

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6533011号  
(P6533011)

(45) 発行日 令和1年6月19日(2019.6.19)

(24) 登録日 令和1年5月31日(2019.5.31)

(51) Int. Cl.	F I
<b>G 1 6 B 30/00</b> (2019.01)	G O 6 F 19/22
C 1 2 Q 1/6806 (2018.01)	C 1 2 Q 1/6806 Z
C 1 2 N 15/12 (2006.01)	C 1 2 N 15/12

請求項の数 32 (全 19 頁)

(21) 出願番号	特願2018-510056 (P2018-510056)	(73) 特許権者	516001591
(86) (22) 出願日	平成28年8月25日 (2016.8.25)		ナントミクス, エルエルシー
(65) 公表番号	特表2018-533111 (P2018-533111A)		アメリカ合衆国, カリフォルニア州 90
(43) 公表日	平成30年11月8日 (2018.11.8)		232, カルバー シティ, 9920 ジ
(86) 国際出願番号	PCT/US2016/048768		ェファーソン ブールバード
(87) 国際公開番号	W02017/035392	(74) 代理人	100114775
(87) 国際公開日	平成29年3月2日 (2017.3.2)		弁理士 高岡 亮一
審査請求日	平成30年4月26日 (2018.4.26)	(74) 代理人	100121511
(31) 優先権主張番号	62/209,858		弁理士 小田 直
(32) 優先日	平成27年8月25日 (2015.8.25)	(74) 代理人	100202751
(33) 優先権主張国	米国 (US)		弁理士 岩堀 明代
早期審査対象出願		(74) 代理人	100191086
			弁理士 高橋 香元

最終頁に続く

(54) 【発明の名称】 高正確度変異体コールのためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項1】

患者についてHLA型をインシリコ予測する方法であって、  
 複数の既知のおよび異なるHLA対立遺伝子の配列を含む参照配列を提供すること；  
 複数の患者配列リードを提供することであって、前記患者配列リードの少なくともいくつかは患者特異的HLAをコードする配列を含む、複数の患者配列リードを提供すること；  
 前記複数の患者配列リードを複数のk-merのそれぞれのセットへと分解すること；  
 前記参照配列および前記複数のk-merのそれぞれのセットを用いてde Bruijnグラフを作成すること；  
 一つの対立遺伝子についてのそれぞれの投票を全て加えることによって前記一つの対立遺伝子に対する前記複数の患者配列リードのそれぞれの投票から算出される複合マッチスコアを使用して前記既知のおよび異なるHLA対立遺伝子のそれぞれをランク付けることであって、各投票が前記既知のおよび異なるHLA対立遺伝子中の対応するセグメントにマッチするk-merを使用する、前記既知のおよび異なるHLA対立遺伝子のそれぞれをランク付けること；  
 最上位HLA対立遺伝子を前記患者の第1の対立遺伝子のHLA型として識別すること；および  
 調節された最上位HLA対立遺伝子を前記患者の第2の対立遺伝子のHLA型として識別するために、調節された複合マッチスコアを使用して残りの非最上位の既知のおよび異

なる H L A 対立遺伝子を再びランク付けること；  
を含み、

前記調節された複合マッチスコアは、第 1 の H L A 対立遺伝子とマッチする k - m e r の重みを削除するのではなく減少させることによって算出される、方法。

【請求項 2】

前記参照配列が、少なくとも 1 % の対立遺伝子頻度を有する少なくとも 1 つの H L A 型についての対立遺伝子を含む、請求項 1 に記載の方法。

【請求項 3】

前記参照配列が、少なくとも 1 つの H L A 型についての少なくとも 10 の異なる対立遺伝子を含む、請求項 1 に記載の方法。

10

【請求項 4】

前記参照配列が、少なくとも 2 つの異なる H L A 型についての対立遺伝子を含む、請求項 1 に記載の方法。

【請求項 5】

前記 H L A 型が、H L A - A 型、H L A - B 型、H L A - C 型、H L A - D R B - 1 型、および / または H L A - D Q B - 1 型である、請求項 1 に記載の方法。

【請求項 6】

前記複数の患者配列リードが、複数の D N A シーケンシングリードと R N A シーケンシングリードの少なくとも 1 つを含む、請求項 1 に記載の方法。

20

【請求項 7】

前記患者配列リードが、染色体 6 p 2 1 . 3 に位置する、請求項 1 に記載の方法。

【請求項 8】

前記患者配列リードが、次世代シーケンシングリードであり、且つメタデータをさらに含む、請求項 1 に記載の方法。

【請求項 9】

前記患者配列リードが、50 塩基と 250 塩基の間の長さを有する、請求項 1 に記載の方法。

【請求項 10】

前記 k - m e r が、10 ~ 20 の長さを有する、請求項 1 に記載の方法。

30

【請求項 11】

前記 k - m e r が、前記患者配列リードの長さの 5 % と 15 % の間の長さを有する、請求項 1 に記載の方法。

【請求項 12】

前記投票が、患者配列リード当たりの k - m e r の合計数に対するマッチング k - m e r の割合を表す値である、請求項 1 に記載の方法。

【請求項 13】

最上位 H L A 対立遺伝子を前記患者の第 1 の H L A 型として識別するステップをさらに含む、請求項 1 に記載の方法。

【請求項 14】

前記参照配列が、少なくとも 1 % の対立遺伝子頻度を有する少なくとも 1 つの H L A 型についての対立遺伝子を含むか、または前記参照配列が、少なくとも 1 つの H L A 型についての少なくとも 10 の異なる対立遺伝子を含むか、または前記参照配列が、少なくとも 2 つの異なる H L A 型についての対立遺伝子を含む、請求項 1 ~ 13 のいずれか 1 項に記載の方法。

40

【請求項 15】

前記 k - m e r が、10 ~ 20 の長さを有するか、または前記 k - m e r が、患者配列リードの長さの 5 % と 15 % の間の長さを有する、請求項 1 ~ 14 のいずれか 1 項に記載の方法。

【請求項 16】

50

前記複合マッチスコアが、前記複数の患者配列リードからのすべての投票の合計である、および/または前記投票が、患者配列リード当たりの  $k$ -mer の合計数に対するマッチング  $k$ -mer の割合を表す値である、請求項 1 ~ 15 のいずれか 1 項に記載の方法。

【請求項 17】

患者について H L A 型をインシリコ予測するためのコンピュータシステムであって、複数の既知のおよび異なる H L A 対立遺伝子の配列を含む参照配列を格納する参照配列データベースと；

複数の患者配列リードを格納するもしくは提供する患者配列データソースであって、前記患者配列リードの少なくともいくつかは、患者特異的 H L A をコードする配列を含む、患者配列データソースと；

( i ) 前記複数の患者配列リードを複数の  $k$ -mer のそれぞれのセットに分解する；

( i i ) 前記参照配列と前記複数の  $k$ -mer のそれぞれのセットとを使用して  $d e B r u i j n$  グラフを作成する；

( i i i ) 1 つの対立遺伝子についてのそれぞれの投票を全て加えることによって前記 1 つの対立遺伝子に対する前記複数の患者配列リードのそれぞれの投票から算出される複合マッチスコアを使用して前記既知のおよび異なる H L A 対立遺伝子のそれぞれをランク付ける；

( i v ) 最上位 H L A 対立遺伝子を前記患者の第 1 の対立遺伝子の H L A 型として識別する；および

( v ) 調節された最上位 H L A 対立遺伝子を前記患者の第 2 の対立遺伝子の H L A 型として識別するために、調節された複合マッチスコアを使用して残りの非最上位の既知のおよび異なる H L A 対立遺伝子を再びランク付ける；

ようにプログラムされた解析エンジンと、  
を含み、

各投票が、前記既知のおよび異なる H L A 対立遺伝子内の対応するセグメントとマッチする  $k$ -mer を使用し、

前記調節された複合マッチスコアは、第 1 の H L A 対立遺伝子とマッチする  $k$ -mer の重みを削除するのではなく減少させることによって算出される、

コンピュータシステム。

【請求項 18】

前記参照配列が、少なくとも 1 % の対立遺伝子頻度を有する少なくとも 1 つの H L A 型についての対立遺伝子を含むか、または前記参照配列が、少なくとも 1 つの H L A 型についての少なくとも 10 の異なる対立遺伝子を含むか、または前記参照配列が、少なくとも 2 つの異なる H L A 型についての対立遺伝子を含む、請求項 17 に記載のコンピュータシステム。

【請求項 19】

前記 H L A 型が、H L A - A 型、H L A - B 型、H L A - C 型、H L A - D R B - 1 型、および/または H L A - D Q B - 1 型である、請求項 17 に記載のコンピュータシステム。

