記文脈木重み付けモデルを示し  
、 $s_x$ が文脈木モデル $S_x$ のパラメータを示す、当該生成するステップと、

クエリDNA又はRNA配列 $y$ に前記CTWモデル $\{S_x, s_x\}$ を適用すること、並びに  
どれだけ良好に各CTWモデルが前記クエリDNA又はRNA配列 $y$ にフィットするかを  
決定することに基づいて前記クエリDNA又はRNA配列 $y$ に最も類似しているものと  
して前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップと、  
を有し、

40

前記生成するステップ及び前記識別するステップが、電子データ処理装置により実行さ  
れる、方法。

【請求項8】

前記識別するステップが、前記CTWモデル $\{S_x, s_x\}$ を使用し、前記データベースに  
記憶された前記DNA又はRNA配列 $x$ を使用しない、請求項7に記載の方法。

【請求項9】

50

前記識別するステップが、

前記クエリDNA又はRNA配列  $y$  に対するCTWモデル $\{S_y, s_y\}$ を計算するステップであって、 $S_y$ が前記クエリDNA又はRNA配列  $y$  に対する文脈木モデルを示し、 $s_y$ が前記文脈木モデル $S_y$ のパラメータを示す、当該計算するステップと、

前記クエリDNA又はRNA配列  $y$  に対する前記CTWモデル $\{S_y, s_y\}$ を使用して前記クエリDNA又はRNA配列  $y$  の圧縮率を測定する圧縮計量の基準値を計算するステップと、

を含み、

前記クエリDNA又はRNA配列  $y$  に前記CTWモデル $\{S_x, s_x\}$ を適用することが、前記圧縮計量の前記基準値と、前記CTWモデル $\{S_x, s_x\}$ を使用して前記クエリDNA又はRNA配列  $y$  の圧縮率を測定する前記圧縮計量の値との間の差に基づいて各CTWモデル $\{S_x, s_x\}$ に対する情報利得を推定することを含む、

請求項7乃至8のいずれか一項に記載の方法。

#### 【請求項10】

前記識別するステップが、

前記クエリDNA又はRNA配列  $y$  に対するCTWモデル $\{S_y, s_y\}$ を計算するステップであって、 $S_y$ が前記クエリDNA又はRNA配列  $y$  に対する文脈木モデルを示し、 $s_y$ が文脈木モデル $S_y$ のパラメータを示す、当該計算するステップと、

前記クエリDNA又はRNA配列  $y$  に対するCTWモデル $\{S_y, s_y\}$ を使用して前記クエリDNA又はRNA配列  $y$  に対する基準符号語長を計算するステップと、

を含み、

前記クエリDNA又はRNA配列  $y$  に前記CTWモデル $\{S_x, s_x\}$ を適用することが、前記基準符号語長と、前記CTWモデル $\{S_x, s_x\}$ を使用して前記クエリDNA又はRNA配列  $y$  に対して計算される符号語長との間の差に基づいて各CTWモデル $\{S_x, s_x\}$ に対する情報利得を推定することを含む、

請求項7乃至8のいずれか一項に記載の方法。

#### 【請求項11】

前記クエリDNA又はRNA配列  $y$  に前記CTWモデル $\{S_x, s_x\}$ を適用することが、

各CTWモデル $\{S_x, s_x\}$ に対して、前記CTWモデル $\{S_x, s_x\}$ を使用して前記クエリDNA又はRNA配列  $y$  に対する符号語長を計算する、

ことを含み、前記識別するステップが好適には、

前記クエリDNA又はRNA配列  $y$  に最も類似しているものとして、前記CTWモデル $\{S_x, s_x\}$ を使用して、前記クエリDNA又はRNA配列  $y$  に対する最も短い符号語長を持つ前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップ、を含む、

請求項7乃至8のいずれか一項に記載の方法。

#### 【請求項12】

データベースに記憶されたデオキシリボ核酸(DNA)又はリボ核酸(RNA)配列をモデル化する配列指標から文脈木重み付け(CTW)モデル $\{S_x, s_x\}$ を検索するステップであって、 $S_x$ が前記DNA又はRNA配列  $x$  に対する文脈木モデルを示し、 $s_x$ が前記文脈木モデル $S_x$ のパラメータを示す、当該検索するステップと、

クエリDNA又はRNA配列に前記検索されたCTWモデル $\{S_x, s_x\}$ を適用すること、並びにどれだけ良好に各CTWモデルが前記クエリDNA又はRNA配列  $y$  にフィットするかを決定することに基づいて前記クエリDNA又はRNA配列に最も類似しているものとして前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップと、

を含む方法を実行するようにプログラムされた電子データ処理装置、を有する装置。

#### 【請求項13】

前記識別するステップが、前記データベースに記憶された前記DNA又はRNA配列を

使用しない、請求項 1 2 に記載の装置。

【請求項 1 4】

前記クエリ DNA 又は RNA 配列  $y$  に前記検索された C T W モデル  $\{S_x, s_x\}$  を適用することが、

各 C T W モデル  $\{S_x, s_x\}$  に対して、前記 C T W モデル  $\{S_x, s_x\}$  を使用して前記クエリ DNA 又は RNA 配列  $y$  に対する符号語長を計算する、

ことを含む、請求項 1 2 に記載の装置。

【請求項 1 5】

前記識別するステップが、前記識別された 1 以上の DNA 又は RNA 配列をモデル化する前記 C T W モデル  $\{S_x, s_x\}$  を使用して前記クエリ DNA 又は RNA 配列  $y$  に対して計算された最も短い符号語長を持つことに基づいて、前記 DNA 又は RNA 配列  $y$  に最も類似しているものとして、前記データベースに記憶された 1 以上の DNA 又は RNA 配列を識別するステップを含む、請求項 1 4 に記載の装置。

10

【発明の詳細な説明】

【技術分野】

【0001】

以下は、ゲノム配列指標付け (indexing)、記憶、検索 (retrieval)、処理、ラベル付け、及び関連するタスク、並びに患者プライバシー及び医療データセキュリティのような態様並びに医療診断及び医療スクリーニング等のような応用に関する。例示的にデオキシリボ核酸 (DNA) 配列を参照して記載されているが、以下は、DNA 配列、及びリボ核酸 (RNA) 配列等のようなゲノム配列と連動した応用を見つける。

