

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6618929号

(P6618929)

(45) 発行日 令和1年12月11日 (2019. 12. 11)

(24) 登録日 令和1年11月22日 (2019. 11. 22)

(51) Int. Cl. F I
G 1 6 B 30/00 (2019. 01)
C 1 2 Q 1/6869 (2018. 01)
 G 1 6 B 30/00
 C 1 2 Q 1/6869 Z

請求項の数 15 (全 36 頁)

(21) 出願番号	特願2016-565058 (P2016-565058)	(73) 特許権者	591003013
(86) (22) 出願日	平成27年5月12日 (2015. 5. 12)		エフ. ホフマン-ラ ロシュ アーゲー
(65) 公表番号	特表2017-520821 (P2017-520821A)		F. HOFFMANN-LA ROCH
(43) 公表日	平成29年7月27日 (2017. 7. 27)		E AKTIENGESSELLSCHAF
(86) 国際出願番号	PCT/EP2015/060442		T
(87) 国際公開番号	W02015/173222		スイス・シーエイチー４０７０バーゼル・
(87) 国際公開日	平成27年11月19日 (2015. 11. 19)		グレンツアーヘルストラッセ１２４
審査請求日	平成30年5月11日 (2018. 5. 11)	(74) 代理人	100099759
(31) 優先権主張番号	61/991, 820		弁理士 青木 篤
(32) 優先日	平成26年5月12日 (2014. 5. 12)	(74) 代理人	100077517
(33) 優先権主張国・地域又は機関	米国 (US)		弁理士 石田 敬
		(74) 代理人	100087871
			弁理士 福本 積
		(74) 代理人	100087413
			弁理士 古賀 哲次

最終頁に続く

(54) 【発明の名称】 ウルトラディープシークエンシングにおける希少バリエーションコール

(57) 【特許請求の範囲】

【請求項 1】

第一の試料中の標的領域における低頻度バリエーションを検出するための、コンピュータで実行される方法であって、コンピュータシステムにおいて、

- 1又は2以上の試料からのDNA断片のシークエンシングから得られる複数の配列リードを受け取り、ここで前記1又は2以上の試料は第一の試料を含み、前記シークエンシングは前記DNA断片中の標的領域を標的化することを含み；

- 前記複数の配列リードを、参照配列の標的領域にアラインメントし；

- 第一の試料の配列リードに基づいて、標的領域の第一の位置において第一のバリエーションアレルを有する第一の候補バリエーションを同定し、ここで前記第一のバリエーションアレルは、前記参照配列の第一の位置の参照アレルとは異なり；

- 前記参照配列の第一の位置にアラインメントする第一の試料の配列リードに基づいて、第一の位置における第一のバリエーションアレルに関する第一のバリエーション頻度を決定し、

- 前記参照アレルを有する前記参照配列の標的領域中の1セットの第二の位置を同定し、ここで前記1セットの第二の位置は前記第一の位置を含み；

前記1セットの第二の位置の各々の位置において、かつ前記1又は2以上の試料の各々に関して；

- 前記参照配列の1セットの第二の位置の各々の位置にアラインメントする試料の配列リードに基づいて、第一のバリエーションアレルの第二のバリエーション頻度を決定し、ここで前記第二のバリエーション頻度は統計分布を形成し；

10

20

- 前記第一のバリエーション頻度を前記統計分布の統計値と比較して、前記統計分布の統計値に対する第一のバリエーション頻度の確率値を決定し；そして

- 前記第一の位置における第一のバリエーションアレルに関して、第一の試料において第一の候補バリエーションが真陽性であるか否かを決定する一部として、前記確率値を閾値と比較すること、ここで前記閾値は前記第一のバリエーションアレルに関して偽陽性と真陽性とを区別する、

含む、方法。

【請求項 2】

前記参照配列は、正常細胞から決定されるとおりのコンセンサス配列に対応する、請求項 1 に記載の方法。

10

【請求項 3】

前記 1 又は 2 以上の試料は無細胞 DNA 断片由来である、請求項 1 又は 2 に記載の方法。

【請求項 4】

前記 1 又は 2 以上の試料は生物試料の RNA 由来である、請求項 1 又は 2 に記載の方法。

【請求項 5】

複数の試料は、単一のシーケンシングランにおいてシーケンシングされる、請求項 1 ~ 4 のいずれか 1 項に記載の方法。

【請求項 6】

前記確率値は、z スコア、修正された z スコア、累積確率、フレッド (Phred) クオリティスコア又は修正されたフレッドクオリティスコアである、請求項 1 ~ 5 のいずれか 1 項に記載の方法。

20

【請求項 7】

前記統計分布は、前記第二のバリエーション頻度の対数変換の統計分布である、請求項 1 ~ 6 のいずれか 1 項に記載の方法。

【請求項 8】

前記閾値は、既知の真陽性と偽陽性とを有するトレーニングデータに基づいて、サポートベクターマシン分類器を用いて決定される、請求項 1 ~ 7 のいずれか 1 項に記載の方法。

30

【請求項 9】

第一の試料中の標的領域における第一の位置に第一のバリエーションアレルを有するバリエーションを検出するための、コンピュータで実行される方法であって、コンピュータシステムにおいて、

- 少なくとも 2 つの試料からの DNA 断片のシーケンシングから得られる複数の配列リードを受け取り、ここで前記少なくとも 2 つの試料は第一の試料を含み、前記シーケンシングは前記 DNA 断片中の標的領域を標的化することを含み；

- 前記複数の配列リードを、参照配列の標的領域にアラインメントし；

- 第一の位置において各試料のアラインメントされた配列リードに基づいて、前記第一のバリエーションアレルが、前記少なくとも 2 つの試料の各々の試料中の第一の位置に存在するか否かを同定し、ここで前記第一のバリエーションアレルは、前記参照配列の第一の位置における参照アレルとは異なり；

40

- 前記少なくとも 2 つの試料の各試料に関する第一の位置にアラインメントされた配列リードを用いて、前記第一の位置における第一のバリエーションアレルのバリエーション数及び前記第一の位置における参照アレルの野生型数を決定し；

- 前記少なくとも 2 つの試料から、1 つ又は 2 つ以上の試料を 1 つ又は 2 つ以上の参照試料として選択し；

- 第一の試料に関する第一の位置における第一のバリエーションアレルの第一のバリエーション数及び第一の位置における参照アレルの第一の野生型数と、前記 1 つ又は 2 つ以上の参照試料に関する第一の位置における第一のバリエーションアレルの第二のバリエーション数及び第一

50

の位置における参照アレルの第二の野生型数とを比較して、第一の試料に関する第一の位置に第一のバリエーションアレルを有するバリエーションが発生する確率値を決定し；そして

- 第一のバリエーションアレルに関して第一の試料中の第一の位置における第一のバリエーションアレルが真陽性を含むか否かを決定する一部として、前記確率値を閾値と比較すること、ここで前記閾値は第一の位置の第一のバリエーションアレルに関して偽陽性と真陽性とを区別する、

含む、方法。

【請求項 10】

前記確率値は、以下の式：

【数 1】

10

$$\chi^2 = n \times (a_1 \times w_2 - a_2 \times w_1)^2 / (n_1 \times n_2 \times a \times w)$$

[式中、

a₂ は第一のバリエーション数であり、w₂ は第一の野生型数であり、a₁ は第二のバリエーション数であり、w₁ は第二の野生型数であり、a は a₁ 及び a₂ の合計であり、w は w₁ 及び w₂ の合計であり、n₁ は a₁ 及び w₁ の合計であり、n₂ は a₂ 及び w₂ の合計であり、n は n₁ 及び n₂ の合計である]

によって定義されるカイ二乗統計値に基づくカイ二乗累積分布関数を用いて決定される、請求項 9 に記載の方法。

20

【請求項 11】

前記確率値は、2 つの比率である p₁ 及び p₂ に基づいて決定され、ここで p₁ = a₁ / n₁ であり、p₂ = a₂ / n₂ であり、a₂ は第一のバリエーション数であり、a₁ は第二のバリエーション数であり、n₂ は第一のバリエーション数及び第一の野生型数の合計であり、そして n₁ は第二のバリエーション数及び第二の野生型数の合計である、請求項 9 に記載の方法。

【請求項 12】

前記確率値は、z スコア、修正された z スコア、p 値、カイ二乗値、累積確率値及び信頼水準を示すクオリティスコアの 1 つ又は 2 つ以上である、請求項 9 に記載の方法。

30

【請求項 13】

前記クオリティスコアはルックアップテーブルを用いて決定され、ここで前記ルックアップテーブルの入力は、カイ二乗値又は正規クオントイル値の少なくとも 1 つである、請求項 12 に記載の方法。

【請求項 14】

前記閾値は、既知の真陽性と偽陽性とを有するトレーニングデータに基づいて、サポートベクターマシン分類器を用いて決定される、請求項 9 ~ 13 のいずれか 1 項に記載の方法。

【請求項 15】

第一の試料の標的領域における真のバリエーションを検出するようにコンピュータシステムを制御するためのコンピュータプログラムであって、以下の指示：

40

- 1 又は 2 以上の試料からの DNA 断片のシーケンシングから得られる複数の配列リードを受容させ、ここで前記 1 又は 2 以上の試料は第一の試料を含み、前記シーケンシングは前記 DNA 断片中の標的領域を標的化することを含み；

- 前記複数の配列リードを、参照配列の標的領域にアラインメントさせ；

- 複数のバリエーションクラス中の 1 つのバリエーションクラス中のバリエーションの参照アレルを各々が有する参照配列の標的領域における 1 セットの配列位置を同定させ、ここで、前記複数のバリエーションクラスの各々は、1 又は 2 以上のバリエーションを含むように定義され、各々のバリエーションは対応する参照アレルとは異なるバリエーションアレルを有し、そして前記 1 セットの配列位置は第一の位置を含み、

50

前記 1 セットの配列位置の各々の位置において、かつ前記 1 又は 2 以上の試料の各々の試料に関して：

- 各試料に関する各位置におけるリード数を決定させ；
- 各試料の配列リードに基づいて、前記バリエーションクラス中のバリエーションに関するバリエーションアレルを有する候補バリエーションを同定させ、ここで前記バリエーションアレルは、前記参照配列の同一の位置における参照アレルとは異なり、そして各試料中の各位置における候補バリエーションの総数は、各試料に関する各位置中のバリエーション数であり；
- 前記リード数及びバリエーション数に基づいてバリエーションクラス中のバリエーションのバリエーション頻度を決定させ、ここで各試料中の各位置に関するバリエーション頻度は、統計分布を形成し、第一の試料に関する前記 1 セットの配列位置中の第一の位置におけるバリエーション頻度は、第一のバリエーション頻度であり；
- 前記第一のバリエーション頻度を統計分布の値と比較して、前記統計分布の値に対する前記第一のバリエーション頻度の確率値を決定させ；そして
- 第一の試料中の候補バリエーションが真陽性であるか否かを決定する一部として、前記確率値を閾値と比較させること、ここで前記閾値は、前記バリエーションクラス中のバリエーションに関して偽陽性と真陽性とを区別する、を含む、前記コンピュータプログラム。

【発明の詳細な説明】

【背景技術】

【0001】

核酸シーケンシングは、ある DNA 分子又は RNA 分子中に存在するヌクレオチドの順序を決定する。より安価で迅速なシーケンシング法に対する要求が、次世代シーケンシング (NGS) 法の開発を推進してきた。NGS プラットフォームは、大規模な平行シーケンシングを行い、この間に、複数の試料からの数百万の DNA 断片をまとめてシーケンシングすることができ、こうして伝統的なサンガー (Sanger) シーケンシングに対してはるかに安価で高スループットな代替法を提供する。NGS は全ゲノムシーケンシング又は標的化シーケンシングで使用する事ができる。標的化シーケンシングでは、ゲノム中の遺伝子又は規定された領域のサブセットがシーケンシングされるか、又は、例えば主に標的領域を増幅することによりシーケンシングされる。

【0002】

ウルトラディープシーケンシング (ultra-deep sequencing) は、一般的な及び希少な配列の変化を同定することを目的とする、広範囲のアンプリコンのシーケンシングである。十分なカバー率を有するため、ウルトラディープシーケンシングは、希少配列バリエーションを 1 % 未満まで完全に特性評価する能力を有する。ウルトラディープシーケンシングは、低頻度 (low-frequency) H I V 薬耐性変異を検出するために、又は複雑な癌試料中の希少体細胞変異を同定するために使用されている。非侵襲性血液検査などの試験に関して、バイオマーカー変異の頻度は 1 % 未満になることがある。しかし NGS はエラーが発生しやすいプロセスであり、シーケンシング深度 (sequencing depth)、試料のタイプ、及びシーケンシングプロトコルによっては、ほぼ 1 % 以上のエラー率になることがある。従って、1 % 未満の頻度を有するバリエーションに関しては偽陽性 (false positives) が現れる可能性があるため、現在の多くの NGS ソフトウェアパッケージは、1 % 以上の頻度を有するバリエーションのみを報告している。それでも、例えば 1 % 未満の低頻度を有するバリエーションに関してさえ、真陽性 (true positives) が存在し得る。従って、例えば約 0.0025 % ~ 約 1 % という低い頻度を有するバリエーションに関して、真陽性を検出する方法とシステムが必要とされる。

【発明の概要】

【0003】

実施態様は、例えば標的化シーケンシングから得られた試料のシーケンシングリード (sequencing reads) に基づいてより多くの正確なバリエーションコール (variant calls) を行うための方法、システム、及び装置を提供することができる。例えば、いったん配列

10

20

30

40

50

リード (sequence reads) が受け取られ (received)、参照配列 (reference sequence) にアラインメントされる (aligned) と、ある位置にバリエーションを有するシーケンシングリードがカウントされる。試料の 1 つの位置で測定される特定のバリエーションの第一のバリエーション頻度を、他の位置で測定される及び / 又は他の試料からの特定のバリエーションの 1 又は 2 以上の第二のバリエーション頻度と比較することができる。第二のバリエーション頻度は、シーケンシングランに関するシーケンシングエラーの予測値に対応することができる。

【 0 0 0 4 】

いくつかの実施態様において、バリエーションがある位置で真陽性である信頼水準を示す確率値 (probability) は、1 又は 2 以上の試料の標的領域中の複数の位置におけるバリエーション数及び総リード数に基づいて計算することができる。確率値は次に閾値レベルと比較されて、検出されたバリエーションが真陽性であるか否かを決定することができる。他の実施態様において、試験試料と参照試料 (例えば、その位置においてシーケンシングエラーのみを有すると推定される) 中の同一の位置におけるバリエーション数と総リード数の差を用いて、試験試料中のバリエーションが真陽性であるかどうかを決定することができる。

【 0 0 0 5 】

ある実施態様において、ある方法は、試験試料の標的領域における希少バリエーションについて真陽性を検出することができる。各試料について、参照配列上で参照アレル (reference allele) が存在する位置における同じバリエーションクラスのバリエーションに関するバリエーション頻度は、バリエーション数と総リード数を用いて計算することができる。同一のクラスのバリエーションに関するバリエーション頻度の分布を使用して、決定されたバリエーション頻度を有する試験試料中のある位置におけるバリエーションの確率値を決定することができる。この確率値に基づいて、試験試料中の位置におけるバリエーションは、真陽性 (変異 (mutation)) 又は偽陽性として分類される。