【請求項 20】

前記複数の患者配列リードが、複数の D N A シーケンシングリードおよび R N A シーケンシングリードの少なくとも 1 つを含む、請求項 17 に記載のコンピュータシステム。

【請求項 21】

前記患者配列リードが、染色体 6 p 2 1 . 3 に位置する、請求項 17 に記載のコンピュータシステム。

【請求項 22】

前記患者配列リードが、次世代シーケンシングリードであり、且つメタデータをさらに含む、または

前記患者配列リードが、50塩基と250塩基の間の長さを有する、

10

20

30

40

50

請求項 17 に記載のコンピュータシステム。

【請求項 23】

前記 k - m e r が、10 ~ 20 の長さを有するか、または前記 k - m e r が、前記患者配列リードの長さの 5 % と 15 % の間の長さを有する、請求項 17 に記載のコンピュータシステム。

【請求項 24】

前記投票が、患者配列リード当たりの k - m e r の合計数に対するマッチング k - m e r の割合を表す値である、請求項 17 に記載のコンピュータシステム。

【請求項 25】

前記解析エンジンが、最上位 H L A 対立遺伝子を前記患者の第 1 の H L A 型として識別するようにさらにプログラムされている、請求項 17 に記載のコンピュータシステム。

10

【請求項 26】

参照配列データベースおよび患者配列データソースが解析エンジンに情報的に連結されているコンピュータシステムに、

複数の既知のおよび異なる H L A 対立遺伝子の配列を含む参照配列を前記参照配列データベースから前記解析エンジンに提供するステップと；

複数の患者配列リードを患者配列データソースから前記解析エンジンに提供するステップであって、前記患者配列リードの少なくともいくつかは、患者特異的 H L A をコードする配列を含む、ステップと；

前記解析エンジンによって前記複数の患者配列リードを複数の k - m e r のそれぞれのセットへと分解するステップと；

20

前記参照配列および前記複数の k - m e r のそれぞれのセットを使用して d e B r u i j n グラフを前記解析エンジンによって作成するステップと；

1 つの対立遺伝子についてのそれぞれの投票を全て加えることによって前記 1 つの対立遺伝子に対する前記複数の患者配列リードのそれぞれの投票から算出される複合マッチスコアを使用して前記解析エンジンによって前記既知のおよび異なる H L A 対立遺伝子のそれぞれをランク付けるステップであって、各投票が、前記既知のおよび異なる H L A 対立遺伝子内の対応するセグメントとマッチする k - m e r を使用する、ステップと；

最上位 H L A 対立遺伝子を前記患者の第 1 の対立遺伝子の H L A 型として識別するステップと；

30

調節された最上位 H L A 対立遺伝子を前記患者の第 2 の対立遺伝子の H L A 型として識別するために、調節された複合マッチスコアを使用して残りの非最上位の既知のおよび異なる H L A 対立遺伝子を再びランク付けるステップであって、前記調節された複合マッチスコアは、第 1 の H L A 対立遺伝子とマッチする k - m e r の重みを削除するのではなく減少させることによって算出される、ステップと；

を含む方法を実施させるためのプログラム命令を含む非一時的なコンピュータ可読媒体。

【請求項 27】

前記参照配列が、少なくとも 1 % の対立遺伝子頻度を有する少なくとも 1 つの H L A 型についての対立遺伝子を含むか、または前記参照配列が、少なくとも 1 つの H L A 型についての少なくとも 10 の異なる対立遺伝子を含むか、または前記参照配列が、少なくとも 2 つの異なる H L A 型についての対立遺伝子を含む、請求項 26 に記載のコンピュータ可読媒体。

40

【請求項 28】

前記 H L A 型が H L A - A 型、H L A - B 型、H L A - C 型、H L A - D R B - 1 型、および / または H L A - D Q B - 1 型である、請求項 26 に記載のコンピュータ可読媒体。

【請求項 29】

前記複数の患者配列リードが、複数の D N A シーケンシングリードおよび R N A シーケンシングリードの少なくとも 1 つを含む、請求項 26 に記載のコンピュータ可読媒体。

【請求項 30】

50

前記患者配列リードが、染色体6 p 2 1 . 3に位置する、または  
 前記患者配列リードが、次世代シーケンシングリードであり、且つメタデータをさらに  
 含む、または  
 前記患者配列リードが、50塩基と250塩基の間の長さを有する、  
 請求項26に記載のコンピュータ可読媒体。

【請求項31】

前記k-merが、10~20の長さを有するか、または前記k-merが、前記患者  
 配列リードの長さの5%と15%の間の長さを有する、請求項26に記載のコンピュータ  
 可読媒体。

【請求項32】

前記投票が、患者配列リード当たりのk-merの合計数に対するマッチングk-mer  
 の割合を表す値である、請求項26に記載のコンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本出願は、2015年8月25日出願の米国仮出願第62/209,858号に対する  
 優先権を主張する。

【0002】

本発明の分野は、ヌクレオチド配列のインシリコ解析のシステムおよび方法であり、特  
 にSNP、マルチヌクレオチド変異体、インデル、構造変異体、およびHLAタイピング  
 の高正確度コールに関する。

【背景技術】

【0003】

本背景技術の記載は、本発明を理解に役立ち得る情報を含む。本明細書で提供する情報  
 のいずれかが従来技術であるもしくは現在主張している発明に関連していること、または  
 具体的もしくは黙示的に参照されたいずれの刊行物が従来技術であることを認めるもの  
 ではない。

【0004】

本明細書のすべての刊行物および特許出願は、それぞれ個々の刊行物または特許出願が  
 参照により具体的におよび個別に組み入れられた場合と同程度に、参照により組み入れら  
 れる。組み入れられた参考文献における用語の定義または使用が本明細書に示すその用語  
 の定義と一致しないまたは相反する場合、本明細書に示すその用語の定義が適用され、参  
 考文献でのその用語の定義は適用されないものとする。

【0005】

配列リード中の小さい変化に起因して正しく整列されないことが多く、変異体情報の不正  
 確さまたは消失のいずれかを引き起こす、関連が高いゲノム配列セグメントを正確にア  
 ライメントするために、ハイスループットシーケンシングデータについての変異体検出は  
 、ますます重要になってきた。関連が高い配列のアライメントを改善するためにいくつか  
 の試みが行われてきた。例えば、「Platypus」(The Wellcome T  
 rust Centre for Human Genetics)は、ハイスループ  
 ットシーケンシングデータ中の比較的効率的で正確な変異体検出のために設計されたツール  
 である。リードのローカルリアライメントおよびローカルアセンブリーを用いることによ  
 り、Platypusは、数kbまでのSNP、MNP、短いインデル、置換および欠失  
 の検出のための比較的高い感受性と高い特異性を達成する。Platypusは従来のア  
 ライメントシステムとしてより正確であることが多いが、それにもかかわらず種々の問題  
 点が残されている。特に、全ゲノムをカバーするゲノムデータの処理は問題であり、類  
 似度が高い複数の配列が存在する場合、所望の正確度に満たないこともある。同様に、D  
 ISCOVAR(Broad Institute)は、配列を構築し変異体を識別する  
 ための比較的正確なツールである。しかし、DISCOVARは概して大量のデータ量の  
 処理に適していない。

10

20

30

40

50

## 【0006】

別の手法において、Big Genomics Inference Engine (BIGGIE; Bioinformatics, vol. 25, pp. 2078-9, 2009)では、最初にゲノムを複雑性が高い領域と低い領域に分類し、続いてそれに応じて情報資源を割り当てることによって、処理速度が上昇する。そのような手法は計算資源に対する要求を減少させる傾向があるが、複雑性が低い領域で変異が起こる場合、変異体コールはそれほど好ましくないことが多い。加えて、次世代のシーケンシングデータのための既知の変異体コーラーの大部分は、変異体を検出しその信頼度を評価するために、確率的フレームワーク(例えば、Bayesian Statisticsを使用する)を利用する。そのような手法は、通常、十分に機能するが、種々の因子、例えば高度のリード深度、プールサンプル、および混入サンプルまたは不純サンプルは、解析を混乱させる傾向がある。そのような問題を解決するために、VarScan (Genome Res. 2012 22:568-576)は、ヒューリスティック/統計的手法を利用して、リード深度、塩基品質、変異体対立遺伝子頻度、および統計的有意性の所望の閾値を満たす変異体コールを行なう。しかし、そのような手法は、通常、単一リードが及ばないゲノム中のより大きな変化を識別しない。