20

【背景技術】

【0002】

DNA シークエンシングは、がん及び他の病気の診断、遺伝性疾患に対する医療スクリーニング、個人用医療、個人用薬物設計、遺伝人類学及び進化研究、系譜的研究、及び法医学人物同定等のような、多くの既存の及び期待される商業的、医療的及び科学的应用を持つ。医療分野において、臨床試験及びゲノムワイド関連研究は、特定の治療、薬物の有効性を評価し、DNA パターンと疾病との間の従属関係等を決定する典型的なツールである。臨床試験において、試験に含める適格性基準は、同様の表現型 (例えば人種) 及び機能性 (例えば遺伝子がオン又はオフである) を持つ DNA 配列を持つ患者を含むことができる。ゲノムワイド関連研究において、試験を行うために、症例群 (例えば突然変異を含む配列) 及び対照群 (突然変位を含まない配列) に分割されることができる DNA 配列が、選択される。遺伝人類学において、ゴールは、一般に、人口移動を追跡する、又は経時的な遺伝的多様性を研究する等のために基準 DNA サンプル (又は基準 DNA サンプルプール) と強い類似性を持つ DNA サンプルを識別することである。これらは、DNA 配列比較を使用する応用の単なる例示的な例である。

30

【0003】

人間の DNA ゲノムは、約 30000 の遺伝子を集合的に暗号化するおおよそ  $3 \times 10^9$  のヌクレオチドからなる。動物、植物及び他の生命体に対するゲノムは、幅広く異なることができるが、典型的には、同等の桁である。臨床試験に対して適格な患者、又は研究目的に対する DNA 配列等を見つけるために、巨大なデータベースが、処理される必要がありうる。したがって、同様な DNA 配列を位置特定する迅速な手順は、有利である。このような検索は、DNA ゲノムの純粋なサイズ並びにギャップ、アライメントエラー、合計配列長の差、及び様々なタイプのノイズを含むことができる実験的に取得された DNA 配列の時々断片的な性質のような多くの問題により複雑にされる。

40

【0004】

人間の DNA に対処する場合、他の検討事項は、対象のプライバシーである。DNA 配列は、遺伝的記録全体を暗号化しており、特定の疾患に対するリスク素因及び祖先情報等のような医療的に又は個人的にセンシティブな情報を明らかにすることができる。DNA 配列は、(一卵性の双生児を例外として) 人間のユニーク識別子でもある。同様の検討事項

50

は、競走馬及び作物等のような商業的に価値のある生命体の非人間ゲノム配列データを処理する際にも生じることができる。このような情報の制御に関する関心は、米国における医療保険会社及び雇用主による個人のDNAから得られた健康情報に基づく差別を禁止することを意図される、2008年の遺伝情報差別禁止法(GINA)により示される。しかしながら、GINAは、生命保険、身体障害保険及び長期ケア保険をカバーしていない。また、DNA配列は、他のタイプの個人医療データと比較してユニークな検討事項を関与させる。人間のゲノムは、全体的に理解されるには程遠く、したがって、DNAから新しい個人的にセンシティブな情報を抽出する新しい技術に対する進行中の可能性が存在する。また、他の医療情報とは異なって、DNA配列は、これら自体が識別子であるので、匿名化されることができない。したがって、DNAマッチングは、好ましくは、データセキュリティを強化する形で行われるべきである。

10

【発明の概要】

【発明が解決しようとする課題】

【0005】

以下は、前述の制限等を克服する改良された装置及び方法を検討する。

【課題を解決するための手段】

【0006】

1つの例示的態様によると、不揮発性記憶媒体は、データベースに記憶されたDNA又はRNA配列に対する配列モデルを有する配列指標を生成するステップであって、有限記憶木ソースモデル及び前記有限記憶木ソースモデルに対するパラメータとして前記データベースに記憶される各DNA又はRNA配列に対する前記配列モデルを計算するステップを含む当該生成するステップと、クエリDNA又はRNA配列に対する前記配列モデルのフィッティングの結果に基づいて前記クエリDNA又はRNA配列に最も類似しているものとして前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップとを含む方法を実行するように電子データ処理装置により実行可能な命令を記憶する。

20

【0007】

他の例示的態様によると、方法は、データベースに記憶されたDNA又はRNA配列に対する文脈木重み付け(CTW、context tree weighting)モデル $\{S_x, s_x\}$ を有する配列指標を生成するステップであって、 $S_x$ は、前記DNA又はRNA配列 $x$ に対する文脈木モデルを示し、 $s_x$ は、文脈木モデル $S_x$ のパラメータを示す、当該生成するステップと、クエリDNA又はRNA配列 $y$ に対するCTWモデル $\{S_x, s_x\}$ のフィッティングに基づいてクエリDNA又はRNA配列 $y$ に最も類似しているものとして前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップとを有する。前記生成するステップ及び前記識別するステップは、電子データ処理装置により適切に実行される。

30

【0008】

他の例示的態様によると、装置は、データベースに記憶されたDNA又はRNA配列をモデル化する配列モデルを配列指標から検索するステップであって、前記データベースに記憶された各DNA又はRNA配列に対する前記検索された配列モデルが、有限記憶木ソースモデル及び前記有限記憶木ソースモデルに対するパラメータを有する、当該検索するステップと、クエリDNA又はRNA配列に対する前記検索された配列モデルのフィッティングに基づいて前記クエリDNA又はRNA配列に最も類似しているものとして前記データベースに記憶された1以上のDNA又はRNA配列を識別するステップとを含む方法を実行するようにプログラムされた電子データ処理装置を有する。

40

【0009】

1つの利点は、ゲノム配列の高速比較を提供することにある。

【0010】

他の利点は、匿名性を維持しながら高速比較を提供する形でゲノム配列に指標付けする指標付け方法を提供することにある。

【0011】

50

他の利点は、指標記録とのクエリゲノム配列の高速比較を容易化するように計算済み有限記憶木ソースモデル及びモデルパラメータを含む前記指標記録を使用してゲノム配列に指標付けする指標付け方法を提供することにある。

【 0 0 1 2 】

多くの追加の利点及び利益は、以下の詳細な記載を読むと当業者に明らかになる。

【 0 0 1 3 】

本発明は、様々なコンポーネント及びコンポーネントの構成並びに様々な処理オペレーション及び処理オペレーションの構成の形を取り得る。図面は、好適な実施例を例示する目的のみであり、本発明を限定すると解釈されるべきではない。

【図面の簡単な説明】

10

【 0 0 1 4 】

【図 1】DNA 配列を記憶及び指標付けするシステムを概略的に示す。

【図 2】クエリ DNA 配列に類似した DNA 配列を識別するように図 1 のシステムにより生成される DNA 配列指標を検索するシステムを概略的に示す。

【図 3】囲みボックスにより示される各クエリ染色体に対する最大相互情報量を持つ、例示的な実際に実行される DNA 検索オペレーションからの相互情報量に対する推定値の表を示す。

