

(19) 日本国特許庁 (JP)

(12) 公表特許公報 (A)

(11) 特許出願公表番号

特表2018-501539

(P2018-501539A)

(43) 公表日 平成30年1月18日 (2018.1.18)

(51) Int.Cl.		F I		テーマコード (参考)
<b>G 0 6 F 19/22</b>	<b>(2011.01)</b>	G 0 6 F 19/22		
<b>C 1 2 N 15/09</b>	<b>(2006.01)</b>	C 1 2 N 15/00	A	

審査請求 未請求 予備審査請求 未請求 (全 23 頁)

(21) 出願番号 特願2017-521153 (P2017-521153) (86) (22) 出願日 平成27年10月15日 (2015.10.15) (85) 翻訳文提出日 平成29年5月26日 (2017.5.26) (86) 国際出願番号 PCT/US2015/055807 (87) 国際公開番号 W02016/061396 (87) 国際公開日 平成28年4月21日 (2016.4.21) (31) 優先権主張番号 62/064,717 (32) 優先日 平成26年10月16日 (2014.10.16) (33) 優先権主張国 米国 (US)	(71) 出願人 517131237 カウンシル, インコーポレイテッド アメリカ合衆国 カリフォルニア 940 80, サウス サンフランシスコ, キ ンボール ウェイ 180 (74) 代理人 100078282 弁理士 山本 秀策 (74) 代理人 100113413 弁理士 森下 夏樹 (74) 代理人 100181674 弁理士 飯田 貴敏 (74) 代理人 100181641 弁理士 石川 大輔 (74) 代理人 230113332 弁護士 山本 健策
--	--

最終頁に続く

(54) 【発明の名称】 バリエントコーラー

## (57) 【要約】

基準ゲノム配列に対してゲノムサンプルからバリエントを読み取るためのプロセスおよびシステムが、提供される。例示的プロセスは、リードのセットを収集することと、リードから  $k$ -mer グラフを生成することを含む。例えば、 $k$ -mer グラフは、収集されたリードの全ての可能なサブストリングを表すように構築されることができる。 $k$ -mer グラフは、連続的グラフにまとめられ、可能なハプロタイプのセットが、連続的グラフから生成されてもよい。プロセスはさらに、共通シーケンサエラーのためのフィルタを提供するエラーテーブルを生成してもよい。プロセスは、次いで、ハプロタイプのセットおよび生成されたエラーテーブルに基づいて、ディプロタイプのセットを生成し、ディプロタイプのセットをスコアリングし、基準ゲノムからバリエントを識別してもよい。

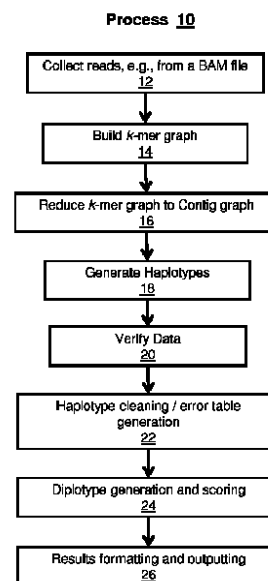


FIG. 1

## 【特許請求の範囲】

## 【請求項 1】

基準ゲノム配列に対してゲノムサンプルからバリエントを判定するためのコンピュータ実装方法であって、

少なくとも 1 つのプロセッサおよびメモリを有する電子デバイスにおいて、

先に配列決定されたサンプルからの配列データのエラーテーブルにアクセスするステップと、

ゲノムサンプルから収集されたリードのセットから可能なハプロタイプのセットを判定するステップと、

前記可能なハプロタイプのセットおよび前記エラーテーブルに基づいて、ディプロタイプのセットを生成するステップであって、前記可能なハプロタイプのセットは、前記エラーテーブルによってフィルタリングされる、ステップと、