10

## 【0007】

さらなる既知の方法では、DeBruijnカラークラフは、比較的長いk-mer(例えばkは少なくとも55)と、グラフを暗黙にコードするハッシュテーブルとを使用して(Nat Genet. 2012; 44(2):226-232)シーケンシングデータから作成される。しかし、単離されたSNP、短いインデル(1~100bp)およびSNPとインデル(1~100bp)の小複合体の組み合わせの場合、わずか80%の検出力でヘテロ接合部位を検出し、90%の検出力でホモ接合変異体部位を検出したことを著者らは報告した。さらに、中等度のサイズ(100~1000bp)のインデルと複合体変異体の場合、ヘテロ接合部位とホモ接合部位に対する検出力はそれぞれ50%と75~80%であり、大きな変異体(1~50kb)の場合、わずかな検出力(35%)でホモ接合変異体部位を検出したことを著者らは報告した。したがって、記載のようにDeBruijnカラークラフは、SNPおよびインデルの解析を少なくともある程度まで容易にするが、正確度と検出力は望ましいものより低い。したがって、その手法の主要な強さは複数のゲノムの同時解析にあり、それは参照ゲノムを必要とせずに変異体検出への強力で正確な手法を可能にする。

20

30

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0008】

このように、変異体コールのための多数のシステムおよび方法が当技術分野で既知であるが、正確度の高い変異体コールのため、特にインシリコHLAタイピングに関するシステムおよび方法の改善の必要性が依然として存在する。

## 【課題を解決するための手段】

## 【0009】

本発明の主題は、患者の配列データからの正確度の高い変異体コールのための種々のシステム、方法および装置に関し、特にシーケンシング装置からのDNAおよび/またはRNA配列を使用するHLAタイピングを対象とする。特に好ましい態様において、複数のHLA対立遺伝子を含む患者配列リードおよび参照配列は、DeBruijnグラフ手法で処理される。各患者配列リードは種々の対立遺伝子に対する重み付き投票を提供し、各対立遺伝子に対する投票総数を次いで使用して対立遺伝子のランクを付ける。ランク付けにおける最上の対立遺伝子が第1のHLA型であり、第1のHLA型にマッチするk-merに対するバイアスを伴う残りの対立遺伝子の再ランク付けが次いで第2のHLA型を与える。

40

## 【0010】

本発明の主題の一態様において、本発明者は患者についてHLA型をインシリコ予測す

50

る方法を検討し、方法では、既知のおよび異なるHLA対立遺伝子の複数の配列を含む参照配列を提供し、および複数の患者配列リードを提供し、ここで患者配列リードの少なくともいくつかは患者特異的HLAをコードする配列を含む。さらなるステップにおいて、患者配列リードは複数のk-merのそれぞれのセットへと分解され、参照配列および複数のk-merのそれぞれのセットを使用して複合de Bruijnグラフが次いで作成される。既知のおよび異なるHLA対立遺伝子のそれぞれが、複数の患者配列リードのそれぞれの投票から算出される複合マッチスコアを使用してランク付けられることがさらに考えられ、ここで各投票は、既知のおよび異なるHLA対立遺伝子中の対応するセグメントにマッチするk-merを使用する。

#### 【0011】

10

最も一般的に、参照配列は少なくとも1%の対立遺伝子頻度を有する少なくとも1つのHLA型についての対立遺伝子を含み、または参照配列は少なくとも1つのHLA型について少なくとも10の異なる対立遺伝子、および/もしくは少なくとも2つの異なるHLA型についての対立遺伝子を含む。HLA型に関して、適切なHLA型はHLA-A型、HLA-B型、HLA-C型、HLA-DRB-1型、および/またはHLA-DQB-1型を含むことが考えられる。

#### 【0012】

患者配列リードは、複数のDNAシーケンシングリードおよびRNAシーケンシングリードの少なくとも1つを一般的に含み、染色体6p21.3に一般的に位置する。最も一般的には、患者配列リードは次世代シーケンシングリードであり、メタデータをさらに含み、および/または50塩基と250塩基の間の長さである。k-merに関して、好ましいk-merは10~20の長さであり、および/または患者配列リード長の5%と15%の間の長さであることが考えられる。本発明の主題に限定されないが、複合マッチスコアは複数の患者配列リードからのすべての投票の合計であることが一般に好ましく、ここで投票は一般的に患者配列リード当たりのk-merの合計数に対するマッチングk-merの割合を表す値である。

20

#### 【0013】

したがって、複合マッチスコアを使用して、意図される方法は、患者の第1のHLA型として最上位HLA対立遺伝子を識別するステップを含み得る。所望される場合、調節された複合マッチスコアを使用して残りの非最上位の既知のおよび異なるHLA対立遺伝子を再ランク付けする追加のステップを実行して、患者の第2のHLA型として調節された最上位HLA対立遺伝子を識別し得る。最も一般的に、調節された複合マッチスコアは複数の患者配列リードのそれぞれの調節された投票から算出されてよく、および調節された投票は第1のHLA型にマッチするk-merの重みを切り下げることによって算出されてよい。

30

#### 【0014】

上記を考慮して、本発明者は、したがって、患者についてHLA型をインシリコン予測するためのコンピュータシステムも検討する。異なる観点から見て、本発明者はまた、参照配列データベースおよび患者配列データソースが解析エンジンに情報的に連結されるコンピュータシステムに実行させるプログラム命令を含む非一時的なコンピュータ可読媒体も検討する。適切な参照配列、患者配列リード、HLA型、k-mer、複合マッチスコア、および追加の再ランク付けステップに関して、上記と同じ考慮が適用される。

40

本

#### 【0015】

発明の主題の種々の目的、特徴、態様および利点は、同様の符号は同様の構成成分を表す添付の図面に加えて、以下の発明を実施するための形態からさらに明らかになる。

#### 【図面の簡単な説明】

#### 【0016】

【図1】本発明の主題による1つの例示的な方法の概念図である。

【図2】本発明の主題による1つの例示的なコンピュータシステムの概念図である。

50

## 【発明を実施するための形態】

## 【0017】

本発明者は、既知の配列情報を有する参照配列、および統計解析とヒューリスティック解析と組み合わせて de Bruijn グラフに基づく方法を使用して配列が処理される手法において、種々の密接に関連している配列の高度に正確なアライメントが容易に達成できることを見いだした。各 HLA 型は多数のしばしば極めて類似した対立遺伝子を有するため、および配列が高い類似度を有する場合に従来のアライメント方法は有意な分別能を有することが一般的にできないので、そのような解析は、DNA および / または RNA シーケンシング情報から HLA を決定するために、特に有利である。

## 【0018】

本発明の主題の 1 つの例示的な態様において、染色体 6 p 2 1 . 3 (またはそこで / あるいはその近くで HLA 対立遺伝子が見いだされるいずれかの他の位置) に位置する比較的多数の患者配列リードは、データベースまたはシーケンシング装置によって提供される。最も一般的に、配列リードは約 100 ~ 300 塩基の長さであり、リード品質、アライメント情報、配向、位置などを包含するメタデータを含む。例えば、適切な形式としては、SAM、BAM、FASTA、GAR などが挙げられる。本発明の主題に限定されないが、患者配列リードは少なくとも 5 x、より一般的に少なくとも 10 x、より一般的に少なくとも 20 x、最も一般的に少なくとも 30 x の深度カバレッジを提供することが一般に好ましい。

## 【0019】

患者配列リードに加えて、意図される方法は、複数の既知のおよび異なる HLA 対立遺伝子の複数の配列を含む 1 または複数の参照配列をさらに利用する。例えば、一般的な参照配列は、その HLA 型の複数の HLA 対立遺伝子を有する少なくとも 1 つの HLA 型の配列セグメントを含む合成の (ヒトまたは他の哺乳類対応物に対応しない) 配列であり得る。例えば、適切な参照配列は、HLA - A の少なくとも 50 の異なる対立遺伝子に関する既知のゲノム配列の一群を含む。または、あるいはさらに、参照配列は HLA - A の少なくとも 50 の異なる対立遺伝子に関する既知の RNA 配列の一群も含む。もちろん、以下に詳述するように、参照配列は HLA - A の 50 の対立遺伝子に限定されないが、HLA 型および対立遺伝子の数 / 組成に関して代替の組成を有することもある。最も一般的に、参照配列はコンピュータ可読形式であり、データベースまたは他のデータ記憶装置から提供される。例えば、適切な参照配列形式としては、FASTA、FASTQ、EMBL、GCG、または GenBank 形式が挙げられ、公開データリポジトリ (例えば、IMGT、International Immunogenetics 情報システム、または The Allele Frequency Net Database, EUROSTAM, www.allele-frequencies.net) のデータから直接取得するまたは構築することができる。または、参照配列は、対立遺伝子頻度、対立遺伝子頻度、民族別対立遺伝子分布、一般的なまたはまれな対立遺伝子型などの 1 または複数の所定の基準に基づいて、個々の既知の HLA 対立遺伝子から構築されてもよい。