【 0 0 0 6 】

他の実施態様において、ある方法は、1 又は 2 以上の参照試料との比較を用いて、試験試料の標的領域における希少バリエーションに関する真陽性を検出することができる。試験試料中の特定の位置における特定のバリエーションに関するバリエーション数と野生型数は、アラインメントされた配列リードから決定することができ、1 又は 2 以上の参照試料中の特定の位置における特定のバリエーションに関するバリエーション数と野生型数と比較して、確率値を決定することができる。この確率値に基づいて、試験試料中の特定の位置における特定のバリエーションは、真陽性又は偽陽性として分類される。

【 0 0 0 7 】

ある実施態様において、第一の試料の標的領域における低頻度バリエーションを検出するためのコンピュータで実行される方法が提供される。ここで、この方法は (コンピュータシステムにおいて)、

1 又は 2 以上の試料からの DNA 断片のシーケンシングから得られる複数の配列リードを受け取り、ここで前記 1 又は 2 以上の試料は第一の試料を含み、前記シーケンシングは前記 DNA 断片中の標的領域を標的化することを含み;

前記複数の配列リードを、参照配列の標的領域にアラインメントし;

第一の試料の配列リードに基づいて、標的領域の第一の位置において第一のアレルを有する第一の候補バリエーションを同定し、ここで前記第一のアレルは、前記参照配列の第一の位置の参照アレルとは異なり;

前記参照配列の第一の位置にアラインメントする第一の試料の配列リードに基づいて、第一の位置における第一のアレルに関する第一のバリエーション頻度を決定し、

複数のバリエーションクラスから選択される第一のバリエーションクラスに対応する第一の候補バリエーションを同定し、ここで前記複数のバリエーションクラスの各バリエーションクラスは、異なるタイプのバリエーションに対応し;

前記参照アレルを有する前記参照配列の標的領域における 1 セットの第二の位置を同定し、ここで前記 1 又は 2 以上の試料中の少なくとも 50 % の他の位置は、第一のアレルに関して偽陽性を示し、そして前記 1 セットの第二の位置は前記第一の位置を含み;

10

20

30

40

50

前記 1 セットの第二の位置の各々において、かつ前記 1 又は 2 以上の試料の各々に関して：

前記参照配列の第二の位置にアラインメントする試料の配列リードに基づいて、第一のアレルの第二のバリエーション頻度を決定し、ここで前記第二のバリエーション頻度は統計分布を形成し；

前記第一のバリエーション頻度を前記統計分布の統計値と比較して、前記統計分布の統計値に対する第一のバリエーション頻度の確率値を決定し；そして

前記第一のアレルに関して、第一の試料において第一の候補バリエーションが真陽性であるか否かを決定する一部として、前記確率値を閾値と比較すること、ここで前記閾値は前記第一のアレルに関して偽陽性と真陽性とを区別する、を含む。

10

【 0 0 0 8 】

ある実施態様において、前記参照配列は、正常細胞から決定されたとおりのコンセンサス配列に対応する。いくつかの実施態様において、前記 1 又は 2 以上の試料は無細胞 DNA 断片由来である。いくつかの実施態様において、前記 1 又は 2 以上の試料は生物試料の RNA 由来である。いくつかの実施態様において、複数の試料は、単一のシーケンシングランにおいてシーケンシングされる。他の実施態様において、前記統計分布の統計値は平均値を含む。他の実施態様において、前記確率値は、z スコア、修正された(modified) z スコア、累積確率、フレッド (Phred) クオリティスコア又は修正された(modified) フレッドクオリティスコアである。他の実施態様において、前記統計分布は、前記第二のバリエーション頻度の対数変換の統計分布である。他の実施態様において、前記閾値は、1 又は 2 以上のシーケンシングランから得られるトレーニングデータに基づくサポートベクターマシン分類器 (support vector machines classifier) を用いて決定される。他の実施態様において、前記閾値はバリエーション頻度の関数である。

20

【 0 0 0 9 】

別の実施態様において、第一の試料中の標的領域における第一の位置に第一のアレルを有するバリエーションを検出するための、コンピュータで実行される方法が提供される。ここでこの方法は (コンピュータシステムにおいて)、

少なくとも 2 つの試料からの DNA 断片のシーケンシングから得られる複数の配列リードを受け取り、ここで前記少なくとも 2 つの試料は第一の試料を含み、前記シーケンシングは前記 DNA 断片中の標的領域を標的化することを含み；

30

前記複数の配列リードを、参照配列の標的領域にアラインメントし；

第一の位置において各試料のアラインメントされた配列リードに基づいて、前記第一のアレルが、前記少なくとも 2 つの試料の各々の試料中の第一の位置に存在するか否かを同定し、ここで前記第一のアレルは、前記参照配列の第一の位置における参照アレルとは異なり；

前記少なくとも 2 つの試料の各試料に関する、第一の位置における第一のアレルのバリエーション数及び第一の位置の参照アレルの野生型数を決定し；

前記少なくとも 2 つの試料から、少なくとも 1 つの試料を参照試料として選択し；

第一の試料に関する第一の位置における第一のアレルの第一のバリエーション数及び第一の位置における参照アレルの第一の野生型数と、前記参照試料に関する第一の位置における第一のアレルの第二のバリエーション数及び第一の位置における参照アレルの第二の野生型数とを比較して、第一の試料に関する第一の位置に第一のアレルを有するバリエーションの確率値を決定し；そして

40

第一のアレルに関して第一の試料中の第一の位置における第一のアレルが真陽性であるか否かを決定する一部として、前記確率値を閾値と比較すること、ここで前記閾値は第一の位置の第一のアレルに関して偽陽性と真陽性とを区別する、を含む。

【 0 0 1 0 】

ある実施態様において、前記参照試料は、第一の試料以外の前記少なくとも 2 つの試料のうち、第一の位置における第一のアレルに関して最も低いバリエーション頻度を有する 2 つの試料を含む。いくつかの実施態様において、前記確率値は、カイ二乗累積分布関数 (ch

50

i-squared cumulative distribution function) を用いて決定される。いくつかの実施態様において、前記確率値は、ピアソン比率検定 (Pearson proportion test) を用いて決定される。いくつかの実施態様において、前記確率値は、 z スコア、修正された (modified) z スコア、 p 値、カイ二乗値、累積確率値及びクオリティスコアの 1 つ又は 2 つ以上である。いくつかの実施態様において、前記クオリティスコアはルックアップテーブル (look-up table) を用いて決定される。いくつかの実施態様において、前記閾値は、1 又は 2 以上のシーケンシングランから得られるトレーニングデータに基づくサポートベクターマシン分類器を用いて決定される。いくつかの実施態様において、前記閾値はバリエーション頻度の関数である。

【0011】

10

別の実施態様において、実施されるときにコンピュータシステムを制御して第一の試料の標的領域における真のバリエーションを検出する複数の指示を記憶する、非一過性の (non-transitory) コンピュータ可読媒体を含むコンピュータ製品が提供される。ここで前記指示は、

1 又は 2 以上の試料からの DNA 断片のシーケンシングから得られる複数の配列リードを受け取り、ここで前記 1 又は 2 以上の試料は第一の試料を含み、前記シーケンシングは前記 DNA 断片中の標的領域を標的化することを含み；

前記複数の配列リードを、参照配列の標的領域にアラインメントし；

バリエーションクラス中のバリエーションの参照アレルを有する参照配列の標的領域における 1 セットの配列位置を同定し、ここで、前記 1 又は 2 以上の試料中の少なくとも 50% の配列位置は、前記配列リード中のバリエーションクラス中のバリエーションに関して偽陽性を示し、そして前記 1 セットの配列位置は第一の位置を含み、

20

前記 1 セットの配列位置の各々の位置において、かつ前記 1 又は 2 以上の試料の各々の試料に関して；

各試料に関する各位置におけるリード数を決定し；

各試料の配列リードに基づいて、前記バリエーションクラス中のバリエーションに関するバリエーションアレルを有する候補バリエーションを同定し、ここで前記バリエーションアレルは、参照配列の同一の位置における参照アレルとは異なり、そして各試料中の各位置における候補バリエーションの総数は、各試料に関する各位置中のバリエーション数であり；

前記リード数及びバリエーション数に基づいてバリエーションクラス中のバリエーションのバリエーション頻度を決定し、ここで各試料中の各位置に関するバリエーション頻度は、統計分布を形成し、第一の試料に関する前記 1 セットの配列位置中の第一の位置におけるバリエーション頻度は、第一のバリエーション頻度であり；

30

前記第一のバリエーション頻度を統計分布の値と比較して、前記統計分布の値に対する前記第一のバリエーション頻度の確率値を決定し；そして

第一の試料中の候補バリエーションが真陽性であるか否かを決定する一部として、前記確率値を閾値と比較すること、ここで前記閾値は、前記バリエーションクラス中のバリエーションに関して偽陽性と真陽性とを区別する、を含む。ある実施態様において、前記統計分布は、各試料に関する各位置におけるバリエーション頻度の対数変換の統計分布である。

【0012】

40

他の実施態様は、本明細書に記載の方法に関連するシステム、装置、及びコンピュータ可読媒体に関する。

【0013】

以下の定義、詳細な説明、及び添付図面を参照することにより、本発明の本質と利点のより良い理解が得られるであろう。

【0014】

定義

本明細書において用語「試料 (sample)」又は「生物試料 (biological sample)」は、核酸を含むか又は含むと推定される任意の組成物を指す。核酸は、動物 (例えば哺乳動物、ヒト)、植物、微生物などに由来してもよい。試料という用語は、細胞、組織、又は血

50

液の、精製されたか又は分離された成分、例えばDNA、RNA、タンパク質、無細胞部分、又は溶解物を含む。試料はまた、他のタイプの生物試料、例えば皮膚、血漿、血清、全血、及び血液成分（パフィーコート）、唾液、尿、涙、精液、膿液、吸引物若しくは洗浄液、組織生検、及び他の体液や組織（パラフィン包埋組織を含む）を指すこともできる。試料はまた、細胞株を含む個体から得られた細胞のインビトロ培養物の成分及び構成要素を含むことができる。「試験試料 (test sample)」は、試料中のバリエーションを検出するための試験中の試料を指す。

【0015】

「ゲノムセグメント (genomic segment)」（「ゲノム断片 (genomic fragment)」とも呼ばれる）は、生物のゲノム由来である完全に又は部分的にシーケンシングされた核酸分子である。これは、DNAセグメント（「DNA断片」とも呼ばれる）又はRNAセグメント（「RNA断片」とも呼ばれる）でもよい。セグメントは、ゲノムの大きな部分を断片化することによって、例えば細胞を音波に供することによって作成することができる。ゲノムセグメントはシーケンシングして、「シーケンシングリード (sequencing read)」（「配列リード (sequence read)」又は単に「リード (read)」とも呼ばれる）を提供することができる。シーケンシングリードは、全ゲノムセグメント又はセグメントの一部であってもよい。

10

【0016】

「参照試料 (reference sample)」（「対照試料 (control sample)」とも呼ばれる）は、試験試料と比較するため、基準、通常既知の基準となる試料を指す。例えば試験試料は、癌又は癌関連変異を有することが疑われる個体から採取することができ、癌のない個体又は癌関連変異のない個体からの参照試料（陰性対照）、又は癌又は癌関連変異を有することがわかっている個体からの参照試料（陽性対照）と比較することができる。対照はまた、多くの試験若しくは結果から集められた平均値又は範囲を表すことができる。

20

【0017】

「標的領域 (target region)」とは、分析される配列中の領域であって、診断的関連性を有し得る領域のことである。一例として、標的領域を含む断片は、プライマー及び増幅プロセスを用いて増幅するか、又はプローブを用いて濃縮することができる。「参照配列 (reference sequence)」（単に「参照 (reference)」とも呼ばれる）は、配列リードがアラインメントされる任意の既知の配列である。種々の実施態様において、参照配列は、生物のゲノム又はトランスクリプトームのすべて又は一部のみに対応することができる。参照配列はまた、2種以上の生物のゲノムを含むことができる。例えば配列リードはまた、試料中に存在し得るウイルスのデータベースと比較することができる。

30

【0018】

バリエーション (variant)（変化 (variation)又は変異 (mutation)とも呼ばれる）は、2つの配列間の差を指す。バリエーションは、例えば1つの塩基の1又は2以上の他の塩基への変化、1又は2以上の塩基の挿入、又は1又は2以上の塩基の欠失でもよい。参照配列中の位置の塩基は参照アレルと呼ばれることがあり、一方、試験試料上の同一の位置の異なる塩基（又は挿入もしくは欠失）はバリエーションアレルと呼ぶことができる。例えばA > Cの単一塩基置換に関して、Aは参照アレルであり、そしてCはバリエーションアレルである。参照アレルは、天然に存在する生物のための最も一般的な遺伝子型を示す野生型アレルであってもよい。配列リードと参照配列の標的領域との差はカウントすることができ、真の変異が同定される可能性がある（例えば、十分な配列リードが変異を示す場合）。

40

【0019】

試料の異なる配列リード上の特定の位置における同一のバリエーションアレルの総数（例えばA > CバリエーションについてのCの数）は、バリエーション数 (variant count)と呼ばれる。ある試料の特定の位置のリードの総数は、リード数 (read count)と呼ばれる。試料の特定の位置におけるバリエーションタイプ又はクラス（例えばA > C）のバリエーション頻度は、試料の特定の位置におけるリード数に対する特定の位置におけるバリエーションに関するバリエーション数の比率として定義される。

50

【 0 0 2 0 】

本明細書において用語「位置 (location)」は、配列中の (例えばゲノムの標的領域中の) 1 又は 2 以上の位置に対応する。例えば多塩基挿入が存在する場合、任意の長さのヌクレオチド (又は塩基対) が位置中に存在してもよい。

【 0 0 2 1 】

特に別の指定がなければ、本明細書で使用される技術用語及び科学用語は、一般に、当業者によって理解されるものと同じ意味を有する。例えば、Pfaffl, Methods: The ongoing evolution of qPCR, vol. 50 (2010); van Pelt-Verkuil et al. Principles and Technical Aspects of PCR Amplification, Springer (2010); Lackie, DICTIONARY OF CELL AND MOLECULAR BIOLOGY, Elsevier (4th ed. 2007); Sambrook et al., MOLECULAR CLONING, A LABORATORY MANUAL, Cold Springs Harbor Press (1989) を参照されたい。

【図面の簡単な説明】

【 0 0 2 2 】

【図 1】図 1 は、本発明の実施態様に従う標的化されたウルトラディープシーケンシングのための次世代シーケンシング (NGS) を用いた、ゲノムシーケンシング及びバリエーションコーリングを示すフローチャートである。

【図 2】参照配列と比較された標的領域の配列リードを示し、ここで、異なる配列位置における同一のクラス及び異なるクラスのバリエーションは、本発明の実施態様に従って示される。

【図 3 A】本発明の実施態様に従う 1 又は 2 以上の試料中の標的領域内の複数の位置のそれぞれの位置における、バリエーションクラスのバリエーションに関するバリエーションの頻度分布の理想的な統計モデルを示す。

【図 3 B】特定の試料上の特定の位置におけるバリエーションのバリエーション頻度が、本発明の実施態様に従う特定の Z 値を有するであろう確率を示す。

【図 3 C】Z 値が、本発明の実施態様に従う z 未満の又は z に等しい値をとる確率の累積分布関数を示す。

【図 3 D】バリエーション頻度値又は Z 値を有するバリエーションが偽陽性であり、本発明の実施態様に従うバリエーションコールを作成するための関連するクオリティスコアである確率を示す。