【発明を実施するための形態】

【 0 0 1 5 】

ここに開示されるのは、（例えば固定又は可変次数）マルコフモデル又は文脈木重み付け（CTW）モデル（ここで使用される例示的アプローチ）等のような有限記憶木ソースモデルを使用して DNA 配列（又は、より一般的に、ゲノム配列、例えば DNA 配列又は RNA 配列等）を指標付けするアプローチである。前記 DNA 配列に対する指標記録が、構築され、前記モデル及びパラメータを含む。この場合、CTW を使用してクエリ DNA 配列の直接的なモデル化により推定される符号語長と比較される、クエリ DNA 配列に対して同じ有限記憶木モデルを使用して得られる推定符号語長は、前記クエリ及び指標 DNA 配列の類似性を定量的に評価する比較計量として機能する。前記符号語長比較は、例えば、エントロピ又は情報利得（IG）又は同様の手段のような相互情報計量を使用して計算される。

20

【 0 0 1 6 】

30

このアプローチは、前記有限記憶木ソースモデル及びパラメータのみが、プレーンテキストで、すなわち暗号化されずに記憶されるので、DNA 配列がデータベースに記憶される患者のプライバシーを保護する。有限長の部分配列の使用は、結果として生じるモデル及びパラメータが元の DNA 配列より大幅に少ない情報を含むので、患者プライバシーを保証し、前記有限記憶木ソースモデルの出力は、実際に本質的に統計的である。前記指標づけされた DNA 配列（のセット）に対する前記モデル及びそのパラメータは、事前に計算されるので、検索は高速である。開示された類似性計量は、相互情報量が検索基準として使用されるので、編集又は設定距離のような他の軽量より柔軟かつ表現豊かである。ここに開示されるように、相互情報量は、ゲノム配列の時間的構造を探索する順次的なユニバーサル圧縮方法に基づいて適切に推定される。

40

【 0 0 1 7 】

図 1 を参照すると、DNA 配列を記憶及び指標付けする例示的システムが、記載される。（ここで  $x^T$  として示され、上付き文字 T が DNA 配列長を示す）指標付けされるべき DNA 配列 10 は、DNA 配列 10 の代表的有限記憶木ソースモデルを生成するように処理される。この事例において、前記有限記憶木ソースモデルは、CTW 方法を使用して計算される文脈木重み付け（CTW）モデルである。DNA 配列  $x^T$  に適用されるモデル化モジュール 12 の出力 14 は、前記有限記憶木ソースモデル及びそのパラメータである。例示的な CTW モデル化において、前記文脈木モデル（すなわち文脈又は部分配列）は、 $S_x$  として（又はモデル化された DNA 配列  $x^T$  のアイデンティティが明らかである場合に、より単純に S として）示され、前記パラメータは、ここで  $s_x$  として（又はモデル化さ

50

れたDNA配列 $x^T$ のアイデンティティが明らかである場合に、より単純にSとして示される、条件付き確率を有する。好ましくは、記述的注釈が、匿名アノテータ16を介して提供される。患者プライバシーが重要である応用において、前記注釈は、匿名であるべきであるが、DNA配列10のソースの関連する記述を構成すべきであり、例えばデモグラフィック情報、又は臨床情報等により前記ソースを記述する。前記応用が、匿名性を必要としない場合、アノテータ16は、前記注釈に対象識別子を含めてもよい。指標記録フォーマット18は、前記モデル及びパラメータ14並びに前記注釈を含む指標記録を構築し、前記指標記録は、電子健康記録(EHR)、又は学問上の目的で採用されるDNAリポジトリ指標等のような、データベース20に記憶される。

#### 【0018】

前記指標記録は、例えばDNA配列 $x^T$ に対する( $S_x$ ,  $s_x$ )として表されるモデル及びパラメータ14を含む。これは、DNA配列 $x^T$ を表すが、近似的表現であり、DNA配列 $x^T$ が導出された対象を識別するには不十分である。したがって、DNA配列 $x^T$ は、適切に安全なフォーマットで別に記憶される。このために、図1の例示的な実施例において、高度暗号規格(AES暗号)に適合する暗号化アルゴリズムを採用する暗号化モジュール24は、DNA配列10を暗号化する。前記暗号化モジュールは、セキュリティ暗号化を実行し、オプションとして、結合された圧縮/暗号化アルゴリズムにより統合的に又は別のオペレーションのいずれかでロスレス圧縮を実行する。データベース記録フォーマット26は、暗号化された(及びオプションとして圧縮された)DNA配列をフォーマット化し、これを暗号化DNA配列データベース28に記憶する。

#### 【0019】

図1を参照し続けると、前記指標付けシステムは、以下のように適切に物理的に実現される。コンピュータ30又は他の電子データ処理装置(例えばコンピュータ、又はセキュア暗号化伝送プロトコルによりリンクされたインターネットベースのサーバ等)は、データ処理モジュール12、18、24、26を実施するように適切にプログラムされる。匿名アノテータ16は、例えば、EHR又は他のデータベースからデモグラフィック又は他の関連情報を抽出する完全自動化システムとして、様々な形で実施されえ、当該情報の匿名化を適切に、又は人間のオペレータが前記関連情報を入力することを可能にするのにユーザインタフェース(例えば例示的なディスプレイ32及びキーボード34)を採用する半自動化システムとして、実行する。DNA配列指標データベース20は、磁気ディスク、個別ディスクの冗長アレイ(RAID)、又は光ディスク等のような非一時的記憶媒体36上で適切に実施される。同様に、暗号化DNA配列データベース28は、磁気ディスク、個別ディスクの冗長アレイ(RAID)、又は光ディスク等のような非一時的記憶媒体38上で適切に実施される。

#### 【0020】

例示的な図1において、同じコンピュータ30が、指標付けモジュール12、18及びアノテータ16又はその自動化された部分、並びに配列暗号化及び記憶モジュール24、26の両方を実施するのに対し、物理的に離れたデータ記憶媒体36、38が、指標20及びデータベース28をそれぞれ記憶する。このアプローチは、(単一のコンピュータ30が適切に使用されるように)ワークフローブロックとして記憶及び指標付けされるべきDNA配列に対して典型的であり、指標20及びデータベース28を別の媒体上で保持することがセキュリティを強化することができるので、有利であることができる。このアプローチにおいて、DNA配列10に対する指標記録は、データベース28に記憶された暗号化DNA配列記録に対するリンクを記憶する(データベース記録フォーマット26を指標記録フォーマット18に接続し、前記指標記録における包含のために前記リンクを後者に伝えることを示す点線矢印により図1に概略的に示される)。