## 【0020】

参照配列を使用して、患者配列リードは、現在、de Bruijn グラフを通して、最良の適合で対立遺伝子を識別できる。この文脈において、各個人が HLA 型について 2 つの対立遺伝子を保有し、これらの対立遺伝子は極めて類似している、または場合によって同一さえあり得ることに留意する必要がある。そのような高類似度は、従来のアライメントスキームに関して重大な問題を提起する。本発明者は、現在、HLA 対立遺伝子、および極めて密接に関連している対立遺伝子さえ、配列リードを比較的小さい k - mer (一般的に 10 ~ 20 塩基の長さを有する) へと分解することにより、および各患者配列リードが対立遺伝子の配列にマッチするその配列リードの k - mer に基づいてそれぞれの対立遺伝子についての投票 (「定量的リードサポート」) を提供する重み付き投票処理を実行することにより de Bruijn グラフが構築される手法を使用して、解決され得ることを見いだした。対立遺伝子についての累積的に最も高い投票は次いで、最も高い可

10

20

30

40

50

能性で予測されるHLA対立遺伝子を示す。加えて、以下でまた詳細に示すように、対立遺伝子にマッチする各フラグメントも用いて全体のカバレッジおよびその対立遺伝子のカバレッジ深度を算出することが一般に好ましい。

【0021】

同じHLA型についての第2の対立遺伝子の識別に関して、本発明者は、比較的類似した第2の対立遺伝子でさえ、最上位HLA対立遺伝子をさらなる考慮から外し、および残りの対立遺伝子を調節された(「スケールされた」)投票を使用して再ランク付ける、よりヒューリスティックな手法で分離できることを見いだした。より具体的には、最上位対立遺伝子とマッチしたk-merの投票値が再ランク付け投票で減少するように再ランク付けが行われる。そのような調節された投票は、最上位対立遺伝子に類似する遺伝子型の重み付け投票を減少させ(しかし削除しない)、したがって遺伝的により関連の少ない対立遺伝子により重みを置く。同時に、類似の対立遺伝子は、無視されない。ランク付けは、全体のカバレッジおよびカバレッジ深度を考慮に入れることによって、さらに改善される。例えば、第1の再ランク付け対立遺伝子は、実質的に低い全体的なカバレッジおよびカバレッジ深度で第2の再ランク付け対立遺伝子よりも高いスコアになり得る。このような場合、第2の再ランク付け対立遺伝子が正しい対立遺伝子である可能性が高い。それゆえ、最上位の再ランク付け対立遺伝子は、同じHLA型の第2の対立遺伝子である。もちろん、上述のように、再ランク付けは全体的なカバレッジおよびカバレッジ深度を考慮に入れることができ、全体的なカバレッジおよび/またはカバレッジ深度が、ユーザが規定する閾値(例えば、94%未満の全体的なカバレッジ、および/または10x未満のカバレッジ深度)を下回るような、対立遺伝子の不適合をもたらすこともあり得る。加えて、投票としてマッチするk-merを使用することで、特定の投票でのユニークなk-merの識別が可能になり、これはその特定の投票が正しい予測でありそうか、そうでないかのさらなるガイダンスとして役立つ。下記の表1は、de Bruijnグラフ手法と、1000 the Genomes Project (IGSR: The International Genome Sample Resource)からの単一ゲノム(YRI)とを使用する、種々のHLA型(HLA-A、HLA-B、HLA-C、DRB1、DQB1)の対立遺伝子の例示的な予測を示す。

【0022】

10

20

【表 1】

重み付きスコア  
 非重み付きスコア  
 %ユニークなk-mer  
 カバレッジ深度  
 カバーされる割合  
 対立遺伝子

NA19238	HLA-A	23721	23721	1.000	21.9	0.925	A*30:01:01
NA19238	HLA-A	15272	22197	0.269	20.5	0.925	A*36:01
NA19238	HLA-A	3595	9609	0.014	11.9	0.938	A*30:18
NA19238	HLA-A	2164	5211	0.031	9.8	0.970	A*30:53
NA19238	HLA-A	1575	5504	0.087	10.3	0.944	A*30:11:01
NA19238	HLA-B	17523	17523	1.000	16.3	0.921	B*53:01:01
NA19238	HLA-B	15938	16485	0.709	15.3	0.913	B*57:03:01
NA19238	HLA-B	5864	12275	0.000	11.4	0.921	B*58:01:01
NA19238	HLA-B	4830	11329	0.000	12.9	0.926	B*57:01:19
NA19238	HLA-B	1913	9496	0.000	8.8	0.921	B*58:01:07
NA19238	HLA-B	762	4835	0.037	9.1	0.925	B*35:27
NA19238	HLA-B	318	4899	0.000	9.2	0.976	B*58:36
NA19238	HLA-C	28463	28463	1.000	26.2	0.924	C*18:02
NA19238	HLA-C	18111	26189	0.317	24.1	0.924	C*04:01:01:01, C*04:01:01:02, C*04:01:01:03, C*04:01:01:04, C*04:01:01:05, C*04:82
NA19238	HLA-C	4320	14117	0.011	17.5	0.927	C*04:41
NA19238	HLA-C	728	3173	0.042	6.0	0.914	C*04:34
NA19238	DRB1	19990	19990	1.000	25.4	0.916	DRB1*16:02:01
NA19238	DRB1	17110	19954	0.599	25.4	0.914	DRB1*11:01:02
NA19238	DRB1	3119	6310	0.000	11.7	0.909	DRB1*11:97
NA19238	DRB1	2411	2790	0.655	10.9	1.000	DRB1*15:96
NA19238	DRB1	1278	4079	0.065	7.6	0.920	DRB1*11:01:08
NA19238	DRB1	893	2165	0.000	8.5	0.934	DRB1*16:23
NA19238	DQB1	17310	17310	1.000	22.4	0.930	DQB1*06:02:01
NA19238	DQB1	16390	16572	0.895	21.5	0.933	DQB1*05:02:01

10

20

表 1

【 0 0 2 3 】

例示的な解析から容易に分かるように、各型の最上位HLA対立遺伝子は容易に区別され、特に重み付けスコアが観察される場合、同じHLA型において第2のランク付け対立遺伝子は残りの対立遺伝子と実質的に異なっている。HLA型の第1と第2のHLA対立遺伝子の選択もまた、有意に高いカバレッジ深度によって、ある程度のカバレッジまで十分にサポートされている。%ユニークなk-mer（最上位と比較して）もまた、本明細書に示すシステムおよび方法の類似性および識別性の良好な指標を提供することも認識すべきである。

30

【 0 0 2 4 】

もちろん、解析およびHLA予測が上記の特定のHLA型に限定される必要はないが、HLA-E、HLA-F、HLA-G、HLA-H、HLA-J、HLA-K、HLA-L、HLA-V、HLA-DQA1、HLA-DMA、HLA-DMB、HLA-DOA、HLA-DOB、HLA-DPA1、HLA-DPB1、HLA-DRA、HLA-DRB345、HLA-MICA、HLA-MICB、HLA-TAP1、HLA-TAP2、およびさらに新たに発見されるHLA型ならびにそれらの対応する対立遺伝子を含む、すべてのHLA型と対立遺伝子変異体が本明細書で検討されることを認識すべきである。さらに、解析が単一HLA型に限定される必要はないが、複数のHLA型が本明細書の使用に適していることを認識すべきである。したがって、それぞれのHLA型についての対立遺伝子の一群とともに、参照配列は2、3、4、またはより多くのHLA型を含み得る。各HLA型はかなりの数の対立遺伝子を有するので、既知の対立遺伝子のすべてを参照配列に包含する必要はないと考えられる。例えば、参照配列は、特定の閾値を上回る対立遺伝子頻度、例えば、少なくとも0.1%、もしくは少なくとも0.5%、もしくは少なくとも1%、もしくは少なくとも2%、もしくは少なくとも5%の対立遺伝子頻度を有する対立遺伝子を含み得る。したがって、異なる観点から見て、適切な参照配列は、少なくとも1つのHLA型について少なくとも10、もしくは少なくとも30、もしくは少