【図 4】本発明の実施態様に従う統計モデルを用いるバリエーションコーリングの方法を示すフローチャートである。

【図 5】本発明の実施態様に従うサポートベクターマシン (SVM) により決定したセパレーターラインを有するエキソン 20 の E G F R T 7 9 0 M のトレーニングデータと試験データに関する統計モデルを用いて決定されたバリエーションクオリティスコア Q_{AMP} を示す。

【図 6】本発明の実施態様に従う SVM により決定したセパレーターラインを有するエキソン 21 の E G F R L 8 5 8 R のトレーニングデータと試験データに関する統計モデルを用いて決定したバリエーションクオリティスコア Q_{AMP} を示す。

【図 7】本発明の実施態様に従う参照試料と試験試料の配列リード上の特定のゲノム位置における特定のバリエーションを示す。

【図 8】本発明の実施態様に従う試験試料と参照試料の配列リードデータを比較することにより、特定の配列位置における特定のバリエーションに関するバリエーションコーリングを示すフローチャートである。

【図 9】本発明の実施態様に従う SVM により決定したセパレーターラインを有するエキソン 20 の E G F R T 7 9 0 M のトレーニングデータと試験データに関する 2 つの試料を比較することにより決定した局所化されたバリエーションクオリティスコア Q_{LOC} を示す。

【図 10】本発明の実施態様に従う SVM により決定したセパレーターラインを有するエキソン 21 の E G F R L 8 5 8 R のトレーニングデータと試験データに関する 2 つの試料を比較することにより決定した局所化されたバリエーションクオリティスコア Q_{LOC} を示す。

【図 1 1】本発明の実施態様に従う SVM により決定したセパレーターラインを有するエキソン 19 の E G F R 15 塩基欠失 2 2 3 5 _ 2 2 4 9 d e l 1 5 のトレーニングデータと試験データに関する 2 つの試料を比較することにより決定した局所化されたバリエーションスコア Q_{LOC} を示す。

【図 1 2】本発明の実施態様に従う効率的なクオリティスコア推定に関するルックアップテーブルを示す。

【図 1 3】本発明の実施態様に従う低頻度バリエーションコーリングに関する例示的コンピュータシステム例のブロック図を示す。

【図 1 4】シーケンシング装置とコンピュータシステムとの関係を示す一般的なブロック図の例である。

【図 1 5 A】本発明の方法及びシステムを実施するために使用できるソフトウェアとハードウェア資源との関係を示す一般的なブロック図の例である。

【図 1 5 B】本発明の方法及びシステムを実施するために使用できるソフトウェアとハードウェア資源との関係を示す一般的なブロック図の例である。

【発明を実施するための形態】

【0023】

シーケンシングは、癌又は他の疾患の突然変異を検出するために使用することができ、またインビトロ診断 (IVD) 検査としても開発することができる。非侵襲的血液検査として、これらの検査を開発することが望ましい。しかし血液試料中のバイオマーカーの変異の頻度は低い。例えば、Kidess and Jeffrey, *Circulating tumor cells versus tumor-derived cellfree DNA: rivals or partners in cancer care in the era of single-cell analysis?* *Genome. Med.*, 5:70 (2013), Diaz and Bardelli, *Liquid biopsies: genotyping circulating tumor DNA*, *J. Clin. Oncol.*, 32:579-586 (2014); and Diehl et al., *Nat Med.*, 14:985-990 (2008) を参照されたい。シーケンシングプロセスに関連するエラーが原因で、閾値が 1 % 以下に設定される時、多くの NGS ソフトウェアパッケージは、1 % 以上の頻度を有するバリエーションのみを報告する。

【0024】

本発明の実施態様は、1 % 未満のバリエーション頻度を有する低頻度バリエーションに関する真陽性を検出するための解決策を提供する。正確なバリエーションコールは、例えば標的化シーケンシングから得られる試料のシーケンシングリードに基づくことができる。例えば、いったん配列リードが受け取られ、参照配列にアラインメントされると、ある位置でバリエーションを有する配列リードはカウントされる。試料の 1 つの位置で測定された特定のバリエーションの第一のバリエーション頻度は、他の位置で測定された及び / 又は他の試料からの特定のバリエーションの 1 又は 2 以上の第二のバリエーション頻度と比較することができる。第二のバリエーション頻度は、シーケンシングランのシーケンシングエラーに関する予測値に対応することができる。

【0025】

いくつかの実施態様において、ある位置でバリエーションが真陽性である信頼水準を示す確率値は、1 又は 2 以上の試料中の標的領域内の複数の位置におけるバリエーション数及び総リード数に基づいて計算することができる。その後、確率値は閾値と比較され、検出されたバリエーションが真陽性であるか否かを決定することができる。他の実施態様において、試験試料と参照試料 (例えば、その位置でシーケンシングエラーのみを有すると推定される) 中の同一の位置におけるバリエーション数と総リード数との差を用いて、試験試料中のバリエーションが真陽性であるか否かを決定することができる。

【0026】

I. 標的化シーケンシングを用いるウルトラディープシーケンシング

ゲノムの特定の領域は、標的化シーケンシングを用いて効率的に分析することができる。例えば生物試料のゲノムセグメントは、標的領域に対応するセグメントをクローニングすることにより (例えば、ポリメラーゼ連鎖反応 (PCR) などの増幅プロセスにおいてプライマーを用いて)、及び / 又は標的領域に対応するセグメントを優先的に捕捉する

10

20

30

40

50

プローブを用いることにより、増加又は増幅させることができる。標的増加試料中のゲノムセグメントは、大規模の平行した次世代シーケンシング (NGS) を用いてシーケンシングし、標的領域内の可能な変異を調査するために分析することができる。

【0027】

しかし、このようなプロセスはエラーを発生させることがある。例えば、増幅又は濃縮の前段階を有する高スループットの次世代シーケンシングを用いるバリエーション検出では、アンプリコン/濃縮ライブラリ (標的増加試料) は偽陽性リードを含む可能性がある。PCR は点突然変異及びインデル (indel) を導入することができ、これはまた、組換え配列又はキメラを生成することができる。更に遺伝的バリエーションの相対頻度は、PCR 中の選択的増幅の偏りにより攪乱されることがある。PCR 中に、追加の単一塩基エラーが起きることがある。シーケンシング自体は、塩基置換エラー及びインデルを導入し得る。これらのエラーは間違った変異報告につながることもあり、疾患の診断のために誤解を招く情報を提供することがある。偽陽性は種々の方法、例えばプライマーの正しい設計及び高忠実度酵素の開発、によって減少させることができる。しかし、それでも偽陽性残って、多くの場合エラー率が約 1 % 以上に大きくなる可能性がある。

【0028】

各個別のヌクレオチドについてのシーケンシングの精度は比較的高くなる可能性があるが、ゲノム中の大多数のヌクレオチドは、個々のゲノムが一度だけシーケンシングされた場合、かなりの数のシーケンシングエラーが存在することを意味する。例えば、1 塩基対当たり 0.2 % のエラー率と 400 塩基対のリード長さの場合、少なくとも 1 つのエラーを有するリードの割合は、 $1 - (1 - 0.002)^{400} = 0.551$ であり、これは 55 % を超える配列リードが少なくとも 1 つのエラーを有する可能性を意味する。したがって、シーケンシングエラーと希少な真の変異を区別するためには、個々のゲノムに多数回シーケンシングすることによりシーケンシング精度を高めることが望ましい。例えば、たとえ各配列リードが 1 % のエラー率を含んでいても、バリエーションの位置をカバーする 8 つの同一のリードの組合せは、エラー率が $(10^{-2})^8$ すなわち 10^{-16} の強く支持されるバリエーション検出を生成するのである。

【0029】

DNA シーケンシングの深度 (depth) は、シーケンシングプロセス中にヌクレオチドが読み取られる回数を意味する。ディープシーケンシングは、リードの総数が調査中の配列の長さよりも何倍も大きいことを示している。カバー率は、再構築された配列内のあるヌクレオチドを表すリードの平均数である。「ディープ (deep)」という用語は、7 倍超などの広範囲の深度について使用されており、用語「ウルトラディープ (ultra-deep)」は、一般に 100 倍超などのより高いカバー率を意味する。シーケンシング深度の要件は、バリエーションのタイプ、疾患モデル、及び関心領域の大きさに依存し得る。すなわち、1 % 以下のバリエーション頻度を有する希少バリエーションについては、より高いカバー率が所望される。大規模な平行 NGS は、真のバリエーション検出のためのそのようなウルトラディープシーケンシングを可能にする。それにもかかわらず、より短いリードの大きい深度を生成することは、必ずしも希少バリエーション検出に関する全ての課題を解決しない。

【0030】

II. ウルトラディープシーケンシングにおけるバリエーションコール

バリエーションコーリングは、試験試料と参照配列の配列リード間の真の差を識別するプロセスである。バリエーションコーリングは、試料の特性評価及び疾患の診断において重要である。しかし、非常に低い頻度でしばしば体細胞バリエーションが発生するため、バリエーションコーリングは本質的に難しい。バリエーションコーリングの 1 つの目標は、謝った偽陽性を最小にするために高い信頼度で体細胞バリエーションを同定することである。

【0031】

図 1 は、標的化ウルトラディープシーケンシングのための次世代シーケンシング (NGS) を用いる、ゲノムシーケンシング及びバリエーションコーリングの方法 100 を示す。他の方法と同様に、実施態様は、記載された工程のすべて又は一部を含むことができ

10

20

30

40

50

、いくつかの工程はコンピュータシステムを用いて行うことができる。方法 100 の結果は、生物の診断を決定する際に医師によって使用することができる。

【0032】

ブロック 110 において、シーケンシングされ、かつ診断されるポリヌクレオチドを含む試料は受け取られ、ここで、前記ポリヌクレオチドは、シーケンシングされるべき標的領域を潜在的に含む。上記で定義したように、用語「試料」は核酸を含むか又は含むと推測される任意の組成物を指す。試料は、そこから試料が得られる生物のゲノムに由来する核酸分子を含む。例えば試料は、染色体中にコードされたゲノムを含有する細胞を含むことができる。試料は、1 又は 2 以上の試験試料を含むことができる。試料はまた、1 又は 2 以上の参照試料又は対照試料を含むことができる。いくつかの試料は、ゲノムの特定の領域における変異について試験されている患者から得ることができる。試料は、癌について試験されている腫瘍の生検から得ることができる。試料は、いくつかの正常細胞、癌進行の初期段階のいくつかの細胞、及び癌の進行の後期ステージのいくつかの細胞を含むことができる。試料は、異なる人や同じ人物（例えば、異なる生検）由来でもよく、異なる実験条件を用いてもよい。

10

【0033】

場合により、ブロック 120 において RNA 又は DNA は、シーケンシング前に試料から分離される。生物試料から核酸を単離するための方法は、例えば Sambrook に記載されるように公知であり、いくつかのキットは市販されており、例えば、DNA Isolation Kit for Cells and Tissues, DNA Isolation Kit for Mammalian Blood, High Pure FFPE DNA Isolation Kit, High Pure RNA Isolation Kit, High Pure Viral Nucleic Acid Kit, and MagNA Pure LC Total Nucleic Acid Isolation Kit があり、全てが Roche から入手可能である。いくつかの実施態様において、単離された核酸はゲノム DNA を含む。いくつかの実施態様において、単離された核酸は、循環遊離 DNA 断片 (circulating free DNA fragments) (cfDNA) を含む。いくつかの実施態様において、単離された核酸は、細胞性 mRNA 又は cfRNA などの RNA を含む。

20

【0034】

RNA の場合、ブロック 130 において、逆転写反応が行われる。例えば RNA は、逆転写酵素を用いて相補的 DNA (cDNA) に変換することができる。

【0035】

30

場合により、ブロック 140 において、シーケンシングのために DNA セグメントを調製することができる。これは DNA を、標的領域を含むより小さな DNA セグメントに断片化し、DNA セグメントの末端にアダプター配列を連結し、そして DNA 断片が由来する試料を同定する固有のバーコード配列を固定することを含むことができる。標的領域は、例えば任意の癌関連変異があるかどうかを調べるための、診断関連性を有する可能性のある DNA 中のセグメントである。例として標的領域は、ほぼ数百塩基、例えば 150 ~ 250 塩基、150 ~ 400 塩基、又は 200 ~ 600 塩基であることができる。別の実施態様において、標的領域に対応するゲノムセグメントを捕捉するためにプローブを使用することができる。例えば、標的領域にハイブリダイズするように設計されたプローブを、表面上に配置することができる。次にゲノムセグメントをその表面の上に配置することができ、標的領域のセグメントが優先的にハイブリダイズされ得る。試料の DNA は、例えば超音波処理又は他の適切な方法によって断片化して、より小さなゲノムセグメントを得ることができる。例えば、200 ~ 500 塩基長さのゲノムセグメントを得ることができる。特定のシーケンシング操作について、ほぼこの長さのゲノムセグメントが好ましい。しかし実施態様は、任意の長さのゲノムセグメントを使用することができる。

40

【0036】

ゲノムセグメントは、バーコードやマルチプレックス識別子 (MID) 配列でマークすることができる。例えば 10 塩基の配列を、ゲノムセグメントの末端を（例えば、リガーゼを用いて）加えることができる。このように、種々の試料からのセグメントは、単一のシーケンシングラン中に並行してシーケンシングすることができる。MID は配列リ

50

ードの一部として読み取ることができ、同じM I Dを有する配列リードは同じ試料に起因し、一緒に分析することが可能である。M I Dは、異なる試料から配列リードを脱多重化又は区別するために使用することができる。

【0037】

ブロック150においてDNAセグメントは、PCR、SDA、及びこれらの派生方法などの増幅法により場合により増幅又は増加させて、DNAセグメントすなわちシーケンシングのための増幅産物を生成することができる。Taqポリメラーゼ又は他の耐熱性ポリメラーゼなどのDNAポリメラーゼを、PCRによる増幅のために使用することができる。例えば、増幅法の総説については、Fakruddin et al., J Pharm Bioallied Sci. 5:245 (2013)を参照されたい。これらの増幅産物は、増幅に使用されるプライマーに基づいて規定される。プライマーは、核酸上の標的領域に特異的である。シーケンシングプライマーが増幅産物内の配列に特異的に特異的である（特異的にハイブリダイズする）ように、シーケンシングプライマーは典型的には、増幅プライマーの選択に基づいて設計される。いくつかの実施態様において標的領域は、標的濃縮工程によって濃縮することができる。増幅及び濃縮プロセスの両方を実行することができる。フォワードプライマー及びリバースプライマーは、標的領域を増幅するために使用することができる。これらのフォワードプライマー及びリバースプライマーは種々の長さ、例えば約15～30塩基長のものでよい。

10

【0038】

いくつかの実施態様において、試料特異的M I Dの添加は、異なる時点で発生し得る。例えばM I Dは、増幅／濃縮後に添加することができ、次に試料と一緒に混合される。こうして、異なる試料は、異なる標的領域について増幅又は濃縮することができるであろう。

20

【0039】

ブロック160において、1又は2以上の試料からのDNAセグメントは、単一のシーケンシングランで大規模に並列様式でシーケンシングされる。シーケンシングプロセスにおいて、増幅過程で作成された同じセグメントのクローンは、別々に決定された配列を有することができる（及び後にカウントされる）。いくつかの実施態様において単一のシーケンシングランは、1テラ塩基（terabase）を超えるデータを生成することができる。いくつかの実施態様において、1試料当たり約3,000リード超を得ることができる。リードの数は、試料のサイズ、標的増加の一部としてどの程度の増幅が行われるか、及びシーケンシングプロセスのバンド幅（すなわち、どの程度のシーケンシングに対して装置が設定されるか、例えばいくつかのビーズが使用されるか）に依存してもよい。ある実施態様において、リードは約150～250塩基長である。