#### 【0021】

代替的な物理的实施が可能であると理解される。例えば、別々のコンピュータが、それぞれ、指標付けオペレーション12、16、18及び暗号化/記憶オペレーション24、26を実施するのに使用されることができる。加えて又は代わりに、前記暗号化されたD

10

20

30

40

50

N A 配列及び対応する指標記録は、同じ物理的非一時的記憶媒体に記憶されることができ  
る。他の変形例として、前記指標記録の要素として前記暗号化された D N A 配列を含める  
ことにより指標 2 0 及び暗号化 D N A 配列データベース 2 8 を結合することが考えられる  
。これは、A E S 又は他の暗号化プロトコルが十分に安全であると見なされる場合に適切  
でありうる。（いかなる事象においても、復号鍵は、別々に、又は何らかの他の安全な形  
で記憶されるべきである）。

#### 【 0 0 2 2 】

以下に、例示的な C T W モデル化モジュール 1 2 のオペレーションが、更に記載される  
。

#### 【 0 0 2 3 】

前記文脈木重み付け ( C T W ) 方法 ( Willems et al., The Context Tree Weighting M  
ethod: Basic Properties, IEEE transactions on Information theory, 1995 ) は、深度  
が指定された最大深度 D を超過しない全ての木モデルに対応する符号化分布 ( coding dis  
tribution ) を計算する。前記分布は、算術的符号化技術を使用して観測された D N A 配  
列 1 0 を圧縮するのに使用されることができ、これは、結果として小さな冗長性を持つ符  
号語を生じる。実際に、実際の圧縮は、実行される必要がなく、むしろ、ここに開示され  
た技術は、前記 D N A 配列を圧縮するのに前記モデルを使用して得られる圧縮の量を示す  
符号語長を推定する。ソース配列の長さにより除算される符号語長は、エントロピの良好  
な推定値を与える。

#### 【 0 0 2 4 】

D N A 配列構造は、アミノ酸に対して及び後で順次的な形でタンパク質に対して符号化  
するようなものである。 $x^T$  が観測された D N A 配列 1 0 を示すとする。（より一般的に  
は、 $x^T$  は、同じ文脈木モデル及びパラメータにより一緒にモデル化される配列のセット  
を示すことができる）。この場合、C T W は、 $P(x^T)$  を推定するのに使用されることが  
でき、ここで  $x^T$  は、アルファベット  $A = \{ 1, 2, 3, 4 \}$  からの値を持つベクトルと  
して適切に表される。（D N A アルファベットが、典型的には  $\{ A, T, G, C \}$  として  
表され、A がアデニンを示し、T がチミンを示し、G がグアニンを示し、C がシトシンを  
示すのに対し、R N A アルファベットは、典型的には  $\{ A, U, G, C \}$  であり、チミン  
がウラシルを表す U により置き換えられることに注意する。アルファベット  $A = \{ 1, 2, 3, 4 \}$  は、一般性を失うことなしにここで使用される。例えばメチル化のような情報  
を取得するように、4 つより多いシンボルを持つアルファベットを採用することも考えら  
れる。） $x^T$  で、観測された配列  $x^T$  内の位置  $t$  におけるアルファベット A からのシンボル  
を示す。前記 D N A 配列に対する統計モデルは、前記文脈木を構築し、前記 C T W アル  
ゴリズムを使用して分布  $P(x^T)$  を、 $P(x_t | \{x_{t-b}, b \in B\})$  として推定することにより推  
定され、ここで B は、適切な整数のセットである。「文脈」 $\{x_{t-b}, b \in B\}$  は、 $x^T$  の  $|B|$   
の異なる場所から得られたアルファベット A からの値のセットからなる。典型的には、  
B は、（最大深度 D までの） $x^T$  に先行する値のセットとして記される。（前記観測され  
た D N A 配列において実際に生じた）全ての可能な文脈は、確率分布  $P(x_t | \{x_{t-b}, b \in B\})$   
と一緒に、それぞれ、文脈木（モデル）及びパラメータを構成する。

#### 【 0 0 2 5 】

前記 C T W アルゴリズムの出力は、前記文脈木モデル及び条件付き確率  $\{S, s\}$  である  
。所定の D N A 配列に対して、前記 D N A 配列が  $\{S, s\}$  を使用して圧縮された場合に得  
られる圧縮の量は、推定された符号語長 L により特徴づけられることができる。ここに開  
示されるように、前記 C T W 方法は、ツープスアプローチで使用されることもでき、第 1  
のステップにおいて、統計モデル  $\{S, s\}$  が、観測された D N A 配列に対して算出され、  
第 2 のステップにおいて、前記モデルを使用して達成可能な前記 D N A 配列の圧縮の量  
を示す前記符号語長が、推定される。前記推定は、第 1 のパスにおいて得られる  $\{S, s\}$  に  
より提供される固定の条件付き確率に基づき、比較すると、従来の（単一パス）C T W に  
おいて、前記符号語長は、各シンボルが処理されると常に更新されている確率に基づいて  
計算される。ここに更に開示されるように、このツープスアプローチは、1 つの D N A 配

10

20

30

40

50



列（一般に一緒にモデル化された基準又は指標配列のセットでありうる、基準又は指標付けされた配列）に前記第 1 のステップを実行し、次いで、結果として生じるモデルを、第 2 の（クエリ）DNA 配列に対する符号語長を推定するのに使用することにより、2 つの異なる DNA 配列に対する類似性計量を規定するように拡張されることができる。前記モデルは、前記指標付けされた DNA 配列から算出されたので、これは、前記指標付けされた DNA 配列に対する最適に短い符号語長を生成すべきである。他方で、前記モデルが、前記クエリ DNA 配列に適用される場合、前記符号語長は、前記クエリ DNA 配列が前記指標付けされた DNA 配列にどれだけ類似しているかに依存する。これらが類似している場合、前記モデルは、良好に「フィット」し、短い推定符号語長に対応する高い度合の圧縮を提供する。他方で、これらが類似していない場合、フィットが貧弱であり、前記クエリ配列に対する推定符号語長は、最適なモデルに対して得られるものより長い。前記クエリ配列から算出されたモデルに対して得られた符号語長は、適切な基準長さを提供する。例示的な定量的定式化は、以下のとおりである。