40

50

なくとも50、もしくは少なくとも100、もしくは少なくとも200もしくは少なくとも500、またはさらに多くの対立遺伝子を含み得る。

【0025】

同様に、患者配列リードの性質および型がかなり変化し得ることを認識すべきである。例えば、検討される患者配列リードはDNA配列とRNA配列を含み、それぞれの配列は当技術分野で既知のすべての方法を使用して取得できる。さらに、そのような配列リードは、データ記憶装置（例えばデータベース）から、またはシーケンシング装置から提供され得る。例えば、DNA配列リードはNGSシーケンシング装置から導き出され、RNA配列はrtPCRシーケンシング装置から導き出され得る。したがって、患者配列リードの長さは、一般的に20塩基超、より一般的に50塩基超、最も一般的に100塩基超であるが、通常は5,000塩基未満、もしくは3,000塩基未満、もしくは1,000塩基未満である。したがって、検討される患者配列リードは、100塩基と500塩基の間または150塩基と1,000塩基の間の長さであり得る。

10

【0026】

計算時間とデータ記憶および/または必要メモリを減らすために、患者配列リードをHLA型遺伝子が位置するゲノム領域にあらかじめ選択しておくことがさらに好ましい。例えば、染色体6p21.3に位置する患者配列リードが特に検討される。同様に、患者配列リードはまた、HLA対立遺伝子座が知られているゲノムに対してありそうな位置を示す1または複数のアノテーションに基づいて選択され得る。代替方法として、アノテーションはまた、HLA対立遺伝子であるという配列の可能性を直接参照することもできる。

20

【0027】

患者配列リードの長さに関係なく、患者配列リードが比較的短い長さのk-merに分解されることが一般に好ましく、特に好ましい長さは一般的に10と30の間である。注目すべきことに、そのような短いk-merの長さは、特にそのようなk-merを含有するフラグメントについての重み付き投票のために、変異体コールにおいてより高度な分解能と正確度を可能にする。したがって、k-mer長は一般的に10~30の間、もしくは15~35の間、もしくは20~40の間である。異なる観点から見て、k-merは、好ましくは60未満の、より好ましくは50未満の、最も好ましくは40未満の、しかし5より長い、より一般的に8より長い、および最も一般的に10より長い長さを有する。例えば、適切なk-merは、したがって、患者配列リードの長さの5%と15%の間の長さである。

30

【0028】

ランク付けおよび複合マッチスコアに関して、最も好ましい態様においてマッチスコアが患者配列リード中に存在するすべてのk-merに基づいて作成され、および各投票（すなわち、マッチング）k-merが同じ投票力を有することに留意する必要がある。その結果、患者配列リードは、参照配列中のそれぞれの対立遺伝子に対して特定の定量的リードサポートを有する。さらに、ほとんどの場合、ゲノム中の各位置は>1のシーケンシング深度を有し、および各患者配列リードは対立遺伝子の全長の一部分だけをカバーするので、各対立遺伝子は複数の患者配列リードから複数の投票を受け取ることができる。最も一般的に、対立遺伝子についての投票のすべては、その対立遺伝子の複合マッチスコアに達するように加えられる。それぞれの対立遺伝子の複合マッチスコアは次いで、ランク付けおよびさらなる解析のために使用される。

40

【0029】

しかし、本発明の主題の別の態様において、複合スコアのスコアリングおよび算出は1または複数の特定の目的を達成するように修正されてもよいことに留意する必要がある。例えば、あるフラグメントのマッチスコアは、マッチングk-merのすべてから算出される必要はないが、k-merの無作為な数または選択だけを計数できる。一方では、完全なマッチに満たないk-mer（例えば14/15マッチング）は、おそらく投票重みが低い投票権を与えられる。同様に、特にメタデータが利用できる場合、投票重みは、k-merに対して減少されてよく、および/またはリード品質が特定の閾値を下回る場合

50

は患者配列リードに対して軽減されてよい。一方では、低いシーケンシング深度が存在する場合、投票は特定のフラグメントに対して多すぎることがあり得る。さらに別の意図される態様において、特にリード深度が比較的高い（例えば、少なくとも15x、もしくは少なくとも20x、もしくは少なくとも30x）場合、同じ位置に対する患者配列リードは投票に基づいて除外され得る、または含まれ得る。したがって、複合マッチスコアは、利用できる投票のすべてに基づいてもよく、または対立遺伝子について利用できる投票の一部分のみに基づいてもよい。

#### 【0030】

ランク付けは累積のマッチスコアに一般的に依存するが、ランク付けは少なくとも1つの因子を使用して補正され得ることも認識すべきである。そのような補正因子としては、カバーされる割合、シーケンシング深度、ユニークなk-merの量、および利用できるフラグメントのメタデータが挙げられる。例えば、投票重みは、対立遺伝子のカバレッジが所定の閾値を下回る（例えば、96%未満、もしくは94%未満、もしくは92%未満など）場合および/またはシーケンシング深度が所定の閾値を下回る（例えば、15x未満、もしくは12x未満、もしくは10x未満など）場合、対立遺伝子について軽減され得る。一方では、投票重みは、例えば、ユニークなk-merのパーセンテージが所定の閾値を上回る（例えば、2%超、もしくは5%超、もしくは10%超）場合、対立遺伝子について増加されてもよい。

#### 【0031】

最上位対立遺伝子は、一般的に所与のHLA型の第1の予測対立遺伝子であり、一方第2のランク付け対立遺伝子は、同じHLA型についての第2の対立遺伝子であり得る。しかし、最上位に続くランクの多くが類似の複合マッチスコアを有する場合（例えば、そのスコアのかなりの部分がk-merの高度に共有されるセットに由来する場合）特に、スコアリングは必要に応じてさらに改善または改良されてもよいことに留意する必要がある。好ましい一例において、スコア改良手法が実行されてよく、それは、最上位k-merとマッチした（完全に、または少なくとも90%、もしくは少なくとも95%、もしくは少なくとも97%、もしくは少なくとも99%の類似度のいずれかで）k-merの重みが補正因子によって軽減される再算出を含む。そのような補正因子は、任意の所定の量によって投票を低減できる。最も一般的に、補正因子は投票を10%、もしくは20~40%、もしくは40~60%、もしくはさらに低減させる。これは最上位対立形質と類似している遺伝子型についての重み付き投票を軽減する効果を有し、異なっている遺伝子型を相対的により重要にする。したがって、第1の対立遺伝子は、すべてのシーケンシングデータからの最高のサポートに基づいて識別され、一方第2の対立遺伝子は、第2の対立遺伝子がデータセット中にサポートを有する（例えば、高いスケールされた重み付き投票および遺伝子型カバレッジ）かどうか、またはゲノムが第1の遺伝子型についてホモ接合性であるか（例えば、高い未処理の重み付き投票、極めて低いスケールされた重み付き投票、適切なカバレッジを有する他の対立遺伝子がない）を決定するために、未処理の重み付き投票、スケールされた重み付き投票の両方と、カバレッジを使用する、よりヒューリスティクスに基づく手法で識別されることを認識すべきである。異なる観点から見て、再ランク付けは、最上位対立遺伝子と類似している対立遺伝子の存在下でも、第2の対立遺伝子のより正確な判別を有利に可能にする。さらに、そのような方法は、ホモ接合HLA型の迅速な識別も可能にする。加えて、そのような方法はハッシュテーブルの使用を必要とせず、配列リードをHLA型へと構築することなく適当なHLA対立遺伝子の識別を可能にすることを認識すべきである。さらに、意図されるシステムおよび方法は、DNAおよび/またはRNAデータの使用も可能にする。

#### 【0032】

意図される方法の代表実施形態を、図1に例示的に示す。ここで、方法100はステップ110を含み、複数の既知のおよび異なるHLA対立遺伝子の配列を含む参照配列が提供される。ステップ120において、複数の患者配列リードが提供され、患者配列の少なくともいくつかは患者特異的HLAをコードする配列を包含し、一方ステップ130に

10

20

30

40

50

において、複数の患者配列リードは複数の  $k$ -mer のそれぞれのセットに分解される（一般的に、各  $k$ -mer は 1 塩基（またはそれほど好ましくないが 2 塩基、もしくは 3 塩基、もしくは 4 塩基）の増分で進む）。ステップ 140 において、de Bruijn グラフは、参照配列と、複数の  $k$ -mer のそれぞれのセットとを使用して作成され、およびステップ 150 において、既知のおよび異なる HLA 対立遺伝子のそれぞれは、複数の患者配列リードのそれぞれの投票から算出される複合マッチスコアを使用してランク付けされ、ここで各投票は、既知のおよび異なる HLA 対立遺伝子中の対応するセグメントとマッチする  $k$ -mer を使用する。