30

【0040】

シーケンシングプロセスは、Roche 454, Illumina GA, 及び ABI SOLiDなどの種々のNGSプラットフォーム上の種々の技術によって行うことができる。ある実施態様においてDNAセグメントは、シーケンシングの一部として増幅を受けることができる。増幅プロセスが標的増加試料を作成するために使用される実施態様において、この増幅は第二の増幅工程であろう。第二の増幅は、第二の増幅が行われなかった場合よりも、強いシグナル（例えば、特定の塩基：A、C、G、又はTに対応する蛍光シグナル）を提供することができる。

40

【0041】

シーケンシング処理の一例において、ブロック150からの増幅されたセグメント（例えば、増幅が溶液中で発生した場合）は、それぞれビーズに付着させることができる。付着したセグメントは、次にビーズ上で増幅することができ、各ビーズから1つの配列リードを得ることができる。表面を使用する実施態様において、セグメントを表面に付着させ、次に表面上で単一のクラスタを作成することができる。各クラスタについて単一の配列リードを得ることができる。配列リードは、ゲノムセグメントの全長又はセグメントの一部についてのものであることができる。

50

【 0 0 4 2 】

ブロック 1 7 0 において、場合により配列リードは濾過されて、低クオリティリードと短いリードが除去され、残りの配列リードは参照配列の標的領域にアラインメントされる。いくつかの実施態様において、同一の塩基を有するリードは、単一の配列リードと見なされるように組み合わせられる。したがって、唯一のユニークリード (unique read) に関するリード数を記録することができる。平均塩基スコアは、全てのユニークリードに関する全ての塩基位置において計算され得る。塩基スコアは、塩基コールが配列リード上で如何に正確であるかを測定することができる。塩基スコアを使用して、低クオリティリードを除去することができる。いくつかの実施態様において、最小値よりも短いリードも同様に除去される。

10

【 0 0 4 3 】

アラインメントすることにより、本方法は配列リードを参照配列の標的領域と比較して、配列リードと参照配列との間の変化の数を決定することができる。アラインメントは 1 又は 2 以上の標的領域にのみ特異的であることができ、ゲノム全体を検索する必要はないため、アラインメントは高速であることができる。また標的領域に対応するセグメントの割合が増加するにつれて、かなりの数のリードが標的領域に良好に一致するであろう（例えば、比較的少数の変化）。

【 0 0 4 4 】

ある実施態様において、複数の標的領域が使用される場合、配列リードは複数の標的領域のすべてと比較することができ、最良のアラインメントを提供する標的領域を同定することができる。異なる標的領域は、異なる遺伝子又は遺伝子を有する異なるエクソンを有することができる。したがって、最良のアラインメントを有するエクソンが同定され得る。

20

【 0 0 4 5 】

バーコード又は M I D が使用される場合、それはアラインメント前に除去され得る。特定の試料に関する全てのリードを 1 つのグループに構成するために、M I D を使用することができる。このようにして、他の試料からの変異は、特定の試料の分析には影響しないであろう。このグループ化は、脱多重化 (de-multiplexing) と呼ばれる。異なる試料は異なる標的領域を有することができるため、アラインメントについて参照配列のどの標的領域を比較すべきかを決定するために M I D は使用され得る。

30

【 0 0 4 6 】

ブロック 1 8 0 において、標的領域からのアラインメントされた配列リードは、標的領域における変異を同定するために使用される。この工程の一部として、バリエーションの数（又はバリエーション数）、参照アレルの数（又は野生型数）、従って各試料に関する配列位置における各バリエーションの頻度を決定することができる。例えば、標的領域内の特定の位置について、通常の A の代わりに G 変異が現れる回数をカウントすることができる。G 変異が見られる回数の割合は、その位置にアラインメントされる全リードから決定することができる。いくつかの実施態様において、一緒に発生する変化を識別することができる、同じ変異の一部として分類することができる。各試料について、標的領域のシーケンシング深度は、その試料について任意のフィルターを通過するリードの数から決定することができる。

40

【 0 0 4 7 】

ブロック 1 9 0 において、バリエーション数、野生型数、及び / 又はバリエーション頻度に基づいてバリエーションコーリングが行われる。ある実施態様において、特定のバリエーションのバリエーション頻度は、実際の変異と見なすためには、閾値（存在フィルタ (abundance filter)）よりも大きいことが必要とされ得る。表 1 は、Illumina MSR 体細胞変異コーラー (caller) がデフォルト設定で報告するポアソン (Poisson) モデルに基づいて計算される最少バリエーション数とバリエーション頻度を示す。

【表 1】

深度	報告される最小数	報告される最小%
100	5	5
200	7	3.5
500	12	2.4
1000	19	1.9
2000	32	1.6
5000	68	1.36
10000	125	1.25
20000	235	1.175
50000	554	1.108
100000	1075	1.075

表 1. 種々の深度について報告される最小バリエーション数と最小バリエーション頻度

【0048】

いくつかの実施態様において、バリエーションが実際に試料中に存在する信頼水準を示すクオリティスコアが提供され、バリエーションコールを行うために使用される。いくつかの実施態様において、クオリティスコアは、バリエーション数、野生型数、及び/又はバリエーション頻度の1又は2以上と組み合わせて使用して、バリエーションコールを行うことができる。医師は同定された変異を使用して、癌の素因を診断するか又は癌を有するとして腫瘍を同定することができる。

【0049】

図2は、参照配列210と比較した試験試料中の標的領域215の配列リードの例を示し、ここで、種々の配列位置における同一のクラスと異なるクラスのバリエーションが示される。図2は、参照配列が塩基Aの参照アレルを有する標的領域内の4つの位置の例を示す。例示を容易にするために5つの配列リードが明示的に示されるが、実際にはより多くのリードが使用される。参照配列210は、位置205、231、255、及び281においてAを有することが示される。

【0050】

位置205について、いくつかの配列リードにおいてAが検出されるが、いくつかの配列リードではCが検出される。Cの検出は、潜在的なA>Cバリエーションを示す。バリエーションA>Cは、特定のバリエーションクラスのものである。他の塩基は、示されていない配列リード中で検出され得る。他のアレルの存在は、他のバリエーションクラスの他のタイプのバリエーションを示す可能性がある。

【0051】

位置231について、いくつかの配列リードにおいてAが検出される；いくつかの配列リードではCが検出される；更にいくつかの他のリードにおいて、位置231において何も検出されない（「0」）。Cの検出は、単一塩基置換A>Cの潜在的なバリエーションを示す。「0」の検出は潜在的な欠失のバリエーションを示す。

【0052】

位置255における塩基Aについて、いくつかの配列リードにおいてAが検出される。いくつかの配列リードではCが検出される；しかし、いくつかの他のリードではTが検出される。Cの検出は、単一塩基置換A>Cの潜在的なバリエーションを示す。Tの検出は、異なる単一塩基置換A>Tの潜在的なバリエーションを示す。

【0053】

位置281における塩基Aについて、いくつかの配列リードにおいて、Aが検出される；いくつかの配列リードにおいて、Cが異なる頻度で検出される。異なる頻度でのCの検出は、異なるバリエーション頻度を有する単一塩基置換A>Cの潜在的なバリエーションを示す。

【 0 0 5 4 】

試験試料についての配列リードに基づいて、それぞれの位置について、野生型塩基 A の数、単一塩基置換 A > C の数、単一塩基置換 A > T の数、及び A の欠失の数をカウントすることができる。図 2 に示されるバリエーションのタイプは、例示のみが目的である。本開示において後述されるように、種々のタイプのバリエーション又は変異が存在し得る。

【 0 0 5 5 】

I I I . 統計分布モデルに基づくバリエーションコーリング

本開示のいくつかの実施態様において、NGS 実験で観察されるすべてのバリエーションを報告することができる。低頻度の真陽性を偽陽性から区別するために、ほとんどの観測される低頻度バリエーションは偽陽性であってもよい。偽陽性バリエーションの分布を用いて、バリエーションコーリングクオリティスコアを確立して、バリエーションが真陽性である可能性を決定することができる。

【 0 0 5 6 】

A . 統計モデルに基づくバリエーションコーリングの数学的理論

図 3 A ~ 3 D は、本発明のいくつかの実施態様に従う統計モデルに基づくバリエーションコーリングの基礎となる数学的理論を提供する。バリエーションコーリングの偽陽性率は配列状況や位置に依存しているため、全ての試料中の種々の位置における A > C などの同一のクラス又はタイプのバリエーションと一緒に比較して、統計分布に基づくバリエーションコールを行うことができる。

【 0 0 5 7 】

いくつかの実施態様において、異なる配列位置でのシーケンシングランにおける単なるバリエーションは、20 のクラスに分けることができる。すべてのクラスにおいて、バリエーションの大部分は偽陽性である。すべてのバリエーションクラスの統計分布のパラメータを計算することができる。バリエーションクラスは以下のように定義することができる：

- (1) A > C、A > G、A > T、C > A、C > G、C > T、G > A、G > C、G > T、T > A、T > C、及び T > G を含む 1 2 の単一塩基置換；
- (2) A C > G A などの多塩基置換；
- (3) A G T > A T 又は G C A T > G T などの 1 ~ 2 塩基の欠失；
- (4) A T C G A > A A などの 3 塩基の欠失；
- (5) G A C C T A > G A 又は T G C G C G A > T A などの 4 ~ 5 塩基の欠失；
- (6) A T C C T C A G > A G などの 6 塩基以上の欠失；
- (7) A T > A A T 又は G C > G T A C などの 1 ~ 2 塩基の挿入；
- (8) G C > G T A A C 又は A C > A G A T G C などの 3 塩基以上の挿入；そして
- (9) 単一塩基置換 A > C などの他の単純な変異に、すぐ続く 1 塩基の欠失、例えば、元々の参照塩基が A T であり、変異塩基が C である、すなわち A T > C。そのような変異 A T > C はまた、A の欠失に単一塩基置換 T > C が続くと解釈することができる。

【 0 0 5 8 】

本明細書において単純な突然変異は、その中に一致する塩基対無しで、2 つの一致する塩基対によって結合された変異である。例えば a A T g や a C g において、単純な変異 A T > C は、一致する対 a - a と一致する対 g - g により結合され、ここで、小文字は一致する対について使用される。しかし a A c G g と a C c T g において、A c G > C c T は、この中に一致する対 c - c が存在するため単純な変異ではない。このように、A c G > C c T は、2 つの単純な変異 A > C と G > T からなる複合変異である。

【 0 0 5 9 】

いくつかの実施態様において、1 又は 2 以上の試料について参照配列中に参照アレル（例えば、バリエーションタイプ A > C についての塩基 A）が存在する標的領域中の種々の位置における、同一のクラスのバリエーション、例えば A > C のバリエーション頻度を用いて、バリエーションクラスについて統計分布を作成することができる。例えば図 2 に示すように、位置 205、位置 231、位置 255、位置 281 のそれぞれ、及び試料の参照配列内に A が存在する標的領域内の他の位置における単一塩基置換 A > C のバリエーション頻度は、バリエア

10

20

30

40

50

トクラス $A > C$ に関する統計分布のためのデータ点であってもよい。位置 205、位置 231、位置 255、位置 281 のそれぞれ、及び試験試料と同一のシーケンシングランでシーケンシングされる他の試料のそれぞれの参照配列内に A が存在する標的領域内の他の位置における単一塩基置換 $A > C$ のバリエーション頻度は、バリエーションクラス $A > C$ に関する統計分布のためのデータ点であってもよい。一方、位置 205、位置 231、位置 255、位置 281 のそれぞれ、及び試験試料と同一のシーケンシングランでシーケンシングされる各試料の参照配列内に A が存在する標的領域内の他の位置における単一塩基置換 $A > T$ 又は単一塩基欠失 $A > 0$ などの異なるバリエーションクラスのバリエーション頻度は、バリエーションクラス $A > C$ について統計分布のために使用されない。

【0060】

いくつかの実施態様において、少なくとも 30 のデータ点が統計分布に含まれる。少なくとも 30 のデータ点は、単一のシーケンシングラン中に 2 以上の試料からののものであってもよい。30 未満のデータ点の場合は、真の分布はデータ点によって表されない場合がある。

【0061】

図 3 A は、同一のクラスの変種についてのバリエーション頻度の理想的な統計分布（正規分布）を示す。図 3 A は例示のみが目的である。バリエーションクラスの変種頻度の実際の統計分布は試料に依存することがあり、二峰性分布のような他の分布形態であってもよい。いくつかの実施態様において、バリエーション頻度の二乗、平方根又は対数のような変換のいくつかの形態は、正規分布に近い分布を形成することができる。

【0062】

図 3 A において、 x 軸はバリエーションクラスの変種頻度値を示し、 y 軸は特定のバリエーション頻度値 f を有するデータ点の数を示す。図 3 A 中の理想的な正規分布に示されるように、平均値 m と標準偏差 s は分布に基づいて決定することができる。

【0063】

図 3 B は、特定の試料上の特定の位置における変種の変種頻度が所定の Z 値を有する確率を示し、ここで、この確率と Z 値は、図 3 A に示された統計分布に由来することができる。いくつかの実施態様において図 3 B は、平均と標準偏差とに基づく図 3 A の正規化された分布であってもよい。いくつかの実施態様において、より複雑な変換又は転換、例えば対数変換を使用することができる。図 3 A 中の斜線部は、 z に等しいか又は z より大きい全ての Z 値の累積確率を示す。

【0064】

図 3 C は、 Z 値が z 未満又はこれに等しい値を取る確率の累積分布関数 F を示す。

【0065】

図 3 D は、特定のバリエーション頻度値又は Z 値を有する変種が、左の主軸上で偽陽性である塩基コーリング誤り確率（ p 値）と、右の 2 次軸上でバリエーションコールを行うための関連するクオリティスコア Q を示す。いくつかの実施態様において p 値は、 $1 - F$ によって計算することができる。いくつかの実施態様においてクオリティスコア Q は、 $Q = -10 \log_{10} p$ によって与えられるフレッド（Phred）クオリティスコア、又はフレッドクオリティスコアの任意の変形であってもよい。

【0066】

B．統計モデルに基づくバリエーションコーリングの方法

図 4 は、統計モデルを用いるバリエーションコーリングの方法 400 を示す。他の方法と同様に実施態様は、記載された操作の全て又は一部を含むことができ、いくつかの操作は追加の操作又はサブ操作を含むことができる。

【0067】

ブロック 410 において、単一のシーケンシングランで 1 又は 2 以上の試料中の標的領域を標的化する配列リードが受け取られる。配列リードデータは受け取られ、読み取り可能なフォーマットで保存され、コンピュータで解析することができる。いくつかの実施態様において、低クオリティリード又はアダプター配列を除去するために、配列リードデ

10

20

30

40

50

ータの予備処理を実行することができる。いくつかの実施態様において、バーコード又は M I D を除去することができ、同じ試料からの配列リードは、標識又はグループ化されてもよい。