前記ディプロタイプのセットをスコアリングするステップと、

前記ディプロタイプのセットをスコアリングするステップに基づいて、バリエントを出力するステップと、を含む、方法。

10

## 【請求項 2】

収集されたリードのセットから  $k - m e r$  グラフを生成するステップと、

前記生成された  $k - m e r$  グラフを連続的グラフにまとめるステップと、

前記連続的グラフから前記可能なハプロタイプのセットを生成するステップと、

をさらに含む、請求項 1 に記載の方法。

20

## 【請求項 3】

前記ディプロタイプのセットをスコアリングするステップはさらに、ディプロタイプ毎に事後確率を判定することを含む、請求項 1 に記載の方法。

## 【請求項 4】

前記エラーテーブルを生成するステップをさらに含み、前記エラーテーブルを生成するステップは、

リードを基準サンプルに対してアライメントすることと、

リードが前記基準サンプルとミスマッチを有する部位を判定することと、

ミスマッチを有する部位を前記エラーテーブルに追加することと、を含む、請求項 1 に記載の方法。

30

## 【請求項 5】

前記エラーテーブルを生成するステップはさらに、シーケンサエラーと関連付けられない部位を前記エラーテーブルからフィルタリングすることを含む、請求項 4 に記載の方法。

## 【請求項 6】

前記エラーテーブルを生成するステップはさらに、Hardy - Weinberg 試験、Bayes Factor 試験、またはStrand Bias 試験の 1 つまたはそれより多くを使用して、閾値に満たない部位を前記エラーテーブルからフィルタリングすることを含む、請求項 4 に記載の方法。

## 【請求項 7】

配列データのエラーテーブルを生成するためのコンピュータ実装方法であって、

少なくとも 1 つのプロセッサおよびメモリを有する電子デバイスにおいて、

ゲノムサンプルから収集されたリードのセットから可能なハプロタイプのセットを判定するステップと、

前記収集されたリードのセットを基準サンプルに対してアライメントするステップと、

前記基準サンプルから前記収集されたリードのセットのリードがミスマッチを有する部位を判定するステップと、

ミスマッチを有する部位をエラーテーブルに追加するステップと、を含む、方法。

40

## 【請求項 8】

前記可能なハプロタイプのセットを判定するステップは、

$k - m e r$  グラフを前記収集されたリードのセットから生成することと、

50

前記生成された  $k$ -mer グラフを連続的グラフにまとめることと、  
前記連続的グラフから前記可能なハプロタイプのセットを判定することと、を含む、請求項 7 に記載の方法。

【請求項 9】

非一過性コンピュータ可読記憶媒体であって、  
先に配列決定されたサンプルからの配列データのエラーテーブルにアクセスするステップと、  
ゲノムサンプルからの収集されたリードのセットから可能なハプロタイプのセットを判定するステップと、  
前記可能なハプロタイプのセットおよび前記エラーテーブルに基づいて、ディプロタイプの  
10     のセットを生成するステップであって、前記可能なハプロタイプのセットは、前記エラー  
テーブルによってフィルタリングされる、ステップと、  
前記ディプロタイプのセットをスコアリングするステップと、  
前記ディプロタイプのセットをスコアリングするステップに基づいて、バリエانتを出力  
するステップとのためのコンピュータ実行可能命令を含む、非一過性コンピュータ可読記憶媒体。

【請求項 10】

収集されたリードのセットから  $k$ -mer グラフを生成するステップと、  
前記生成された  $k$ -mer グラフを連続的グラフにまとめるステップと、  
前記連続的グラフから前記可能なハプロタイプのセットを生成するステップとをさらに含  
20     む、請求項 9 に記載の非一過性コンピュータ可読記憶媒体。

【請求項 11】

前記ディプロタイプのセットをスコアリングするステップはさらに、ディプロタイプ毎に  
事後確率を判定することを含む、請求項 9 に記載の非一過性コンピュータ可読記憶媒体。

【請求項 12】

前記エラーテーブルを生成するステップをさらに含み、前記エラーテーブルを生成するス  
テップは、  
リードを基準サンプルに対してアライメントすることと、  
リードが前記基準サンプルとミスマッチを有する部位を判定することと、  
ミスマッチを有する部位を前記エラーテーブルに追加することとを含む、請求項 9 に記載  
30     の非一過性コンピュータ可読記憶媒体。

【請求項 13】

前記エラーテーブルを生成するステップはさらに、シーケンサエラーと関連付けられない  
部位を前記エラーテーブルからフィルタリングすることを含む、請求項 12 に記載の非一  
過性コンピュータ可読記憶媒体。