#### 【0033】

そのような方法のための例示的なシステムを図 2 に示す。ここで、システム 200 は、参照配列データベース 202（例えば、複数の既知のおよび異なる HLA 対立遺伝子の配列を含む参照配列を格納するデータベースまたはファイル）を含み、ならびに患者配列データソース 204（例えば、複数の患者配列リードを格納もしくは提供する配列データベースまたはシーケンシング装置であって、患者配列リードの少なくともいくつかは患者特異的 HLA をコードする配列を含む）も含み、ここで両者はネットワーク 206（例えば、LAN、WAN、イーサネット、インターネット）を介して解析エンジン 208 に情報的に連結されており、解析エンジンは、(i) 複数の患者配列リードを複数の  $k$ -mer のそれぞれのセットへと分解する；(ii) 参照配列と  $k$ -mer の複数のそれぞれのセットを使用して複合体 de Bruijn グラフを作成する；および (iii) 複数の患者配列リードのそれぞれの投票から算出される複合マッチスコアを使用して既知のおよび異なる HLA 対立遺伝子のそれぞれのランク付けを行うようにプログラムされており、ここで各投票は既知のおよび異なる HLA 対立遺伝子中の対応するセグメントとマッチする  $k$ -mer を使用する。

#### 【0034】

コンピュータに向けられるいずれかの言語は、サーバ、インターフェイス、システム、データベース、エージェント、ピア、エンジン、コントローラ、または個々にもしくは集合的に作動する他の種類の計算装置を包含する、計算装置の任意の適切な組み合わせを含むように読み取られる必要があることに留意すべきである。計算装置が有形の、非一時的なコンピュータ可読記憶媒体（例えば、ハードドライブ、ソリッドステートドライブ、RAM、フラッシュ、ROM など）に格納されるソフトウェア命令を実行するように構成されるプロセッサを含むことを認識すべきである。ソフトウェア命令は、開示される装置に関して後述のとおり、役割、責任、または他の機能性を提供するように計算装置を好ましく構成する。特に好ましい実施形態において、種々のサーバ、システム、データベース、またはインターフェイスは、おそらく HTTP、HTTPS、AES、公開鍵/秘密鍵交換、ウェブサービス API、既知の金融取引プロトコル、または他の電子情報交換方法に基づき、標準化プロトコルまたはアルゴリズムを使用してデータを交換する。データ交換は好ましくは、パケット交換ネットワーク、インターネット、LAN、WAN、VPN、または他の種類のパケット交換ネットワーク上で行なわれる。

#### 【0035】

さらに、本明細書に提示されるシステムおよび方法は、従来のデータ形式および処理方式と比較して、de Bruijn グラフエレメントの構築およびランク付け（および重み付け）が正確度および速度を大幅に上昇させるので、コンピュータ機能を改善することに留意すべきである。さらに、本発明者によって解決される問題はバイオインフォマティクス分野に特異的であり、オミクス情報のコンピューティングなしでは存在さえしないことを認識されたい。最後に、解析エンジンによって実行されるタスクは、コンピュータシステムの支援なしに人の一生のうちに合理的に遂行され得ないことを認識すべきである。

#### 【0036】

上記から容易にわかるように、意図されるシステムおよび方法は、各 HLA 型に、第 2 位のスコアとして実質的により高く格付け/重み付けされる最上位スコアを提供する。したがって、De Bruijn グラフ型解析に基づき、HLA 型は非常に高い正確度で予

10

20

30

40

50

測され得ることを認識すべきである。さらに、本明細書に提示するシステムおよび方法は、種々の他のタスク、例えば、病原体変異体が参照配列の一部を形成する場合の病原体（例えば、HPVなどのウイルス病原体、マイコバクテリアなどの細菌性病原体、または熱帯熱マラリア原虫などの寄生性病原体）のタイピング、または腫瘍の多様性のタイピングなどにも適していることを認識すべきである。

#### 【0037】

本発明の主題のさらなる態様において、de Bruijnグラフに基づく意図されるシステムおよび方法を利用して、構造変異体を識別し分類することもできる。ここでは、参照および未処理のシーケンシングデータを2つのゲノム領域（例えば、推定上の構造的変異の両側、例えば、bcr-abl融合）から取得し、これを使用してグラフを構築する。バブルが次いで、境界参照エッジがユーザ定義の最小ゲノム距離を超えて分離されるようなまたは境界参照エッジが異なる染色体上に位置するような、可能な構造的変異として識別される。そのような手法はほとんどの場合、疑われる構造的変異についての先験的な位置の知識（参照エッジの位置は、構造的変異が疑われるゲノム中の正確な位置を提供する）を必要とするが、そのような知識は通常、境界での正確な配列の識別に役立たない。現在De Bruijnグラフ手法を使用することで、構造的変異のさらに多くの正確な再構築が可能になり、かつ分岐点近くのまたは分岐点内の何らかの新規の配列に役立つ。そのような方法は構造的変異（例えば、挿入、重複など）が同じ鎖上に位置する場合に機能するだけでなく、グラフの構築が算出された逆相補k-merの使用も含む場合に反転を識別するのにもまた同様により有用であることに留意すべきである。すでに前述したように、そのように識別された構造的変異は、続いてvcf形式または他の適切な形式で報告され得る。

#### 【0038】

例えば、腫瘍からの収集された配列情報はDe Bruijnカラーグラフで表わされ、そこではエッジが、k-merが見いだされる入力ソース（例えば、参照、正常サンプル、および/または腫瘍サンプル、様々な時期または年齢で採取されたサンプル、異なる患者または対象群由来のサンプルなど）を識別する「カラー」を有するk-mer（例えば、k=15）であり、および各エッジが隣接するエッジに連結される。もちろん、配列はDNA配列ならびにRNA配列であってよく、このことは発現された体細胞変異、RNA編集および選択的スプライシング（例えば、DNAとRNAが同じ組織に由来する場合）の識別を有利に可能にすることに留意すべきである。最も一般的には、本発明の主題の好ましい一態様において、ゲノム中にk-mer位置を保管するために第1のグラフが参照配列から構築される。好ましくは、必要とされる特定のタスクに応じて、k-merは3塩基と300塩基との間、より好ましくは10~100塩基の長さを有する。例えば、インデル解析が所望される場合、k-mer長は20~50の間（例えばk=30）であり得る。したがって、別の観点から見て、k-mer長は、配列リードの平均長の5%から15%の間であり得る。一旦第1のグラフが確立されると、ゲノムの所与の領域（マップされていないアンカーリードを含む）に位置する腫瘍のおよび正常な未処理シーケンシングデータからのk-merが加えられる。必要に応じて、そのための最大サポートがユーザ定義の特定の閾値（例えば、k=13の場合、閾値は8である）を下回る弱いエッジをグラフから剪定してリードを除去できる。そのような剪定は、配列予測/アライメントの正確度を一般的に高める。

#### 【0039】

de Bruijnグラフ（k=5）内の2つの隣接するエッジについてのデータ構造の例を後述する。

Edge0. 配列 = ATATC

Edge0. 外向き = [TATCG, TATCC]

Edge0. 内向き = [TATAT]

Edge0. サポート = { '参照': 1, '腫瘍': T0, '正常': N0 }

Edge0. quality\_sum = { 'tumor': TQ0, 'normal': NQ0 }

10

20

30

40

50

```

Edge1.配列 = TATCG
Edge1.外向き = [ATCGG]
Edge1.内向き = [ATATC]
Edge1.support = { 'reference' : 0, 'tumor' : T1, 'normal' : N1}
Edge1.品質_sum = { '腫瘍' : TQ1, '正常' : NQ1}