【 0 0 6 8 】

例えば方法 1 0 0 のブロック 1 7 0 に記載されるように、ブロック 4 2 0 において配列リードは参照配列の標的領域にアラインメントされる。

【 0 0 6 9 】

ブロック 4 3 0 において、試験試料のアラインメントされた配列リード上の特定の配列位置における同じバリエーションクラスのバリエーションアレルは、同定され、カウントされて、バリエーション数を決定することができる。試験試料のアラインメントされた配列リード上の特定の配列位置についてのリード数も、同様に決定することができる。例えば図 2 に示すように、試験試料の配列リード中の位置 2 0 5 における C の総数は、位置 2 0 5 におけるバリエーションクラス A > C についてのバリエーション数であり、試験試料の配列リード中の位置 2 0 5 におけるリードの総数は、位置 2 0 5 におけるバリエーションクラス A > C についてのリード数である。いくつかの実施態様において、試験試料についての特定の位置のリード数は、別の操作で決定することができる。

【 0 0 7 0 】

ブロック 4 4 0 において、特定の位置における同一のクラスのバリエーションのバリエーション頻度が決定される。ある実施態様においてバリエーション頻度は、試験試料中の特定の位置におけるバリエーション数をリード数で割ることにより決定することができる。別の実施態様においてバリエーション頻度は、試験試料中の特定の位置においてバリエーション数を非バリエーション数（例えば、リード数 - バリエーション数）で割ることによって決定することができる。当業者であれば、使用することができるバリエーション頻度の種々のタイプの形態を理解し得る。

【 0 0 7 1 】

ブロック 4 5 0 において、試験試料として同一のシーケンシングランでシーケンシングされた各試料について、同一のクラスのバリエーション、例えば A > C は、バリエーションクラスの参照アレル、例えば A が、標的領域における参照配列上に存在する複数の位置のそれぞれの位置で、同定されカウントされる。同様に、同一のシーケンシングランにおける各試料について、参照配列上でバリエーションクラスの参照アレルが見つかる複数の位置のそれぞれに関するリード数を決定することができる。

【 0 0 7 2 】

ブロック 4 6 0 において、試験試料として同一のシーケンシングランでシーケンシングされた各試料について、複数の位置のそれぞれの位置における、同一のバリエーションクラス、例えば A > C のバリエーション頻度は、各位置についてのバリエーション数をその位置のリード数により割ることによって決定することができる。すなわち、もし、例えば 3 つの試料がシーケンシングランで一緒にシーケンシングされ、標的領域内の参照配列上の 3 0 の位置が、バリエーションクラスについての参照アレルを有する場合、各試料上の各位置について 1 つで最大 9 0 のバリエーション頻度を計算することができる。これらのバリエーション頻度を用いて、同一のシーケンシングランで同一のクラスのバリエーションについてバリエーション頻度の統計分布を決定することができる。他のバリエーションクラスのバリエーション頻度は、統計分布を決定するために含まれていないことに注目されたい。更に、分布モデルの正確性に影響を与える可能性のあるシーケンシングラン間の変動の影響を低減するために、他のシーケンシングランから得られたデータ点は、統計分布を決定するために含まれない。

【 0 0 7 3 】

ブロック 4 7 0 において、試験試料中の特定の位置における同一のクラスのバリエーションに関するバリエーション頻度に対応する確率値は、バリエーション頻度をブロック 4 6 0 において形成された統計分布のパラメータと比較することにより決定される。いくつかの実施態様において、確率値は実際の確率、累積分布、又はクオリティスコアでもよい。いくつかの実施態様において、統計分布のパラメータは、平均値と標準偏差の 1 又は 2 以上でもよい

10

20

30

40

50

。

【 0 0 7 4 】

ブロック 4 8 0 において、試験試料上の特定の位置におけるバリエーションクラスの変異が真陽性か否かを決定するために、確率値と閾値に基づいてバリエーションコールが行われる。いくつかの実施態様において、閾値は単一の値でもよい。いくつかの実施態様において閾値は、例えばバリエーション頻度の関数でもよい。いくつかの実施態様において閾値は、サポートベクターマシン (SVM) などのマシン学習アルゴリズムを用いて、トレーニングデータセットに基づいて決定することができる。いくつかの実施態様において、閾値は、異なるシーケンシングランから得られるトレーニングデータに基づいて決定することができる。

10

【 0 0 7 5 】

上記した方法は、以下の例に照らしてより良く理解することができる。

【 0 0 7 6 】

C. 例

以下の例は、このセクションで上記した方法を示す。以下の例において、バリエーション頻度が正規分布ではなく、一方で対数バリエーション頻度の分布が後述されるように正規分布に近い場合、バリエーションクラスに関する対数バリエーション頻度の統計分布に基づくモデルが使用される。

【 0 0 7 7 】

表 2 は、エクソン 20 の置換 T 7 9 0 M (2 3 6 9 で C > T) とエクソン 21 の置換 L 8 5 8 R (2 5 7 3 で T > G) の偽陽性を有する野生型データに関する元々のバリエーション頻度 f とその対数変換 x に適用される、Lilliefors 検定及び Shapiro-Wilk 検定などの正規性検定の結果を示す。この結果は、正規分布の仮定が使用される時、 x が観察された試料の結果を得るためのより大きな確率 (P 値) (> 0.08) を有することを示し、これは、実際の分布と正規分布との間のより小さな相違を示し、 f がより小さい P 値 (< 0.016) を有することを示す。従って、 x は f よりも正規分布に近い。

20

【表 2】

偽陽性	変数	Lilliefors 検定の P 値	Shapiro-Wilk 検定の P 値
T790M	f	0.008805	0.001830
	$\log_{10}(f+1e-06)$	0.348639	0.084104
L858R	f	0.014024	0.015862
	$\log_{10}(f+1e-06)$	0.602277	0.520155

30

表 2. f 及び x に関する正規性検査の P 値

【 0 0 7 8 】

変換されたバリエーション頻度は、ほとんどノイズについて元のバリエーション頻度よりも正規分布に近い場合、統計分析を行うのに通常の近似を使用するために、最初にバリエーション頻度の対数変換が行われる。いくつかの実施態様において、 $f = 0$ のときの負の無限大値を避けるために、以下の対数変換が使用される。

40

【数 1】

$$x = \log_{10}(f + e)$$

【 0 0 7 9 】

ここで、 e は負の無限大値を回避するための調整定数である。調整定数 e は、任意の適切な値に設定することができる。例えば、いくつかの実施態様において、 e は 10^{-6} に設定することができ、従って最小の x 値は -6 である。

50

【0080】

対数変換した後、正規分布近似の平均値 m と標準偏差 s を計算することができる。次に正規分布の近似を用いて、配列の位置で検出されるバリエーションの確率値を計算することができる。例えば、バリエーション頻度 f_1 、対数バリエーション頻度 $x_1 = \log_{10}(f_1 + e)$ 、十分な深度（総リード数）を有する、ある位置におけるバリエーションクラス中の観察されたバリエーションに関して、統計的確率値 z スコアは、以下によって計算することができる：

$$z = (x_1 - m) / (s / \sqrt{n})$$

10

【0081】

ここで、 n は、 s と m の推定に用いられる参照データ点の数である。計算結果は、 z スコアが大きい n に対して大きいことを示し、これは非常に小さい塩基コーリングエラー確率（ p 値）を生成し、したがって非常に大きなクオリティスコアを生成することができることを示す。したがっていくつかの実施態様において、 z スコアは z 様スコアで置換することができ、これは、上記式中の n を $\min(n, N)$ で置換することにより計算される。 N は、任意の適切な値に設定することができる。いくつかの実施態様において、 N は 36 に設定される。いくつかの実施態様において、下限 s_2 はまた、 s が小さすぎる状況では、 $s / \sqrt{\min(n, N)}$ に設定することができる。 s_2 は、例えばデフォルト値 0.01 のような任意の適切な値に設定することができる。したがって、いくつかの実施態様において、 z 様スコアは以下により表すことができる。

20

【数3】

$$z' = (x_1 - m) / \max(s_2, s / \sqrt{\min(n, N)})$$

【0082】

z スコア又は z 様スコア z' を用いて、塩基コーリングエラー確率 p 値は、 $p = 1 - F(z)$ 又は $p = 1 - F(z')$ により決定することができ、ここで F は標準正規分布の累積分布関数である。次にバリエーションコーリングクオリティスコア Q_{AMP} は、フレッドスコアを用いて決定することができる。いくつかの実施態様において、 Q_{AMP} はフレッド様スコア：

30

【数4】

$$Q_{AMP} = -10 \log_{10}(\max(p, \min P))$$

【0083】

として定義することができる。ここで、 $\min P$ は $10^{-\max Q/10}$ である。 $\max Q$ は任意の適切な値に設定することができる。例えばいくつかの実施態様において、 $\max Q$ は、80 又は 130 に設定してもよい。

40

【0084】

いくつかの実施態様において、クオリティスコアを計算するために、試料平均と試料標準偏差の代わりに、データの中心位置と変動のロバスト推定 (robust estimations) を使用することができる。

【0085】

いくつかの実施態様において、線形カーネル (linear kernel) を有するサポートベクターマシン (SVM) などの分類法は、既知の真陽性と偽陽性とを有するトレーニングデー

50

タセットを用いて、偽陽性から真陽性を分離するために使用することができる。いくつかの実施態様において、閾値はデータを視覚化することによって設定することができる。

【0086】

図5及び図6は、実際の試料から配列リードデータに適用された上記方法の結果の例を示す。図5は、SVMによって決定されるセパレーターラインを有するエクソン20のEGFR T790M(2369でC>T)の異なるトレーニングデータと試験データに関する、最大対N=4を有するバリエーションコーリングクオリティスコア Q_{AMP} を示す。図5は、バリエーションと野生型データが十分に分離されていないことを示し、従って、0.1%以下のバリエーション頻度で真の変異と偽陽性を区別することは困難であるかも知れない。しかし、0.5%以上のバリエーション頻度を有するすべての試験データと少なくとも0.2%のバリエーション頻度を有するほとんどの試験データに関しては、真の陽性と偽陽性は正確に区別することができる。

10

【0087】

図6は、バリエーションは、SVMによって決定されるセパレーターラインを有するエクソン21のEGFR L858R(2573でT>G)のトレーニングデータに関するmax N=4を有するバリエーションコーリングクオリティスコア Q_{AMP} を示す。図6は、0.1%のバリエーション頻度を有するものを含むすべての試験データが、正しく分類することができることを示す。

【0088】

IV. 1又は2以上の参照試料との比較を用いる、特定の位置における特定のバリエーションに関するバリエーションコーリング

20

本発明のいくつかの実施態様において、異なる試料中の同一の位置におけるバリエーション及びバリエーションの野生型数を比較して、バリエーションコールを行うことができる。この方法は、シーケンシングランにおいて陰性対照として野生型(通常は正常)試料が利用可能である場合に特に有用である。

【0089】

A. 特定の位置における特定のバリエーションを検出するために試験試料と参照試料とを比較する方法

この方法は、異なる試料について特定の位置における特定のバリエーションを比較するために使用することができ、2つという少ないデータ点に適用することができる。

30

【0090】

図7は、参照試料と試験試料の配列リード上の特定の位置112における特定のバリエーションC>Tを示す。図7に示されるように、参照配列の位置112における参照アレルはCであり、そして参照試料の位置112における配列リードはほとんどCであるが、シーケンシングエラーに起因するバリエーションC>Tを有することがある。試験試料について、位置112における配列リードは、低いバリエーション頻度に起因してCであり、真の変異に起因していくつかのTであり、及びシーケンシングエラーに起因していくつかのTでもよい。

【0091】

参照試料は、理論的には真の変異は無いが、図7に示すようにシーケンシングエラーが小さいバリエーション数を引き起こす可能性がある。特定の位置における特定のバリエーションに関する参照試料と試験試料のバリエーション数及び参照試料と試験試料の野生型数は、配列リードに基づいて決定することができ、以下の表3に示される表に入れられる。試験試料及び参照試料のカウントデータは、試験試料のバリエーション頻度が、同一の位置における参照試料のバリエーション頻度よりも有意に大きいかどうかを決定するために使用することができる。

40

【0092】

表3において、a1は参照試料中の特定の位置における特定のバリエーションの数であり、n1は参照試料の配列リードの深度であり、w1 = n1 - a1は参照試料中の特定の位置における野生型数を示す。a2、n2、及びw2は、試験試料に関する対応するバリエーション

50

ト数、深度、及び野生型数である。表 3 はまた、行の合計 $a = a_1 + a_2$ 、 $w = w_1 + w_2$ 、及び総数 $n = n_1 + n_2$ を列記する。

【 0 0 9 3 】

【表 3】

カウントテーブル	参照試料	試験試料	行の合計
バリエーション数	a_1	a_2	a
野生型数	w_1	w_2	w
総数（深度）	n_1	n_2	n

表 3. 参照試料と試験試料のカウントテーブル

10

【 0 0 9 4 】

(a_1, w_1) と (a_2, w_2) の比に有意差があるかどうかを試験するために、いくつかの方法がある。いくつかの実施態様において、 n_1 及び n_2 はウルトラディープシーケンシングにおいて非常に大きくなる可能性があるため、片側カイ二乗検定を用いることが好ましい。片側カイ二乗検定では、比率 $f_1 = a_1 / n_1$ と $f_2 = a_2 / n_2$ が最初に計算される。もし $f_2 < f_1$ の場合、すなわち、試験試料の割合が参照試料の割合より大きくない場合（偽陽性であることが知られている）、エラー率 $p = 0.63$ に対応する 2 などの非常に小さいクオリティスコアを設定することができ、更なる分析を必要としない。しかし、 $f_2 > f_1$ の場合は、カイ二乗統計値 (chi-squared statistic) は次のように計算することができる：

20

【数 5】

$$\chi^2 = n \times (a_1 \times w_2 - a_2 \times w_1)^2 / (n_1 \times n_2 \times a \times w)$$

【 0 0 9 5 】

片側バリエーションコーリングエラー確率 p 値は、 $p = 0.5 \times (1 - pchisq(\chi^2, d))$ として計算することができ、ここで $pchisq$ は、自由度 d を有するカイ二乗累積分布関数である。いくつかの実施態様において、自由度 d は 1 である。

30

【 0 0 9 6 】

(a_1, w_1) 及び (a_2, w_2) の比率が有意に異なるかどうかを試験するための別の方法は、大規模な試料のためのピアソン比率検定 (Pearson proportion test) である。ピアソン比率検定において、2 つの比率、 $p_1_hat = a_1 / n_1$ と $p_2_hat = a_2 / n_2$ が最初に計算される。Z スコアは、

【数 6】

$$Z = (p_2_hat - p_1_hat) / \sqrt{V}$$

40

【 0 0 9 7 】

により計算することができ、ここで、 V は次の 2 つの式のうちの少なくとも一つを用いて計算することができる。

【数 7】

$$V = p1_hat \times (1 - p1_hat) / n1 + p2_hat \times (1 - p2_hat) / n2、及び$$

$$V = p_hat \times (1 - p_hat) \times ((1/n1) + (1/n2))$$

【0098】

ここで、 $p_hat = (a1 + a2) / (n1 + n2)$ である。次に片側 p 値は、 $p = 1 - p_{norm}(Z)$ として計算することができ、ここで p_{norm} は累積確率分布関数である。

10

【0099】

いくつかの実施態様において、フィッシャーの正確確率検定 (Fisher's exact test) を用いて、 $(a1, w1)$ 及び $(a2, w2)$ の比率が有意に異なるかどうかを決定することができる。フィッシャーの正確確率検定は、低幾何学的 (hypogeometric) 分布を用いる。フィッシャーの正確確率検定のための計算はより複雑になることがあり、大規模な試料についてオーバーフローを引き起こす可能性がある。