【 0 0 2 6 】

観測された DNA 配列  $x^T$  を検討する。 $\{S, \sigma\}$  は、D より大きくない深度の木ソースを記述するモデル（文脈）及びパラメータセット（条件付き確率）であると仮定する。この例において、 $\{S, \sigma\}$  が必ずしも  $x^T$  から算出されないことに注意する。パラメータ  $\{S, \sigma\}$  を持つモデルが、DNA 配列  $x^T$  を圧縮するのに使用される場合、圧縮された配列の長さは、

$$L(x^T | x_{-D}^1, S, \Theta_S) = - \sum_{t=1}^T \log_2 P(x_t | x_{-D}^{t-1}, S, \Theta_S) = - \sum_{t=1}^T \log_2 \theta_{\sigma_{\{x_{-D}^{t-1}\}}}^{x_t} \quad (1)$$

により与えられ、式（ 1 ）において、

$$\sigma_{\{x_{-D}^{t-1}\}}$$

は、S から文脈への

$$x_{-D}^{t-1}$$

のマッピングであり、

$$P(x_t | x_{-D}^{t-1}, S, \Theta_S) = \theta_{\sigma_{\{x_{-D}^{t-1}\}}}^{x_t} \in \Theta$$

は、部分配列

$$\sigma_{\{x_{-D}^{t-1}\}}$$

が  $x^T$  において観測された後に生じるシンボル  $x^T$  の確率である。 $\{S, \sigma\}$  が、 $x^T$  を生成した実際のソースを記述する場合（例えば、上の例において、 $x^T$  が前記指標付けされた DNA 配列である場合）、 $L(x^T | x_{-D}^1, S, \sigma)$  は、最小の符号語長である理想的な符号

語長に対応する。しかしながら、 $\{S_x, s_x\}$ が、何らかの他のソースを記述する場合（例えば、上の例において、 $x^T$ が前記クエリ配列である場合）、 $L(x^T|x_{-D}^1, S_x, s_x)$ は、（少なくとも一般的には）前記モデルが他のDNA配列に対して算出され、観測されたDNA配列 $x^T$ を効果的に記述しないので、前記理想的な符号語長より大幅に大きい。前記CTW方法が、観測された（DNA）配列のモデル及びパラメータを推定するのに使用される場合、結果として生じる符号語長は、前記理想的な符号語長から最小の距離（冗長性）を持つ。

【0027】

類似性計量は、前記符号語長が、どれだけ良好に前記モデルが前記DNA配列にフィットするかを示し、前記DNA配列の符号語長が、式（1）の符号語長推定を使用して推定され、この概念を使用して規定されることができる。 $y^N$ 及び $x^T$ が、必ずしも同じ長さではない2つの観測されたDNA配列であると仮定する。前の例に対する類推において、 $x^T$ が長さTの指標付けされたDNA配列であるとし、 $y^N$ が長さNのクエリDNA配列であるとする。 $\{S_x, s_x\}$ が、前記CTW方法を使用して $x^T$ に対して算出されたモデル及びパラメータセットであるとする。有利には、 $\{S_x, s_x\}$ は、指標付けされたDNA配列 $x^T$  10に対して事前に計算され、図1を参照して記載されるようにDNA指標20に記憶されてもよい。更に、 $L_{ctw}(y^N)$ が、前記CTW方法を使用して推定される（クエリ）DNA配列 $y^N$ に対する符号語長であるとする。換言すると、 $L_{ctw}(y^N)$ は、クエリDNA配列 $y^N$ に対して算出されたモデル $\{S_y, s_y\}$ を使用して得られる符号語長である。したがって、 $L_{ctw}(y^N)$ は、前記CTW方法を使用して $y^N$ に対して取得可能な最適な（すなわち最短の）符号語長である。この場合、差

$$\begin{aligned} & \frac{1}{N}L_{ctw}(y^N) - \frac{1}{N}L(y^N|S_x, \Theta_{S_x}) \\ &= -\frac{1}{N}\sum_{t=1}^N \log_2 P_{ctw}(y_t|y_{-D}^{t-1}) + \frac{1}{N}\sum_{t=1}^N \log_2 P(y_t|y_{-D}^{t-1}, S_x, \Theta_{S_x}) \\ &= -\frac{1}{N}\sum_{t=1}^N \log_2 \frac{P_{ctw}(y_t|y_{-D}^{t-1})}{P(y_t|y_{-D}^{t-1}, S_x, \Theta_{S_x})} \\ &= -\frac{1}{N}\sum_{t=1}^N \log_2 \frac{P_{ctw}(y_t|y_{-D}^{t-1})}{\theta_{S_x, \sigma_{\{y_{-D}^{t-1}\}}}^{y_t}} \end{aligned} \quad (2)$$

が、計算されることができる。式（2）の差は、 $x^T$ の分布が $y^N$ を記述（圧縮）するために $y^N$ の代わりに使用される場合に、どれだけ得られることができるかを示すことができる。利得が高い場合、 $\{S_x, s_x\}$ は、 $y^N$ に良好にフィットするソースを記述し、したがって、我々は、 $y^N$ 及び $x^T$ の両方が同じソースにより生成されることを仮定し、これらが類似していると思えることができる。利得が低い場合、 $\{S_x, s_x\}$ を使用して推定される $y^N$ に対する符号語長は、非常に高い冗長性を持ち、 $\{S_x, s_x\}$ は、 $y^N$ を圧縮する助けにならず、これは、他のタイプの（DNA）配列を生成する他のソースに対応することを意味する。したがって、我々は、 $y^N$ 及び $x^T$ が異なるソースにより生成され、これが類似していないと示すことができる。一般に、利得が高いほど、モデル及びパラメータセット $\{S_x, s_x\}$ が、配列 $y^N$ を、より良好に記述する。したがって、 $\{S_x, s_x\}$ を持つソースが $y^N$ を生成したことは、更にもっともらしい。

【0028】

前記CTW方法を使用して推定されたソースシンボルごとの符号語長は、前記DNAソース配列のエントロピの推定値を与える。したがって、式（2）の類似性計量は、DNA

配列  $y^N$  と DNA 配列  $x^T$  を生成した DNA ソースとの間の相互情報量の推定値でもある。式 (2) により提供される相互情報量の推定値は、過小評価である。これは、相互情報量が真に非負であるので、見られることができる。対照的に、式 (2) は、最適な (最小の) 符号語長である  $L_{ctw}(y^N)$  と、非最適な (したがってより大きい) 符号語長である  $L(y^N | S_x, s_x)$  との間の  $(1/N)$  によりスケーリングされた) 差を取る。後に続くのは、式 (2) が、一般的に、厳密に非負の真の相互情報値より一般的に小さい、負の値を取り上げることができる。式 (2) により与えられる相互情報量の過小評価は、部分的に、第 2 項の符号化冗長性の結果として生じる。前記過小評価は、類似性計量としての式 (2) の有用性を否定しないが、しかしながら、より高い類似性 (すなわちより大きな情報利得) が、式 (2) の類似性計量により出される「より小さい負」値により示される。