```

## 【 0 0 4 0 】

この例において、Edge 0 データ構造は、それらの k m e r 配列 T A T C G と T A T C C によって定義される 2 つの外向きエッジを有し、配列の前者は、後の E d g e 1 データ構造中に記載される。E d g e 1 の内向きエッジは E d g e 0 へ戻って連結する。上記のデータ構造に記載されるサポートは、シーケンシングデータ（「腫瘍」または「正常」）または参照ゲノム（「参照」）中にエッジ配列が見られた回数をまとめる。上記のエッジ中のサポートに基づいて、E d g e 0 は参照ゲノム中にサポートを有するが、一方 E d g e 1 に連結される外向きエッジはサポートをもたない。これは、E d g e 1 が非参照変異体の始まりであり得ることを示すが、その接続形態が真の変異体（例えば、S N V に起因する、または参照ゲノム中に存在するエッジに囲まれている小さい挿入 / 欠失に起因する d e B r u i j n グラフ中の「バブル」）または人為的変異体（例えば、ジャンクまたはランダムなシーケンシングデータに起因することもあり得る、参照ゲノム中のエッジに再連結しないグラフ中の「チップ」）と一致するかどうかを決定するために、後続のエッジのさらなる内観が必要である。「腫瘍」および「正常」シーケンシングデータ（例えば T 0、N 0、T 1、および N 1）中のサポートのレベルに応じて、非参照変異体の体細胞のまたは生殖系列の分類が決定され得る。分類の 1 つの単純な方法では、変異体は、T 1 > 0 および N 1 > 0 ならば生殖系列として、T 1 > 0 および N 1 = 0 ならば体細胞として、または T 1 = 0 および N 1 > 0 ならば L O H として分類されるが、ほとんどすべての実際の形では、体細胞または生殖系列の状態は、非参照変異体を記述するパス全体の概略分析（すなわち、非参照パス内の平均 / 最小 / 最大サポートおよびエッジの塩基品質）を介して決定される。

## 【 0 0 4 1 】

さらなるステップでは、そのように構築された複合グラフが次いで、腫瘍および参照が分岐する分岐点について解析される。各分岐について、深さ優先探索を使用して、参照に収束する腫瘍をもたらず腫瘍エッジを介するすべてのユニークなパスを識別し、これは d e B r u i j n グラフ中のバブルとして一般的に示される。ブレッドグラムを用いてループを回避できる。複合グラフが次いで、追加の配列で確立される。ここでは、一配列は、同じ患者のマッチする正常組織を表わすことがあり、そこから 2 つの他の配列、腫瘍 D N A および腫瘍 R N A を取得する。そのような例において、腫瘍 D N A および腫瘍 R N A は、同一である（これは必ずしもいつもそうとは限らない）。分岐点および収束点は、k - m e r を使用する配列情報での相違によって決定される。上述の通り、分岐の領域は、グラフ中で「バブル」を生成する。したがって、別の観点から見て、腫瘍配列は分岐点と再収束点の両方を有し得ることを認識すべきである。また留意すべきであるが、腫瘍 D N A および R N A グラフは互いに同等であってよく、このことは D N A とその対応する転写物の配列同一性を示す。

## 【 0 0 4 2 】

各バブル解の終わりから統計解析を次いで利用して、最も可能性の高いアライメントおよび / または配列を識別できる。最も一般的な実施形態において、配列は単なる未処理配列リードではなくアノテーション付きの S A M または B A M ファイルであるので、統計解析は各リードについてのメタベースに基づくリード特異的パラメータを含み得る。したがって、統計解析は、最大のサポート、k - m e r のマッピング / 塩基品質、マッチした正常でのサポートなどを含み得る。結果として、参照配列を再構築するための参照エッジに沿ったバックトラッキングおよびゲノム中の位置の決定は、一般的にユーザ定義の基準（例えば、最小サポート > X リード、正常での最大サポート < Y リードなど）を満たすグラフにおけるパスについて実行され得ることを認識すべきである。そのように構築された配

10

20

30

40

50

列および/または構造を次いで使用して、特定の変異体を分類できる。好ましくは、変異体分類はvcf形式で提示されるが、他の形式も考えられる。

【実施例】

【0043】

HLA予測を確認するために、3種の独立した既知の患者記録とサンプルを1000 Genome project (NA19238、NA19239およびNA19240) から取得し、次いで上述のようにHLA型を予測した。注目すべきことに、かつ予想外に、上述のようにDe Brujnグラフ方法を使用するHLAの決定および予測は、以下の表2Aおよび2Bに見られるように、HLA-C (NA19238について)、DRB1 (NA19239について) およびHLA-C (NA19240について) を除いてほぼ完全にマッチした。

10

【0044】

【表2】

予測：

NA19238 HLA-A	23721	23721	1.000	21.9	0.925	A*30:01:01
NA19238 HLA-A	15272	22197	0.269	20.5	0.925	A*36:01
NA19238 HLA-B	17523	17523	1.000	16.3	0.921	B*53:01:01
NA19238 HLA-B	15938	16485	0.709	15.3	0.913	B*57:03:01
NA19238 HLA-C	28463	28463	1.000	26.2	0.924	C*18:02
NA19238 HLA-C	18111	26189	0.317	24.1	0.924	C*04:01:01
NA19238 DRB1	19990	19990	1.000	25.4	0.916	DRB1*16:02:01
NA19238 DRB1	17110	19954	0.599	25.4	0.914	DRB1*11:01:02
NA19238 DQB1	17310	17310	1.000	22.4	0.930	DQB1*06:02:01
NA19238 DQB1	16390	16572	0.895	21.5	0.933	DQB1*05:02:01
NA19239 HLA-A	24093	24093	1.000	22.2	0.926	A*02:01:01
NA19239 HLA-A	17596	20701	0.537	19.1	0.927	A*68:02:01
NA19239 HLA-B	21308	21308	1.000	19.8	0.918	B*35:01:01
NA19239 HLA-B	15080	20254	0.286	18.8	0.912	B*52:01:02
NA19239 HLA-C	18529	18529	1.000	17.0	0.920	C*04:01:01
NA19239 HLA-C	17846	18484	0.707	17.0	0.919	C*16:01:01
NA19239 DRB1	26014	26014	1.000	33.1	0.914	DRB1*13:01:01
NA19239 DRB1	16174	24412	0.178	31.0	0.914	DRB1*12:01:01
NA19239 DQB1	18503	18503	1.000	24.0	0.930	DQB1*05:01:01
NA19239 DQB1	13459	13510	0.939	17.5	0.931	DQB1*03:01:01
NA19240 HLA-A	21944	21944	1.000	20.2	0.924	A*30:01:01
NA19240 HLA-A	20059	20512	0.800	18.9	0.929	A*68:02:01
NA19240 HLA-B	18637	18637	1.000	17.3	0.927	B*35:01:01
NA19240 HLA-B	17850	18550	0.682	17.3	0.926	B*57:03:01
NA19240 HLA-C	28054	28054	1.000	25.8	0.923	C*18:02
NA19240 HLA-C	20132	27609	0.390	25.3	0.919	C*04:01:01
NA19240 DRB1	22869	22869	1.000	29.1	0.917	DRB1*16:02:01
NA19240 DRB1	17094	20016	0.591	25.4	0.915	DRB1*12:01:01
NA19240 DQB1	14654	14654	1.000	19.0	0.930	DQB1*05:02:01
NA19240 DQB1	10926	10959	0.951	14.2	0.931	DQB1*03:01:01

表 2 A

真実：

NA19238 HLA-A	A*30:01
NA19238 HLA-A	A*36:01
NA19238 HLA-B	B*53:01
NA19238 HLA-B	B*57:03
NA19238 HLA-C	C*18:01
NA19238 HLA-C	C*04:01
NA19238 DRB1	DRB1*16:02
NA19238 DRB1	DRB1*11:01
NA19238 DQB1	DQB1*06:02
NA19238 DQB1	DQB1*05:02
NA19239 HLA-A	A*02:01
NA19239 HLA-A	A*68:02
NA19239 HLA-B	B*35:01
NA19239 HLA-B	B*52:01
NA19239 HLA-C	C*04:01
NA19239 HLA-C	C*16:01
NA19239 DRB1	DRB1*13:01
NA19239 DRB1	DRB1*13:01
NA19239 DQB1	DQB1*05:01
NA19239 DQB1	DQB1*03:01
NA19240 HLA-A	A*30:01
NA19240 HLA-A	A*68:02
NA19240 HLA-B	B*35:01
NA19240 HLA-B	B*57:03
NA19240 HLA-C	C*18:01
NA19240 HLA-C	C*04:01
NA19240 DRB1	DRB1*16:02
NA19240 DRB1	DRB1*12:01
NA19240 DQB1	DQB1*05:02
NA19240 DQB1	DQB1*03:01

20

30

表 2 B

【0045】

ここでは、不明瞭な数字を上記の対立遺伝子から除いた。例えば、予測がA\*04:02:01とA\*04:02:02の場合、最後の不明瞭な数字(ここでは01または02)を除き、したがって予測A\*04:02を得た。

予測されたHLA型と実験に基づいて決定されたHLA型(「真実」と)の間の相違をさらに調査することにより、以下にさらに詳細に検討するように、NA19238とNA19239がNA19240の両親であった場合、実験に基づいて決定されたHLAは予想された遺伝パターンと一致しなかったことが驚くべきことに明らかになった。