【0100】

p 値が計算された後、対応するクオリティスコアは $Q_{LOC} = -10 \times \log_{10}(p)$ として定義することができる。片側カイ二乗検定において、p は $(0, 0.5)$ の範囲内であることに留意されたい。いくつかの実施態様において、p が 0 に近い場合に、数値計算の困難さを回避するために、 $Q_{LOC} = -10 \times \log_{10}(\max(p, \min P))$ が使用され、ここで $\min P$ は任意の適切な値、例えば 10^{-13} に設定することができ、これは最大クオリティスコアを 130 に設定することと同等である。

20

【0101】

B. 特定の位置における特定のバリエントを検出するための参照試料を選択する方法

複数の試料のシーケンシングラン中の特定の位置における特定のバリエントに関する参照カウントを設定するために、様々な方法を使用することができる。一つの方法は、特定の位置における特定のバリエントの最も低いバリエント頻度と、最小値 $\min D$ 以上の深度を有する、同一のシーケンシングランにおいて、2つの試料のバリエント数の合計と深度の合計を使用する。いくつかの実施態様において、 $\min D$ は 3000 に設定することができる。いくつかの実施態様において、参照割合が $f0$ (これは、例えば 0.01 又は 1% に設定してもよい) より大きい時、全ての試料が特定の位置における特定のバリエントに関して高バリエント頻度を有するまれな可能性を回避するために、 $a1$ は $f0 \times n1$ に設定される。すなわち、使用される $a1$ 値は、実際の $a1$ 値か又は $f0 \times n1$ のいずれか小さい方である。この方法では、野生型試料にバリエントが混入している場合、バリエント混入 (variant contamination) を有する野生型試料は高バリエント頻度を示し、したがって特定のバリエントのための参照試料として選択されることはないであろう；従って、他の試料のクオリティスコア Q_{LOC} は通常は影響を受けない。いくつかの複雑な変異は、複数の単純な変異で構成されている。このような状況では、 Q_{LOC} は、複雑な変異のすべての単純な変異成分のクオリティスコア Q_{LOC} の中央値として定義することができる。

30

40

【0102】

既知の野生型試料はまた、参照試料として使用することもできる。しかし野生型試料にバリエントが混入している場合は、他の試料のクオリティスコア Q_{LOC} は小さくてもよい。

【0103】

C. 特定の位置における特定のバリエントを検出するために試験試料を参照試料と比較することによるデータ解析

図 8 は、特定の位置における特定のバリエントを分類するために試験試料を 1 又は 2 以上の参照試料と比較することによる、バリエントコーリングの方法 800 を示す。他の方

50

法と同様に、実施態様は、記載された操作の全て又は一部を含むことができ、いくつかの操作は追加の操作又はサブ操作を含むことができる。

【0104】

ブロック810において、単一のシーケンシングラン中の1又は2以上の試料からのDNAセグメント中の標的領域を標的化する配列リードが受け取られる。配列リードデータは受け取られ、読み取り可能な任意のフォーマットで記憶され、コンピュータにより解析することができる。いくつかの実施態様において、低クオリティリード又はアダプター配列を除去するために、配列リードデータの前処理が行うことができる。いくつかの実施態様において、バーコード又はMIDは除去してもよく、同一の試料からの配列リードは標識又はグループ化することができる。

10

【0105】

ブロック820において、配列リードは、方法100のブロック170に記載されたように、参照配列の標的領域にアラインメントされる。

【0106】

ブロック830において、アラインメントされた配列リード上の特定の配列位置における特定のバリエーションのバリエーションアレルは、アラインメントされた配列リードを参照配列と比較することにより、すべての試料について同定することができる。当業者に知られているように、任意の適切なアラインメント技術を使用することができる。

【0107】

ブロック840において、全ての試料について特定の配列位置における特定のバリエーションに関するバリエーション数とリード数を決定することができる。バリエーション数は、試料の異なる配列リード上の特定の位置における同じバリエーションアレルの、例えばA>CバリエーションのCの、合計数である。リード数は、試料の特定の位置のリードの総数である。

20

【0108】

ブロック850において、少なくとも1つの試料が参照試料として選択される。上記したように、いくつかの実施態様において、既知の野生型試料を参照試料として使用することができる。いくつかの実施態様において、シーケンシングラン中で最小バリエーション頻度を有する2つの試料を参照試料として使用することができる。このような実施態様において、最小バリエーション頻度を有する2つの試料のバリエーション数の合計とリード数の合計は、計算において参照試料のバリエーション数a1及びリード数n1として使用することができる。

30

【0109】

ブロック860において、セクションIV(A)で上記した方法を用いて、試験試料と参照試料のための特定の配列位置における特定のバリエーションのバリエーション数とリード数を比較して、確率値を決定される。確率値は、カイ二乗値、累積確率分布値、p値、Z値、及びクオリティスコアの1又は2以上であってもよい。

【0110】

ブロック870において、試験試料上の特定の位置における特定のバリエーションが真の陽性であるかどうかを決定するために、確率値と閾値とに基づいてバリエーションコールが行われる。いくつかの実施態様において、閾値は単一の値であってもよい。いくつかの実施態様において、閾値は例えばバリエーション頻度の関数であってもよい。いくつかの実施態様において、閾値はトレーニングデータセットに基づいて、例えばサポートベクターマシン(SVM)などのマシン学習アルゴリズムを用いて決定することができる。いくつかの実施態様において、閾値は、異なるシーケンシングランから得られたトレーニングデータに基づいて決定することができる。

40

【0111】

D. 例

以下の例は、特定の位置における特定のバリエーションを検出するための、試験試料を参照試料と比較することによるバリエーションコーリングの結果を示す。

【0112】

50

図 9 は、SVM によって決定されるセパレーターラインを有するエクソン 20 の EGF R T790M のトレーニングデータと試験データに関する局所化されたバリエーションコーリングクオリティスコア Q_{LOC} を示す。図 9 から、SVM により決定されるセパレータが、真陽性として T790M のコーリング 0.1% に低下されても、野生型の試験データの誤分類が存在しないことがわかる。更に、例えば $f = 0.1\%$ 又は $Q_{LOC} = 18$ の単一の閾値は、T790M バリエーションの良好な判断ポイントとすることができる。

【0113】

図 10 は、SVM によって決定されるセパレーターラインを有するエクソン 21 の EGF R L858R のトレーニングデータと試験データに関する局所化されたバリエーションコーリングクオリティスコア Q_{LOC} を示す。0.1% のバリエーション頻度を有するものを含むすべての試験データが、正しく分類されていることがわかる。更に例えば $Q_{LOC} = 18$ の単一の閾値は、L858R バリエーションの良好な判断ポイントとすることができる。

10

【0114】

図 11 は、SVM によって決定されるセパレーターラインを有するエクソン 19 の EGF R 15 塩基欠失 2235_2249 del 115 のトレーニングデータと試験データに関するバリエーションコーリングクオリティスコア Q_{LOC} を示す。0.1% のバリエーション頻度を有するものを含むすべての試験データが、正しく分類されていることがわかる。例えば $Q_{LOC} = 18$ 又は 20 の単一の閾値は、SVM を使用することなく、真の陽性と偽陽性を分離するように設定することができる。

20

【0115】

図 9 ~ 11 はまた、局所化されたバリエーションコーリングスコア Q_{LOC} が、モデルに基づくバリエーションコーリングスコア Q_{AMP} より、真陽性と偽陽性の間により広いマージンを有することを示す。

【0116】

V. 単純化されたクオリティスコアの推定

いくつかの適用において、すべてのバリエーションについて p 値とクオリティスコアを直接計算することは、時間がかかる。いくつかの実施態様において、クオリティスコアのみを整数として報告する必要があるため、 Q_{LOC} と Q_{AMP} の値を離散化することができる。例えば $f_2 = f_1$ である場合、クオリティスコアは 2 に設定することができる； $f_2 > f_1$ である場合、クオリティスコアは 3、4、...、又は $\max Q$ に設定することができ、これは、例えばいくつかの実施態様において 130 に設定することができる。

30

【0117】

いくつかの実施態様において、クオリティスコアは、例えば z 値 q_{chisq} 又は正規クオンタイル値 (normal quantile value) q_{norm} 、及び図 12 に示されるルックアップテーブルを用いて決定することができる。図 12 において、 $Q = 3.5, 4.5, \dots, 129.5$ について z 値と q_{norm} 値が計算され、ルックアップテーブルに示される。すなわち、二分探索アルゴリズムなどの探索アルゴリズムは、 z 値又は q_{norm} 値に基づいて、3、4、...、130 の最良近似整数値を決定するために使用することができる。

【0118】

VI. ゼロイベントを回避するために必要な試料量

40

血液検査の一つの実用的な問題は、低頻度変異を検出することができるように、バリエーションを検出するために十分な gDNA の量を決定することである。本発明のいくつかの実施態様において、ゼロイベント検出の確率を用いて、必要な試料の量を推定することができる。

【0119】

6.022×10^{23} / モルのアボガドロ定数、塩基対あたり 650 ダルトン (g/mol) の重量平均分子量、及びヒトゲノムあたり 3.096×10^9 塩基対に基づいて、1 ナノグラム (ng) のヒト gDNA は、 $6.022 \times 10^{23} / (650 \times 3.096 \times 10^9 \times 10^9) = 300$ 分子を含有すると計算される。

【0120】

50

変異を検出するために必要とされる g D N A の量は、変異頻度に依存し、ゼロイベントを回避する統計的問題を解決することによって決定することができる。例えば、Lachin, Biostatistical Methods: The Assessment of Relative Risks, p.19, Wiley (2000)を参照されたい。血液試料中の変異体コピー数が B であり、D N A コピーの総数が N であり、変異確率が $p = B / N$ であると仮定する。二項分布に従うと、ランダム試験において変異体コピーが得られない確率は $(1 - p)$ であり、N 回のランダム試験において変異体コピーが得られない確率は $(1 - p)^N$ である。従って以下の不等式を設定することができる。

【数 8】

$$(1 - p)^N \leq \alpha$$

10

【0 1 2 1】

ここで、 α は、変異が検出されない最大許容確率（最大許容失敗率）であり、 $1 - \alpha$ は、上側信頼限界である。したがって、試料サイズ N は、以下の不等式を解くことによって推定することができる。

【数 9】

$$N \geq \ln(\alpha) / \ln(1 - p)$$

20

【0 1 2 2】

$p \ll 1$ である稀な変異の場合、推定はテイラー展開

【数 1 0】

$$\ln(1 - p) \doteq -p$$

【0 1 2 3】

を用いて単純化することができ、そして試料サイズの推定値は以下の通りとなる。

【数 1 1】

$$N \geq -\ln(\alpha) / p$$

30

【0 1 2 4】

$-\ln(0.05) = 2.9957$ 、そして $-\ln(0.005) = 5.2983$ であるため、 $3 / p$ 又は $5.3 / p$ は、それぞれ 0.95 と 0.995 の上側信頼限界を持つ稀な変異の試料サイズ N を推定するために使用することができる。

40

【0 1 2 5】

表 4 は、最大許容失敗率が 0.05 と 0.005 を有する少なくとも一つの変異体コピーを含むのに必要な g D N A 分子の推定数を示す。例えば、試料中で少なくとも一つの変異体コピーを得るための 95 % の上側信頼水準 ($\alpha = 0.05$) を有する 0.1 % ($p = 0.001$) の変異を検出するためには、2995 の g D N A コピーが必要であり、これは約 $10 \ln g$ の g D N A 分子と同等である。

【0 1 2 6】

【表 4】

p	$\alpha=0.05$		$\alpha=0.005$	
	コピー数	ng	コピー数	ng
0.01	300	1.0	530	1.8
0.005	600	2.0	1060	3.5
0.002	1500	5.0	2650	8.8
0.001	2995	10.0	5300	17.7

表 4. 希少変異検出のための gDNA 分子のコピー数と重量の推定

10

【0127】

VII. 適用と検証

セクション III 及び IV において上記した方法は、判定基準として使用されるバリエーション頻度の閾値を決定する補助となり得る。この方法は、十分な入力 (input) DNA 量で 0.1 ~ 0.3 % の頻度で、置換をうまく検出することができる。偽陽性率は変異の状態と位置に依存するため、特定の位置における特定の置換について、0.03 % という低いバリエーション頻度を有するバリエーションを正しく検出することができる。

【0128】

20

適度なサイズの挿入、欠失、及び 15 塩基の欠失などの複雑な変異について、シーケンシングにおいてこれらのタイプの変異をランダムに発生させることは困難であり、誤差の主な原因は他の試料からのキャリーオーバー混入 (carry-over contamination) である。すなわち、ラン間で十分に確立された洗浄プロトコルを用いることにより、0.0025 % と低いバリエーション頻度を有するこれらのタイプのバリエーションを正確に検出することができる。

【0129】

Illumina MiSeq Reporter (MSR) は、本開示に記載された方法により検出される低頻度バリエーションを確認するための非標準的な方法で使用することができる。MSR は、組み込みポアソン (Poisson) モデルを用いる体細胞バリエーションコーラーを使用して、低頻度バリエーションを報告する。MSR が報告する最低頻度は、深度に依存する。ポアソンモデルに基づくと、MSR 体細胞バリエーションコーラーが報告する最低のバリエーション数と頻度を計算することができ、表 1 に示されるようにデフォルト設定される。例えば深度が 100 である場合、最低の報告される頻度は 5 % である；深度が 500 である場合、最低の報告される頻度は 1.36 % である；深度が更に大きい場合、最低の報告される頻度は上記の 1 % に近くなる。

30

【0130】

いくつかの実施態様において、既知のバリエーションを含む試料を参照試料として用いる MSR は、MSR が、野生型アレルを参照試料の「バリエーションアレル」として報告し、実際のバリエーションアレルを「野生型アレル」として報告するように実行することができる。こうして、本開示に記載された方法を用いるバリエーションコーリングを検証することができる。MSR のこの非標準的な使用は、いくつかの欠点を有する。第一に、これは既知のバリエーションを確認するためにのみ使用することができる。第二に、MSR が報告するバリエーションコーリングクオリティスコアは、実際のバリエーションのためというより野生型のものである。第三に、複数の重複する既知のバリエーションが存在する場合、この方法を使用することが面倒又は困難になる。しかし、上記の欠点を考慮した後、MSR は、既知のバリエーションのための検証ツールとして使用することができる。これは、ゲノム全体が参照配列として使用される場合、MSR マッピング/アラインメントソフトウェアが、マッピングされていないリードとして報告する適度なサイズのインデルのために特に有用である。

40

【0131】

50

V I I I . コンピュータシステムとシーケンシングシステム

本明細書に記載の任意のコンピュータシステムは、任意の適切な数のサブシステムを利用することができる。そのようなサブシステムの例は、図 13 でコンピュータ装置 1300 内に示されている。いくつかの実施態様において、コンピュータシステムは単一のコンピュータ装置を含み、ここでサブシステムはコンピュータ装置の構成要素とすることができる。他の実施態様においてコンピュータシステムは、内部構成要素を含む、それぞれがサブシステムである複数のコンピュータ装置を含むことができる。コンピュータシステムは、デスクトップコンピュータ及びラップトップコンピュータ、タブレット、携帯電話、及び他のモバイルデバイスを含むことができる。