10

【0029】

先行する記載の観点から、クエリ DNA 配列  $y^N$  と、モデル及びパラメータセット  $\{S_x, s_x\}$  が事前に計算され、指標データベース 20 に記憶される、指標付けされた DNA 配列  $x^T$  との間の類似性を測定する類似性計量  $I$  は、式 (2) を使用して適切に計算される、又は換言すると  $I(y^N; x^T, \{S_x, s_x\})$  は、式 (2) を使用して適切に推定される。

【0030】

一例として、クエリ DNA 配列  $y^N$  に最も類似している DNA 配列指標 20 内の指標付けされた DNA 配列  $x^T$  を見つける問題を検討する。これは、

$$\max_{P(x^T)} I(Y^N; X^T)$$

20

を見つけることになる。 $\{S_x, s_x\}$  が、 $x^T$  の関数である場合、データ処理不等式、

$$\begin{aligned} \max_{P(x^T)} I(y^N; x^T) &= \max_{P(x^T)} I(y^N; x^T, \{S_x, \Theta_{S_x}\}) \\ &= \max_{P(x^T)} (I(y^N; \{S_x, \Theta_{S_x}\}) + I(y^N; x^T | \{S_x, \Theta_{S_x}\})) \\ &\geq \max_{P(x^T)} I(y^N; \{S_x, \Theta_{S_x}\}), \end{aligned} \quad (3)$$

30

による。 $\{S_x, s_x\}$  が、 $y^N$  を生成したソースにマッチする場合、前記不等式は、等式になる。最も類似している指標付けされた DNA 配列は、 $I(Y^N; \{S_x, s_x\})$  を最大化するものである。

【0031】

ここで図 2 を参照すると、クエリ DNA 配列  $y^N$  に類似している DNA 配列を識別するように図 1 のシステムにより生成された DNA 配列指標 20 を検索するシステムが、記載される。クエリ DNA 配列  $y^N$  が、受け取られる。文脈木重み付け (CTW) モジュール 12 (図 1 の指標付けシステムと併せて既に記載されている) は、クエリ DNA 配列  $y^N$  に対するモデル及びパラメータ  $\{S_y, s_y\}$  を算出するのに使用され (これはツーパスバージョンの CTW の第 1 のパスである)、符号語長推定器モジュール 42 は、 $\{S_y, s_y\}$  を使用して得られた最適な (最小の) 符号語長  $L_{ctw}(y^N)$  を推定するのに式 (1) を使用する。

40

【0032】

各指標付けされた DNA 配列  $x^T$  は、次いで、現在試験下の指標付けされた DNA 配列  $x^T$  に対する指標エントリを検索する検索モジュール 52 を起動することにより開始する、試験ループ 50 の反復により試験される。この指標エントリは、CTW を使用して (すなわち、図 1 を参照して記載された CTW モジュール 12 により)  $x^T$  に対して算出され

50

たモデル及びパラメータセット $\{S_x, s_x\}$ を提供する。オペレーション54において、式(1)は、 $x^T$ に対して算出されたモデル及びパラメータセット $\{S_x, s_x\}$ を使用してモデル化されたクエリ配列 $y^N$ に対して(非最適、及び一般的により大きい)符号語長 $L(y^N | S_x, s_x)$ を推定するのに再び使用される。換言すると、オペレーション54は、ツープラスCTWアルゴリズムの第2のパスを実行するが、 $x^T$ に対して算出されたモデル及びパラメータセット $\{S_x, s_x\}$ を使用する。試験ループ50は、相互情報量の推定値 $(1/N)L_{ctw}(y^N) - (1/N)L(y^N | S_x, s_x)$ を計算することにより終了する。

#### 【0033】

代案として、オペレーション54は、省略されることができ、式(2)の最後の表現が、 $(1/N)L_{ctw}(y^N) - (1/N)L(y^N | S_x, s_x)$ を直接的に計算するのに、代わりに使用されることができる。

#### 【0034】

試験ループ50は、試験下の各指標付けされたDNA配列 $x^T$ に対して繰り返される。(これは、DNA指標20において指標付けされたあらゆるDNA配列であってもよく、又は代わりに、匿名化された注釈に基づいてフィルタリングすることにより生成される前記指標のサブセットであってもよい)。セレクトモジュール60は、次いで、クエリDNA配列 $y^N$ に最も類似している1つ(又はそれ以上)の指標付けされたDNA配列を選択する。これは、例えば式(3)により、単一の最も類似している指標付けされたDNA配列を選択してもよく、又は「上位K」の最も類似している指標付けされたDNA配列が、選択されてもよく(すなわち、最も高い相互情報量を持つKの指標付けされたDNA配列)、「上位K」の最も類似している指標付けされたDNA配列は、相互情報計量により測定される類似性によりランク付けされ、又は閾値が使用されてもよく、例えば相互情報計量が閾値を超過する全ての指標付けされたDNA配列が、選択される、又はその他である。出力モジュール62は、次いで、セレクトモジュール60により選択された前記1以上の最も類似している指標付けされたDNA配列を表示する又は他の形で人間知覚可能形式で提示する。

#### 【0035】

図2の説明的な例において、処理コンポーネント12、42、50、60、62は、処理コンポーネント12、42、50、60、62の機能を実施する適切なソフトウェアにより、指標付けモジュール12、18、24、26を実施する同じコンピュータ30又は他の電子データ処理装置により実施される。代わりに、異なるコンピュータが、それぞれ図1及び2のシステムにより実行される指標付け及び検索オペレーションに対して使用されてもよい。出力モジュール62は、前記選択された指標付けされたDNA配列に関する情報をディスプレイ32上に表示してもよく、又はこの情報を他のコンピュータ(例えば暗号化DNA配列データベース28に対するアクセスを制御するリポジトリコンピュータ)に送信してもよく、又は(プリンタ又は他のマーキングエンジンと連動して)印刷されたレポートを生成してもよく、又はその他であってもよい。これが、データセキュリティ及び対象プライバシーを危険にさらすので、出力モジュール62が、典型的には、実際の指標付けされたDNA配列を実際に符号及び提供しないと理解されるべきである。むしろ、前記出力モジュールは、(クエリDNA配列 $y^N$ に対する類似性に基づいて)関心配列を識別子、実際の配列は、適切なセキュリティ検査処理が実行された後に復号され、認可された個人に提供される。