40

【0046】

C\*18:01と決定される「真実」および予測されるC\*18:02に関して、これらの2つの対立遺伝子形態間にわずか一塩基の変化があることが注目される。特に、C\*18:01は、WGSデータ中にリードサポートがゼロであるCTGGTTGTC(関連する配列部分のみ)の配列を有するが、C\*18:02はWGSデータ中にそれをサポートする33のリードがあるCTGGCTGTC(関連する配列部分のみ)の配列を有する。このデータによれば、「真実」C\*18:01に対するサポートはないが、予測されるC\*18:02に対しては多数のサポートがある。

50

## 【 0 0 4 7 】

DRB1\*13:01と決定される「真実」および予測されるDRB1\*12:01に関して:NA19240は両親NA19238とNA19239の子供であることが注目される。子供は各親から各HLA型についての対立遺伝子を1つだけ受け継ぐので、真の対立遺伝子は単純な基本的メンデル遺伝から決定できる:

親1 (NA19238): 16:02、 11:01

親2 (NA19239): 13:01、 ? 問題の対立遺伝子

子供 (NA19240): 16:02、 12:01

## 【 0 0 4 8 】

上記からわかるように、子供は親1から16:02を受け継がなければならない、このことは対立遺伝子12:01が親2から来なければならないことを意味する。特に、「真実」は13:01として親2についての第2の対立遺伝子を記載するが、これは遺伝に基づくとは不可能である。親2の予測される対立遺伝子は12:01である。しかし、これは、まさに遺伝に基づいて予想するものである。したがって、上記の例に基づいて、「不正確な」予測は、実際のところ「真実」における誤りに起因していた。このように、本明細書に示すHLA予測方法は、3つの個々のデータセットにおいて5つのHLAのそれぞれ異なるパネルにわたり100%の正確度を示した。上記の予測が平均的カバレッジのWGSサンプルを用いて行われたことを、さらに認識すべきである。本方法の正確度は、腫瘍によって発現される対立遺伝子の識別を可能にするRNS配列データを用いてさらにより改善され得る。これは、時には、DNAに存在する2つの対立遺伝子の1つだけであり得る。意図されるシステムおよび方法のさらなる有利な態様において、DNAもしくはRNA、またはDNAとRNAの両方の組み合わせを処理して高度に正確であるHLA予測を行うことができ、かつそれらを腫瘍または血液のDNAもしくはRNAから得ることができる。さらに、意図される方法は、26のすべてのHLA型についての予測を極めて迅速に(実行時間は一般的に5分未満)取得し、かつ新たに発見された、または極めてまれなHLA対立遺伝子が自明な方法で追加できる。最後に、集団に基づくヒューリスティクスは正確な結果を出すために必要とされないことに留意する必要がある。

## 【 0 0 4 9 】

したがって、本明細書に示すシステムおよび方法を用いて、ゲノム解析で明らかにされる異なるオブジェクトを確証または確認できることを認識すべきである。さらに、同じグラフでRNA情報を使用する場合、変異対立遺伝子発現を直ちに識別できる。さらに、上記の結果と考察に基づいて、システムおよび方法が、RNA-Seqを使用して遺伝子融合を、特に「実施可能な融合」(例えばBCR-ABL)または発癌遺伝子のアイソフォーム(例えばEGFRvIII)をコールできるであろうとも考えられる。

## 【 0 0 5 0 】

本明細書で用いる場合、文脈が明らかに指示しない限り、「に連結される」という用語は、直接連結(互いに連結される2つのエレメントが互いに接触する)および間接連結(少なくとも1つの追加のエレメントが2つのエレメント間に位置する)の両方を含むと意図される。したがって、「に連結される」および「と連結される」という用語は同義的に使用される。さらに、本明細書に開示される本発明の代替エレメントまたは実施形態のグループ化は、限定として解釈されるべきではない。各グループメンバーは、個別に、または本明細書に見られるグループの他のメンバーもしくは他のエレメントとの任意の組み合わせで参照され、または主張され得る。グループの1または複数のメンバーは、利便性および/または特許性の理由から、グループに包含され、もしくはグループから削除され得る。そのような包含または削除が行なわれた場合、本明細書は、修正されたグループを含み、したがって、添付の特許請求の範囲に使用されるすべてのマーカッシュグループの記載を満たすと本明細書ではみなされる。

## 【 0 0 5 1 】

すでに記述されているものの他にさらに多くの修正が本明細書の発明概念を逸脱しない範囲で可能であることは当業者にとって明らかである。したがって、本発明の主題は、添

10

20

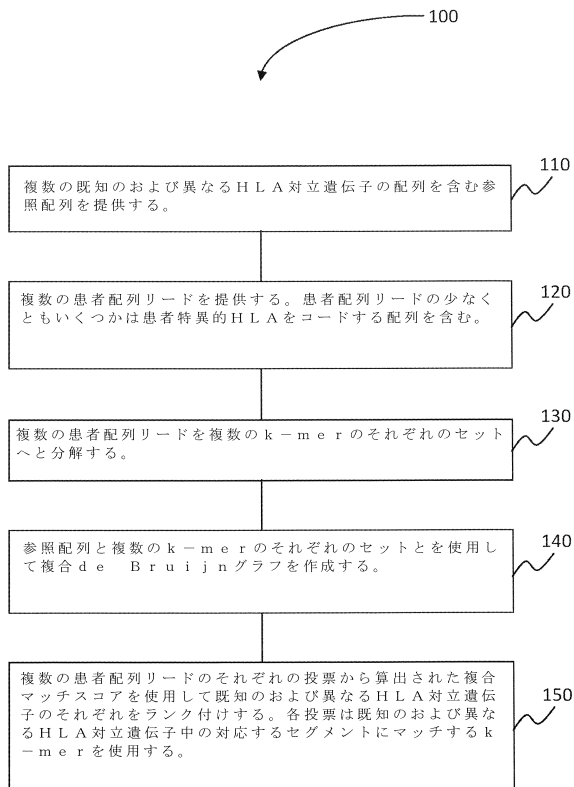
30

40

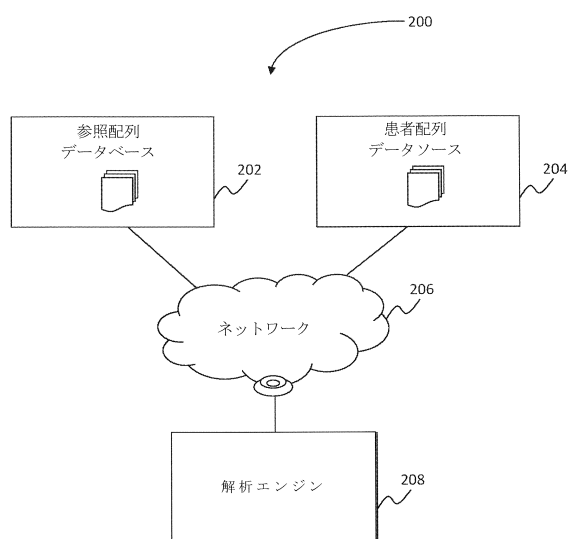
50

付の特許請求の範囲を除いて限定されるべきではない。さらに、本明細書および特許請求の範囲の両方を解釈する際に、すべての用語は、文脈と一致する最も広い可能な方法で解釈されなければならない。特に、「含む」および「含んでいる」という用語は、エレメント、成分、またはステップを参照して、非排他的方法で解釈すべきであり、言及するエレメント、成分、またはステップが、明白に参照していない他のエレメント、成分、またはステップとともに存在し、もしくは利用され、もしくは組み合わせられてもよいことを指示している。本明細書、特許請求の範囲がA、B、C . . . .、およびNからなる群から選択されるもののうちの少なくとも1つを指す場合、本文は、A + NまたはB + Nなどではないその群からの唯一のエレメントを要求していると解釈すべきである。

【図1】



【図2】



---

フロントページの続き

(72)発明者 サンボーン, ジョン ザキャリー  
アメリカ合衆国, カリフォルニア州 95065, サンタ クルス, 195 ケニー アベニュー

審査官 松野 広一

(56)参考文献 米国特許出願公開第2015/0110754 (US, A1)  
米国特許出願公開第2014/0114584 (US, A1)  
特開2015-035212 (JP, A)  
米国特許出願公開第2013/0267429 (US, A1)  
Alexander DILTHEY et al., Improved genome inference in the MHC using a population reference graph, NATURE GENETICS, 2015年 6月, Vol.47 No.6, pp.682-688

(58)調査した分野(Int.Cl., DB名)

G16B 5/00 - 99/00

C12N 15/12

C12Q 1/6806

PubMed