【0132】

図 13 に示されるサブシステムは、システムバス 1305 を介して相互接続される。プリンタ 1340、キーボード 1370、記憶装置 1380、モニタ 1352 (これはディスプレイアダプタ 1350 に接続されている) などの追加のサブシステムが示されている。周辺機器及び入力/出力 (I/O) 装置 (これらは、I/O コントローラ 1310 に接続されている) は、当技術分野で任意の数の公知の手段、例えば、シリアルポート 1360 などによりコンピュータシステムに接続することができる。例えばシリアルポート 1360 又は外部インタフェース 1390 (例えば、イーサネット (登録商標)、Wi-Fi など) は、コンピュータシステム 1300 をインターネットなどの広域ネットワーク、マウス入力装置、又はスキャナに接続するために使用することができる。システムバス 1305 を介する相互接続は、中央プロセッサ 1330 が各サブシステムと通信し、システムメモリ 1320 又は記憶装置 1380 (例えば、固定ディスク) からの命令の実行、ならびにサブシステム間の情報の交換を制御することを可能にする。システムメモリ 1320 及び/又は記憶装置 1380 は、コンピュータ読み取り可能媒体を具体化することができる。本明細書に記載の任意の値は、一つの構成要素から別の構成要素に出力することができる、ユーザに出力することができる。

【0133】

コンピュータシステムは、例えば、外部インタフェース 1390 又は内部インタフェースによって接続される、複数の同一の構成要素又はサブシステムを含むことができる。いくつかの実施態様において、コンピュータシステム、サブシステム、又は装置は、ネットワーク上で通信することができる。このような例では一つのコンピュータはクライアントとして、別のコンピュータはサーバと見なすことができ、それぞれは、同じコンピュータシステムの一部であることができる。クライアントとサーバは、それぞれ複数のシステム、サブシステム、又は構成要素を含むことができる。

【0134】

なお、本発明の任意の実施態様は、ハードウェア (例えば、特定用途向け集積回路又はフィールドプログラマブルゲートアレイ) を使用するか、及び/又はコンピュータソフトウェアを使用して、一般にプログラム可能なプロセッサを用いてモジュラー又は統合的方法で、制御ロジックの形態で実施することができることを理解すべきである。本明細書において、プロセッサは、同一の集積チップ上でシングルコアプロセッサ、マルチコアプロセッサ、又は単一の回路基板上の又はネットワーク化された複数の処理ユニットを含む。本明細書で提供される開示及び教示に基づき、ハードウェア及びハードウェアとソフトウェアの組み合わせを用いて、本発明の実施態様を実施する他の手法及び/又は方法を、当業者は周知しており理解しているであろう。

【0135】

本出願に記載されている任意のソフトウェア構成要素又は機能は、例えば Java (登録商標)、C、C++、C#、Objective-C、Swift などのコンピュータ言語、又は Perl や Python などのスクリプト言語などの任意の適切なコンピュータ言語を用いて、例えば、従来型又はオブジェクト指向技術を用いて、プロセッサによって実行されるソフトウェアコードとして実行することができる。ソフトウェアコードは、記憶及び/又は送信のためのコンピュータ読み取り可能な媒体上の、一連の支持又は命令

として記憶することができる。適切な非一時的コンピュータ可読媒体としては、ランダムアクセスメモリ（RAM）、読み出し専用メモリ（ROM）、磁気媒体、例えばハードドライブ又はフロッピー（登録商標）ディスク、又は光学媒体、例えばコンパクトディスク（CD）若しくはDVD（デジタル多用途ディスク）、フラッシュメモリなどが挙げることができる。コンピュータ可読媒体は、このような記憶又は送信装置の任意の組み合わせであってもよい。

【0136】

このようなプログラムは、コード化され、インターネットを含む種々のプロトコルに適合する有線、光、及び/又は無線ネットワークを介して送信するために適合されたキャリア信号を用いて送信することができる。このように、本発明の実施態様に係るコンピュータ可読媒体は、そのようなプログラムでエンコードされたデータ信号を用いて作成することができる。プログラムコードでコード化されたコンピュータ可読媒体は、互換性のある装置と共にパッケージされるか、又は他の装置とは別に（例えば、インターネットダウンロードを介して）提供されてもよい。任意のそのようなコンピュータ可読媒体は、単一のコンピュータ製品（例えばハードドライブ、CD、又はコンピュータシステム全体）上に又はその中に常駐することができ、システム又はネットワーク内の異なるコンピュータ製品上又はその中に存在してもよい。コンピュータシステムは、ユーザに本明細書に記載の結果のいずれかを提供するための、モニタ、プリンタ、又は他の適切なディスプレイを含むことができる。

【0137】

本明細書に記載される任意の方法は、全体的に又は部分的に、工程を実行するように構成することができる1又は2以上のプロセッサを含むコンピュータシステムを用いて行うことができる。すなわち実施態様は、潜在的に各工程又は工程のそれぞれの群を実行する異なる構成要素を用いて、本明細書に記載された任意の方法の工程を実行するように構成されたコンピュータシステムに関する。番号付きの工程として提示されているが、本明細書の方法の工程は、同時に又は異なる順序で行うことができる。さらに、これらの工程の一部は、他の方法の他の工程の一部とともに使用することができる。また、工程の全て又は一部が任意であってよい。また、任意の方法の任意の工程は、これらの工程を実行するためのモジュール、回路、又は他の手段を用いて行うことができる。

【0138】

ある態様において本発明はまた、シーケンシングシステムも提供する。典型的なシーケンシングシステムは図14に表示される。図14に示されるシステムは、シーケンシング装置内に位置することができるシーケンシング分析モジュールと、コンピュータシステムの一部であるインテリジェンスモジュールとを含む。データセット（シーケンシングデータセット）は、ネットワーク接続又は直接接続を介して、分析モジュールからインテリジェンスモジュールに、又はその逆に、転送される。データセットは、例えば図4又は8に示すようにフローチャートに従って処理することができる。フローチャートに提供された工程は、コンピュータシステムのハードウェアに格納されたソフトウェアによって、例えば図15A及び15Bに記載されたフローチャートに従って便利に実行することができる。図15Aを参照して、コンピュータシステム（1100）は、例えば複数の配列リードから得られるデータを受け取るための受け取り手段（1110）、前記複数の配列リードを参照配列の標的領域にアラインメントするためのアラインメント手段（1120）、第一の試料の配列リードに基づいて、標的領域の第一の位置において、前記参照配列の第一の位置の参照アレルとは異なる第一のアレルを有する第一の候補バリエーションを同定するための同定手段（1130）、前記参照配列の第一の位置にアラインメントする第一の試料の配列リードに基づいて、第一の位置における第一のアレルに関して第一のバリエーション頻度を決定するための決定手段（1140）、複数のバリエーションクラスから選択される第一のバリエーションクラスに対応する第一の候補バリエーションを同定するための同定手段（1150）であって、ここで前記複数のバリエーションクラスの各バリエーションクラスは、異なるタイプのバリエーションに対応する上記手段、前記参照アレルを有する前記参照配列の

標的領域中の１セットの第二の位置を同定するための同定手段（１１６０）であって、前記１又は２以上の試料中の少なくとも５０％の他の位置は、第一のアレルに関して偽陽性を示し、そして前記１セットの第二の位置は前記第一の位置を含む上記手段を含み、前記１セットの第二の位置の各々において、かつ前記１又は２以上の試料の各々に関して、前記参照配列の第二の位置にアラインメントする試料の配列リードに基づいて、第一のアレルの第二のバリエーション頻度を決定するための決定手段（１１７０）であって、ここで前記第二のバリエーション頻度は統計分布を形成する上記手段、前記第一のバリエーション頻度を前記統計分布の統計値と比較して、前記統計分布の統計値に対する第一のバリエーション頻度の確率値を決定するための比較手段（１１８０）、そして、前記第一のアレルに関して、第一の試料において第一の候補バリエーションが真陽性であるか否かを決定する一部として、前記確率値を閾値と比較するための比較手段（１１９０）であって、ここで前記閾値は前記第一のアレルに関して偽陽性と真陽性とを区別する上記手段、を含むことができる。図１５Ｂを参照して、コンピュータシステム（２１００）は、例えば複数の配列リードから得られるデータを受け取るための受け取り手段（２１１０）、前記複数の配列リードを参照配列の標的領域にアラインメントするためのアラインメント手段（２１２０）、第一の位置における各試料のアラインメントされた配列リードに基づいて、少なくとも２つの試料の各試料中の第一の位置に、前記参照配列の第一の位置における参照アレルとは異なる第一の位置における第一のアレルが存在するか否かと同定するための同定手段（２１３０）、前記少なくとも２つの試料の各試料の、第一の位置における第一のアレルのバリエーション数と、第一の位置における参照アレルの野生型数を決定するための決定手段（２１４０）、前記少なくとも２つの試料から少なくとも１つの試料を参照試料として選択するための選択手段（２１５０）、前記第一の試料に関する第一の位置における第一のアレルの第一のバリエーション数と第一の位置における参照アレルの第一の野生型数とを、前記参照試料に関する第一の位置における第一のアレルの第二のバリエーション数と第一の位置における参照アレルの第二の野生型数とを比較して、第一の試料に関する第一の位置において第一のアレルを有するバリエーションの確率値を決定するための比較手段（２１６０）、そして前記第一の試料中の第一の位置における第一のアレルが、第一のアレルについて真陽性であるかどうかを決定する一部として、確率値を閾値と比較するための比較手段（２１７０）であって、ここで前記閾値は、前記第一の位置における第一のアレルに関して偽陽性と真陽性とを区別する上記手段、をさらに備えることができる。

【０１３９】

ある実施態様において、システムはまた、結果をコンピュータスクリーン上に表示するための表示手段を含むこともできる。図１４は、シーケンシング装置とコンピュータシステムとの間の相互作用を示す。システムは、シーケンシング装置内に位置することができる配列分析モジュールと、コンピュータシステムの一部であるインテリジェンスモジュールとを含む。データセット（シーケンシングデータセット）は、ネットワーク接続又は直接接続を介して、分析モジュールからインテリジェンスモジュールに又はその逆に、転送される。データセットは、プロセッサ上で作動しインテリジェンスモジュールの記憶装置に記憶されるコンピュータコードにより、図１５Ａ又は１５Ｂに従って処理することができ、処理後、分析モジュールの記憶装置に転送されて戻され、ここで修正されたデータは表示装置上に表示することができる。いくつかの実施態様において、インテリジェンスモジュールはまたシーケンシング装置で実行することができる。

【０１４０】

特定の実施態様の具体的な詳細は、本発明の実施態様の精神及び範囲から逸脱することなく、任意の適切な方法で組み合わせることができる。しかし、本発明の他の実施態様は、個々の態様、又はこれらの個々の態様の特定の組み合わせに関連する特定の実施態様に関してもよい。

【０１４１】

本発明の例示的な実施態様の上記記載は、例示と説明のために提示されている。これは網羅的であること又は記載した正確な形態に本発明を限定することを意図するものでもな

10

20

30

40

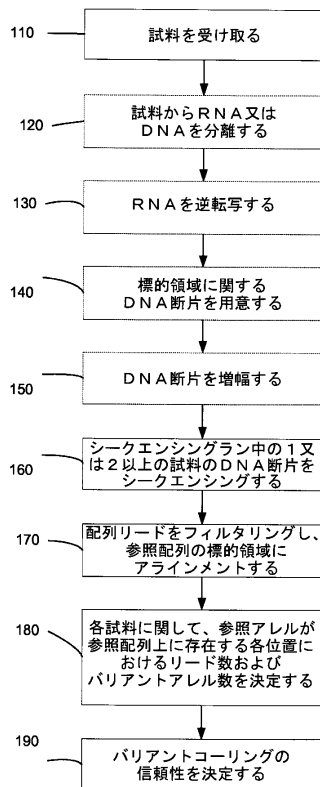
50

く、多くの修正及び変更が上記の教示に照らして可能である。

【0142】

「a」、「an」又は「the」の列挙は、特に別の指定がなければ、「1又は2以上」を意味することが意図される。特に別の指定がなければ、「又は」の使用は「含んで又は」を意味し、「含まなくて又は」ではないことを意味する。

【図1】



100

FIG. 1

【図2】

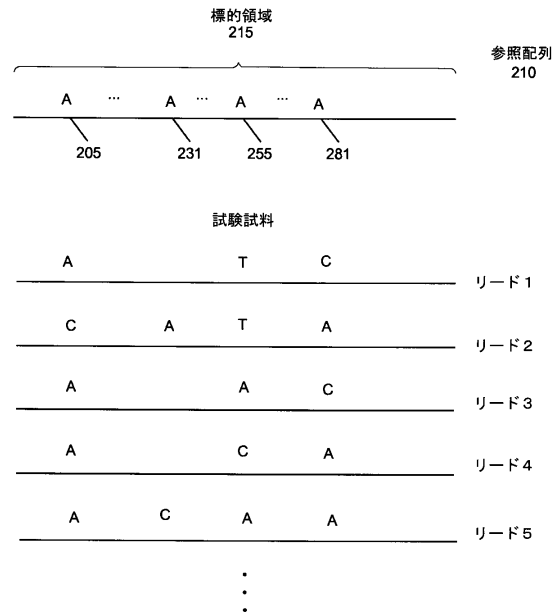


FIG. 2

【図 3 A】

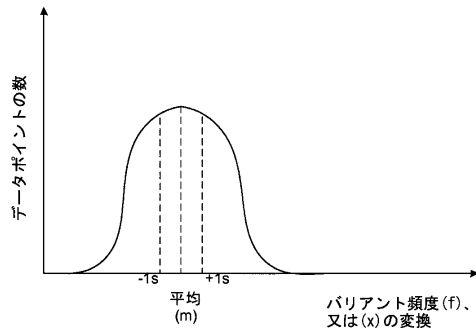


FIG. 3A

【図 3 C】

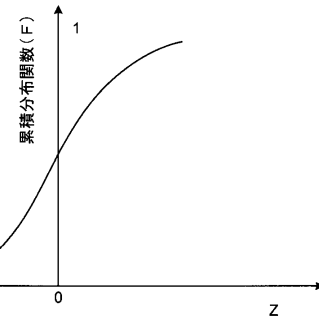


FIG. 3C

【図 3 B】

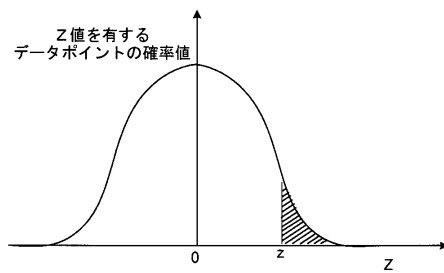


FIG. 3B

【図 3 D】

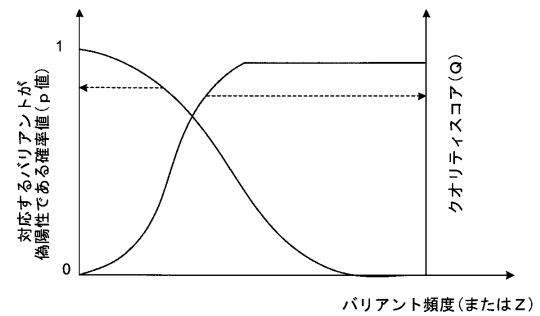
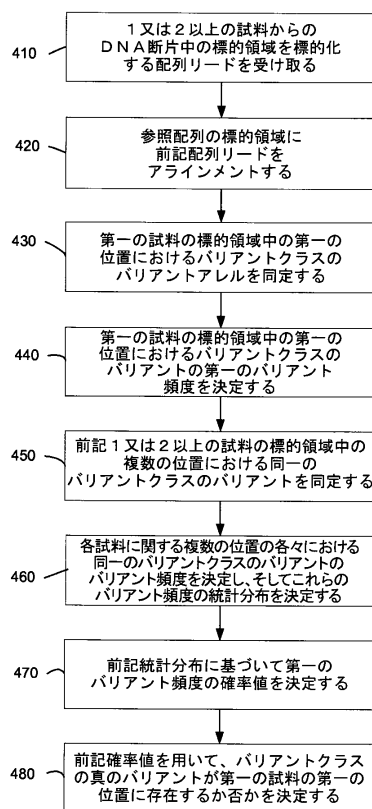