#### 【0036】

DNA配列指標付けモジュール12、18、24、26及び/又はDNA配列検索モジュール12、42、50、60、62が、指標付けモジュール12、18、24、26及び/又は検索モジュール12、42、50、60、62の機能を実行するようにコンピュータ30により実行可能な命令(すなわちソフトウェア)を符号化する非一時的記憶媒体として実施されうるとも理解されるべきである。前記非一時的記憶媒体は、例えば、ハードディスクドライブ又は他の磁気記憶媒体、ランダムアクセスメモリ(RAM)、読取専用メモリ(ROM)、フラッシュメモリ又は他の電子記憶媒体、光ディスク又は他の光記

10

20

30

40

50

憶媒体、又はこれらの様々な組み合わせ等の1以上を有してもよい。

【0037】

簡潔な総括のために、図1の例示的な指標付けシステムの実施例は、DNA配列(のセット)  $x_i^T, i = 1, 2, \dots, n$  のDNAデータベース28及び対応する匿名化されたDNA配列指標20を作成することを含む指標付けを実行する。これを行うために、モデル及びパラメータ  $\{S_{x_i}, s_{x_i}\}$  は、前記CTW方法を適用することにより各DNA配列(のセット)  $x_i^T, i = 1, 2, \dots, n$  に対して推定され、 $\{S_{x_i}, s_{x_i}\}$  セットは、他の関連情報(すなわち、注釈、オプションとして匿名化される)と一緒に指標データベース20に記憶される。

【0038】

図2の検索プロセスは、クエリ(例) DNA配列  $y^N$  40を与えられる。前記CTWアルゴリズムが、適用され、ソースシンボルごとの符号語長  $(1/N) L_{ctw}(y^N)$  が、モジュール12、42を使用して  $y^N$  に対して推定される。指標データベース20内の各DNA指標記録  $i, i = 1, 2, \dots, n$  に対して、前記符号語長は、 $\{S_{x_i}, s_{x_i}\}$  を仮定して、 $y^N$  内の部分配列を  $S_{x_i}$  からの文脈にマッピングし、対応するパラメータを使用して

$$\frac{1}{N} L(y^N | S_{x_i}, \Theta_{S_{x_i}}) = - \sum_{t=1}^N \log_2 \theta_{S_{x_i}, \sigma_{\{y_{-D}^{t-1}\}}}^{y_t}$$

を計算する(CTW第2パスモジュール54)ことにより  $y^N$  に対して推定される。(  $y^N$  からのある部分配列に対する  $S_{x_i}$  内に文脈が存在しない場合、対応するパラメータは、  $1/2$  のような何らかの適切な値に適切にセットされる。) 情報利得推定値  $(1/N) L_{ctw}(y^N) - (1/N) L(y^N | S_{x_i}, s_{x_i})$  を最大化するDNA配列を指標付けする記録

$\hat{i}$

が、選択され(モジュール60)、前記関連情報が、クエリを行っているパーティに返される(モジュール62)。

【0039】

指標データベース20において、DNA配列(のセット)に対応するモデル及びパラメータセット  $\{S_{x_i}, s_{x_i}\}$  を記憶することのみを必要とすることが理解される。この情報は、実際の配列を生成したソースの確率的特徴のみを提供するので、単独では、前記DNA配列を再構成するのに使用されることができない。

【0040】

図3を参照すると、開示された検索プロセスの説明的な例が、記載される。この例は、GenBankからの14のDNA配列を使用する。ゴールは、染色体ごとにデータベースを構成することである。この例において、前記CTW方法は、各染色体、すなわち本例において染色体1, 2, 3, 5, 8, 9, 10, 14に対して前記モデル及びパラメータセットを推定するのに深度  $D = 9$  (3つのコドンに対応する)を使用する。これらのモデル及びパラメータセットは、前記指標データベースに記憶される。前記クエリDNA配列は、人間のDNA配列フラグメントであり、ゴールは、これがいずれの染色体から来るのかを決定することである。染色体1, 2, 3, 5, 8, 9, 10, 14に対応する前記指標付けされたDNA配列とともに図2の検索システムを使用して、前記クエリDNA配列フラグメントと異なる(指標付けされた)染色体に対応する前記モデル及びパラメータとの間の相互情報計量の推定値が、計算され、前記相互情報計量を最大化する染色体が、返される。図3は、複数のクエリ配列に対するこのような推定値の結果を提示する。図3において観測されるのは、提案された方法が、DNAのクエリピースがいずれの染色体からくるのかを正しく検出したことである。注意すべきは、前記クエリDNAフラグメントが、完全

な染色体ではなく、むしろ、DNA配列長 $N$ のクエリフラグメント $y^N$ が、長さ $T$ の指標付けされた（完全な染色体）DNA配列 $x^T$ の小さな一部であることである。

#### 【0041】

例示的な実施例は、例として意図され、多くの変形例が考えられる。例えば、CTWが、例示的な実施例において採用されているが、様々な有限長マルコフ連鎖モデル又は可変次数マルコフモデルのような、他の有限記憶木ソースモデルが、採用されることができる。一般に、前記アプローチは、（好ましくは暗号化された）データベース28に記憶されたDNA（又はRNA）配列に対する配列モデルを有する配列指標20を生成する。データベース28に記憶された各DNA（又はRNA）配列に対する配列モデルは、有限記憶木ソースモデル及び前記有限記憶木ソースモデルに対するパラメータを有する。説明用の例において、各指標付けされたDNA配列 $x^T$ に対する前記配列モデルは、CTWを使用して $x^T$ から算出されたモデル及びパラメータセット $\{S_{xi}, s_{xi}\}$ である。

10

#### 【0042】

検索フェーズにおいて、データベース28に記憶された1以上のDNA（又はRNA）配列は、クエリDNA（又はRNA）配列40に対する前記配列モデルのフィッティングに基づいて前記クエリDNA（又はRNA）配列に最も類似しているとして識別される。例示的な実施例において、符号語長は、前記クエリDNA配列に対する前記配列モデルのフィッティングを評価するのに使用される。より一般的には、前記有限記憶木ソースモデルを使用して達成可能な前記クエリDNA配列の圧縮の量を測定するいかなる圧縮計量も、モデルフィットを評価するのに使用されることができる。前記圧縮計量が、より高いレベルの圧縮が前記クエリDNA（又はRNA）配列に前記モデルを適用することにより達成可能であることを示す場合に、前記配列モデルは、前記クエリDNA（又はRNA）配列に、より良好にフィットする。