FIG. 3D

【図 4】



400

FIG. 4

【図 5】

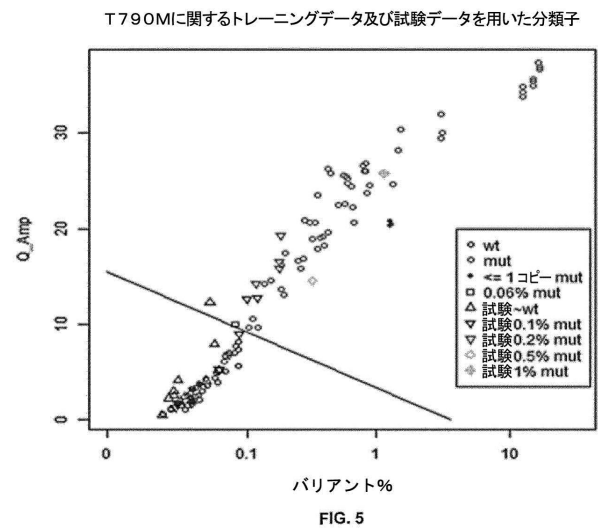
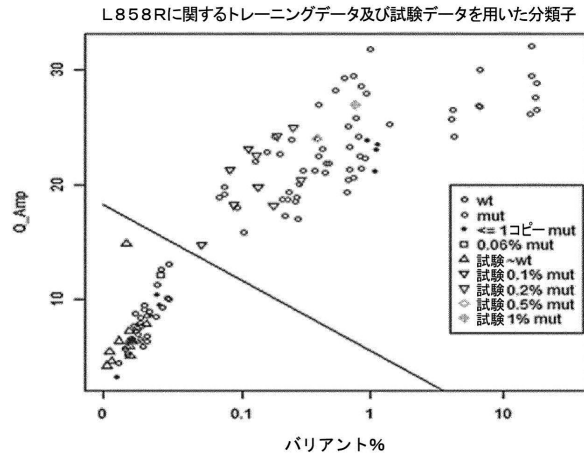
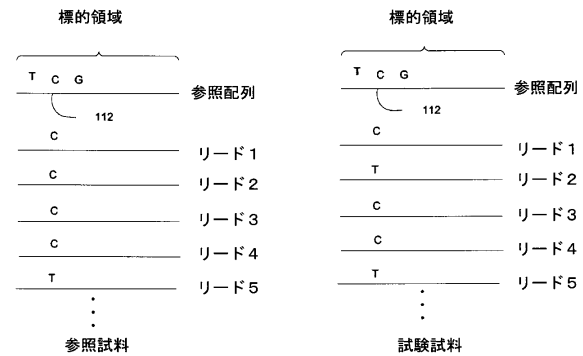


FIG. 5

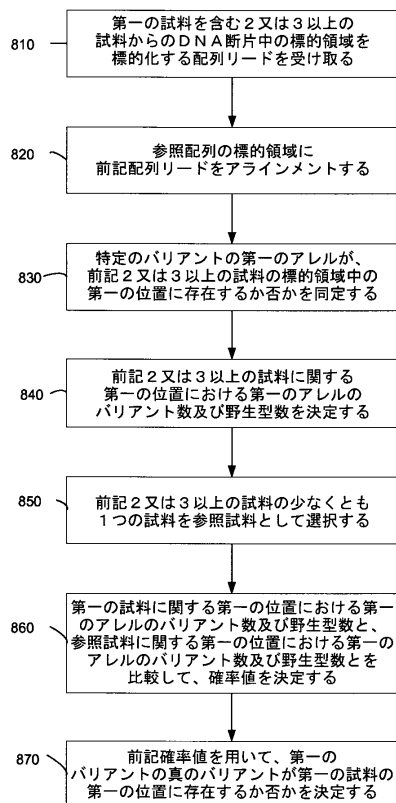
【図 6】



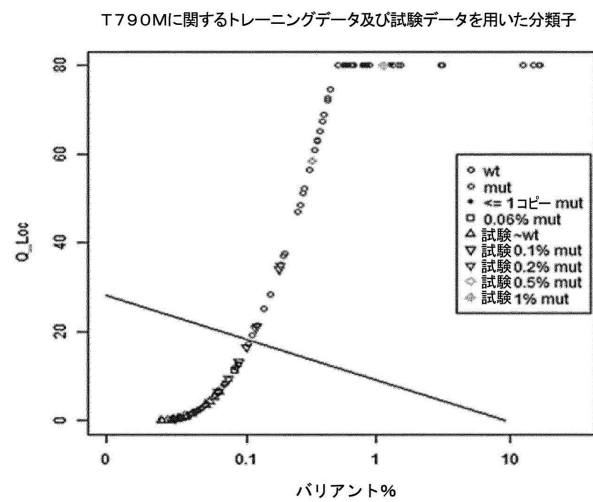
【図 7】



【図 8】



【図 9】



【図 10】

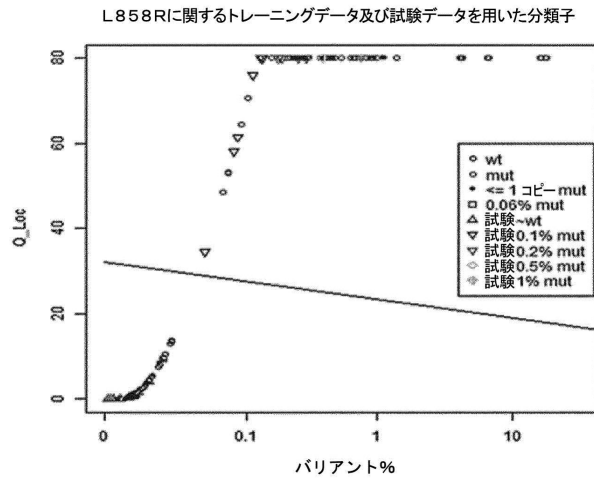


FIG. 10

【図 11】

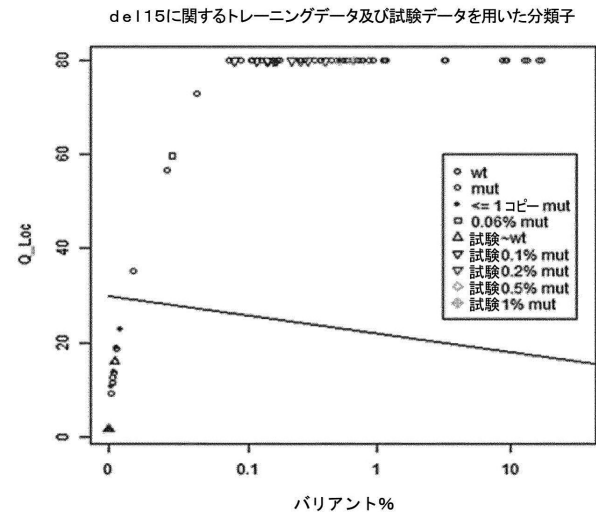


FIG. 11

【図 12】

FIG.12

Q	qnorm	qchisq	Q	qnorm	qchisq	Q	qnorm	qchisq	Q	qnorm	qchisq
3.5	0.1340	0.0180	35.5	3.4485	11.8923	67.5	5.0913	25.9211	99.5	6.3436	40.2417
4.5	0.3724	0.1387	36.5	3.5102	12.3216	68.5	5.1348	26.3657	100.5	6.3790	40.6916
5.5	0.5774	0.3334	37.5	3.5710	12.7519	69.5	5.1779	26.8106	101.5	6.4142	41.1416
6.5	0.7592	0.5764	38.5	3.6308	13.1830	70.5	5.2207	27.2557	102.5	6.4492	41.5917
7.5	0.9237	0.8532	39.5	3.6898	13.6148	71.5	5.2632	27.7010	103.5	6.4840	42.0419
8.5	1.0747	1.1550	40.5	3.7480	14.0475	72.5	5.3053	28.1465	104.5	6.5186	42.4923
9.5	1.2149	1.4760	41.5	3.8054	14.4809	73.5	5.3472	28.5923	105.5	6.5531	42.9427
10.5	1.3462	1.8122	42.5	3.8620	14.9149	74.5	5.3887	29.0383	106.5	6.5874	43.3932
11.5	1.4699	2.1606	43.5	3.9179	15.3497	75.5	5.4300	29.4844	107.5	6.6215	43.8438
12.5	1.5872	2.5192	44.5	3.9730	15.7850	76.5	5.4709	29.9308	108.5	6.6554	44.2945
13.5	1.6989	2.8863	45.5	4.0275	16.2209	77.5	5.5116	30.3774	109.5	6.6892	44.7453
14.5	1.8057	3.2606	46.5	4.0814	16.6574	78.5	5.5519	30.8241	110.5	6.7228	45.1962
15.5	1.9082	3.6412	47.5	4.1345	17.0945	79.5	5.5921	31.2710	111.5	6.7563	45.6472
16.5	2.0068	4.0271	48.5	4.1871	17.5320	80.5	5.6319	31.7181	112.5	6.7896	46.0983
17.5	2.1019	4.4178	49.5	4.2391	17.9701	81.5	5.6715	32.1654	113.5	6.8227	46.5494
18.5	2.1938	4.8127	50.5	4.2905	18.4086	82.5	5.7108	32.6129	114.5	6.8557	47.0007
19.5	2.2828	5.2113	51.5	4.3414	18.8476	83.5	5.7498	33.0605	115.5	6.8885	47.4520
20.5	2.3692	5.6133	52.5	4.3917	19.2870	84.5	5.7886	33.5082	116.5	6.9212	47.9034
21.5	2.4532	6.0182	53.5	4.4415	19.7269	85.5	5.8272	33.9562	117.5	6.9538	48.3549
22.5	2.5349	6.4259	54.5	4.4908	20.1671	86.5	5.8655	34.4043	118.5	6.9862	48.8065
23.5	2.6146	6.8360	55.5	4.5396	20.6077	87.5	5.9036	34.8525	119.5	7.0184	49.2582
24.5	2.6923	7.2484	56.5	4.5879	21.0488	88.5	5.9415	35.3009	120.5	7.0505	49.7099
25.5	2.7682	7.6629	57.5	4.6357	21.4901	89.5	5.9791	35.7494	121.5	7.0825	50.1617
26.5	2.8424	8.0793	58.5	4.6831	21.9318	90.5	6.0165	36.1980	122.5	7.1143	50.6136
27.5	2.9150	8.4974	59.5	4.7301	22.3739	91.5	6.0537	36.6468	123.5	7.1460	51.0654
28.5	2.9862	8.9171	60.5	4.7766	22.8163	92.5	6.0906	37.0958	124.5	7.1776	51.5174
29.5	3.0559	9.3384	61.5	4.8228	23.2590	93.5	6.1274	37.5448	125.5	7.2090	51.9696
30.5	3.1243	9.7610	62.5	4.8685	23.7020	94.5	6.1639	37.9940	126.5	7.2403	52.4217
31.5	3.1914	10.1850	63.5	4.9138	24.1452	95.5	6.2003	38.4433	127.5	7.2714	52.8742
32.5	3.2573	10.6102	64.5	4.9587	24.5888	96.5	6.2364	38.8927	128.5	7.3025	53.3259
33.5	3.3221	11.0365	65.5	5.0033	25.0327	97.5	6.2723	39.3423	129.5	7.3333	53.7788
34.5	3.3858	11.4639	66.5	5.0475	25.4768	98.5	6.3081	39.7919			

【図 13】

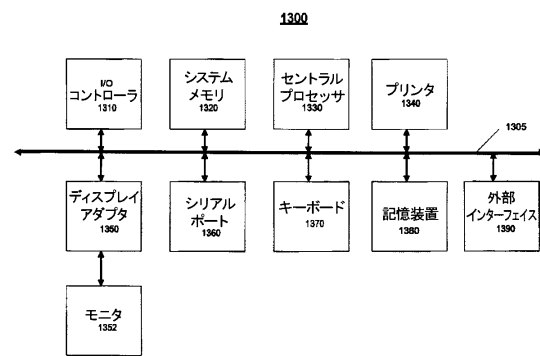


FIG. 13

【図 14】

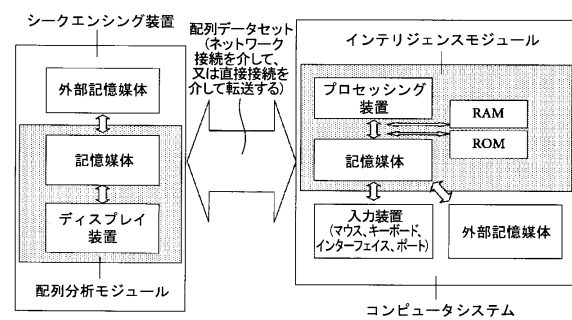
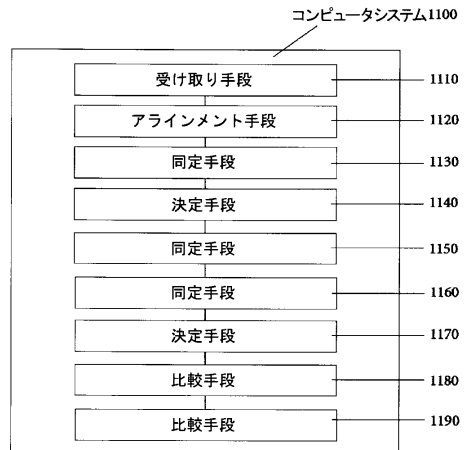


FIG. 14

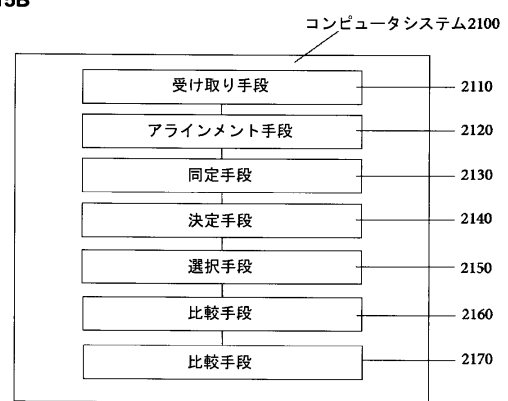
【図 15 A】

FIG. 15A



【図 15 B】

FIG. 15B



フロントページの続き

(74)代理人 100117019

弁理士 渡辺 陽一

(74)代理人 100150810

弁理士 武居 良太郎

(74)代理人 100164563

弁理士 佐々木 貴英

(72)発明者 ウエイ - ミン リウ

アメリカ合衆国, カリフォルニア 94568, ダブリン, シェルトン ストリート 4929

審査官 田付 徳雄

(56)参考文献 国際公開第2014/083023(WO, A1)

米国特許出願公開第2014/0143188(US, A1)

(58)調査した分野(Int.Cl., DB名)

G16B 5/00 - 99/00

C12Q 1/6869