20

#### 【0043】

例示的な類似性（又は比較）計量は、（近似）情報利得（又は、同等に、相互情報量又はエントロピーの変化）表現として定式化される。式（2）は、一例である。しかしながら、これらは、場合により単純化されることができる。例えば、 $N$ による正規化は、1つのクエリDNA配列のみが存在する（したがって $N$ が全ての場合において同じである）場合には、式（2）において省略されてもよい。実際に、1つのクエリDNA配列のみが、前記検索において採用されている場合、前記類似性計量は、 $L_{ctw}(y^N)$ 項がこの場合に一定のオフセットであるので、 $L(y^N | S_{xi}, s_{xi})$ 単独で与えられる推定符号語（すなわち圧縮計量）にされることができる。近似情報利得を得るために、前記類似性又は比較計量は、前記クエリDNA（又はRNA）配列から算出された有限記憶木ソースモデルを使用して前記クエリDNA（又はRNA）配列を圧縮するために得られた（CTW符号語長推定値のような）圧縮計量の値（これは説明的な例において $(1/N)L_{ctw}(y^N)$ である）を、前記データベースの前記DNA（又はRNA）配列から算出された前記配列モデルを使用して前記クエリDNA（又はRNA）配列に対して得られた前記比較計量の値（これらは説明的な例において $(1/N)L(y^N | S_{xi}, s_{xi})$ である）と適切に比較する。

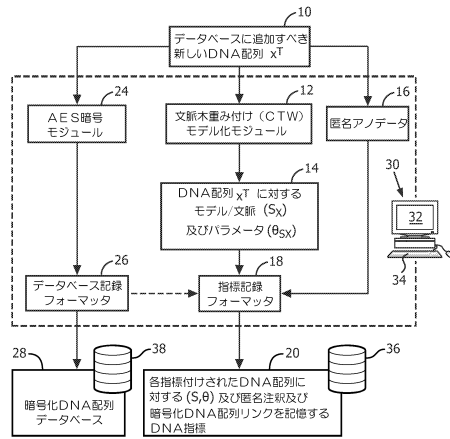
30

#### 【0044】

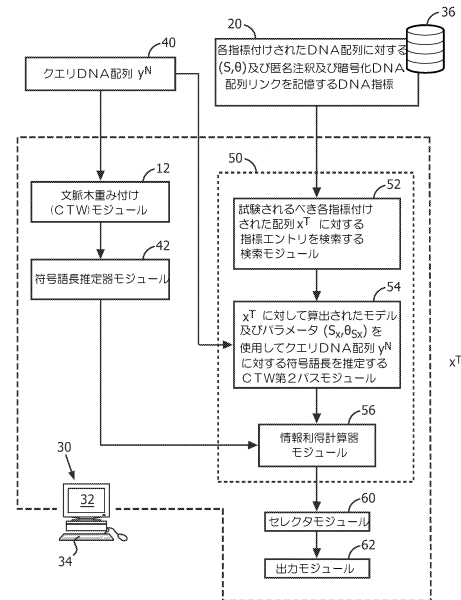
本発明は、好適な実施例を参照して記載されている。明らかに、修正例及び変更例は、先行する詳細な記載を読み、理解すると他者が思いつく。本発明が、添付の請求項又はその同等物の範囲内に入る限り、全てのこのような修正例及び変更例を含むと解釈されるべきである。

40

【 図 1 】



【 図 2 】



【 図 3 】

	染色体に対する相互情報量推定値									
クエリ染色体	1	2	3	5	8	9	10	14		
1	-0.16883	-0.17662	-0.19062	-0.19116	-0.17525	-0.17702	-0.18031	-0.18617		
1	-0.0237	-0.02721	-0.03196	-0.03766	-0.03254	-0.02697	-0.03561	-0.03784		
1	-0.20133	-0.21518	-0.22017	-0.2222	-0.21824	-0.21132	-0.21977	-0.22331		
2	-0.00613	0.085012	-0.00542	-0.00982	-0.00741	-0.00223	-0.00618	-0.00994		
3	-0.02269	-0.01402	0.041464	-0.01767	-0.01713	-0.00675	-0.0218	-0.01881		
5	-0.07684	-0.06272	-0.06162	-0.00854	-0.06846	-0.05452	-0.07161	-0.07257		
5	-0.0971	-0.0648	-0.07114	-0.06463	-0.07918	-0.05804	-0.08502	-0.08603		
8	-0.01266	-0.01229	-0.01528	-0.01913	0.05676	-0.0103	-0.01446	-0.01544		
8	-0.02475	-0.02566	-0.04455	-0.03514	-0.02306	-0.02315	-0.02324	-0.03107		
9	-0.02467	-0.01264	-0.01365	-0.01563	-0.01858	0.07308	-0.02074	-0.02068		
10	-0.04395	-0.02693	-0.03615	-0.03762	-0.03394	-0.02079	-0.00575	-0.04144		
10	-0.04919	-0.03164	-0.03858	-0.04385	-0.0395	-0.02606	-0.00923	-0.04534		
10	-0.04458	-0.02924	-0.0371	-0.04147	-0.03617	-0.02454	-0.0071	-0.0417		
14	-0.05247	-0.05541	-0.05562	-0.05696	-0.05506	-0.05147	-0.05438	0.04525		

---

フロントページの続き

(72)発明者 イグナテンコ ターニャ

オランダ国 5 6 5 6 アーエー アインドーフェン ハイ テック キャンパス ビルディング  
5

審査官 塩田 徳彦

(56)参考文献 米国特許出願公開第2 0 0 4 / 0 0 6 8 3 3 2 ( U S , A 1 )

Z. Dawy et al, Mutual information based distance measures for classification and content recognition with applications to genetics, Communications, 2005. ICC 2005. 2005 IEEE International Conference on, IEEE, 2 0 0 5年 5月16日, pages 820-824, DOI: 10.1109/ICC.2005.1494466,

Kertesz-Farkas A, The Application of Data Compression-Based Distances to Biological Sequences, Springer, Boston, MA, 2 0 0 9年, pp 83-100, DOI: [https://doi.org/10.1007/978-0-387-84816-7\\_4](https://doi.org/10.1007/978-0-387-84816-7_4), Print ISBN: 978-0-387-84815-0, Online ISBN: 978-0-387-84816-7

(58)調査した分野(Int.Cl., DB名)

G 0 6 F 1 9 / 1 0 - 1 9 / 2 8

C 1 2 N 1 5 / 0 0

C 1 2 Q 1 / 6 8