【請求項 14】

前記エラーテーブルを生成するステップはさらに、Hardy - Weinberg 試験、  
Bayes Factor 試験、または Strand Bias 試験の 1 つまたはそれよ  
り多くを使用して、閾値に満たない部位を前記エラーテーブルからフィルタリングするこ  
とを含む、請求項 12 に記載の非一過性コンピュータ可読記憶媒体。  
40

【請求項 15】

システムであって、  
1 つまたはそれより多くのプロセッサと、  
メモリと、  
1 つまたはそれより多くのプログラムであって、前記 1 つまたはそれより多くのプログラ  
ムは、前記メモリ内に記憶され、前記 1 つまたはそれより多くのプロセッサによって実行  
されるように構成され、  
先に配列決定されたサンプルからの配列データのエラーテーブルにアクセスするステップ  
と、  
ゲノムサンプルから収集されたリードのセットから可能なハプロタイプのセットを判定す  
50

るステップと、  
前記可能なハプロタイプのセットおよび前記エラーテーブルに基づいて、ディプロタイプの  
のセットを生成するステップであって、前記可能なハプロタイプのセットは、前記エラー  
テーブルによってフィルタリングされる、ステップと、  
前記ディプロタイプのセットをスコアリングするステップと、  
前記ディプロタイプのセットをスコアリングするステップに基づいて、バリエーションを出力  
するステップとのための命令を含む、1つまたはそれより多くのプログラムと、を備える  
、システム。

【請求項16】

収集されたリードのセットから  $k$ -mer グラフを生成するステップと、  
前記生成された  $k$ -mer グラフを連続的グラフにまとめるステップと、  
前記連続的グラフから前記可能なハプロタイプのセットを生成するステップと、  
をさらに含む、請求項9に記載のシステム。

10

【請求項17】

前記ディプロタイプのセットをスコアリングするステップはさらに、ディプロタイプ毎に  
事後確率を判定することを含む、請求項9に記載のシステム。

【請求項18】

前記エラーテーブルを生成するステップをさらに含み、前記エラーテーブルを生成するス  
テップは、  
リードを基準サンプルに対してアライメントすることと、  
リードが前記基準サンプルとミスマッチを有する部位を判定することと、  
ミスマッチを有する部位を前記エラーテーブルに追加することを含む、請求項9に記載  
のシステム。

20

【請求項19】

前記エラーテーブルを生成するステップはさらに、シーケンサエラーと関連付けられない  
部位を前記エラーテーブルからフィルタリングすることを含む、請求項18に記載のシス  
テム。

【請求項20】

前記エラーテーブルを生成するステップはさらに、  
Hardy-Weinberg 試験、Bayes Factor 試験、または Strand Bias 試験の1つまたはそれより多くを使用して、閾値に満たない部位を前記エラ  
ーテーブルからフィルタリングすることを含む、請求項18に記載のシステム。

30

【発明の詳細な説明】

【技術分野】

【0001】

(関連出願への相互参照)

本出願は、2014年10月16日に出願された、「VARIANT CALLER」  
と題する米国仮出願番号第62/064,717号に基づく優先権を主張しており、その  
内容は、すべての目的のためにその全体が参考として本明細書によって援用される。

40

【0002】

(分野)

本願は、概して、DNAシーケンサリード(read)におけるバリエーションを識別およ  
び定量化するためのプロセスおよびシステムに関し、一実施例では、エラーテーブルの使  
用を通して基準ゲノム配列からバリエーションを識別し、ハプロタイプエラーを除去し、次い  
で、ディプロタイプ(対のハプロタイプ)を生成およびスコアリングし、バリエーションを判  
定するためのバリエーションコーラープロセスおよびシステムに関する。

【背景技術】

【0003】

50

( 背景 )

バリエントコーラー ( variant caller ) は、概して、基準ゲノム配列に対して DNA 配列リード内のヌクレオチド差異が存在することを判定する。Platypus、Genome Analysis Toolkit「GATK」、およびFreeBayesとして知られるものを含む、いくつかの公知のバリエントコーラーが存在する。例えば、Platypusは、主に、リードの局所再アライメントおよびその局所アセンブリに依拠する、高スループット配列決定データ内のバリエント検出のためのシステムである。Platypusは、「Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications」においてより詳細に説明されており、参照することによってその全体が本明細書に援用される。

10

【発明の概要】

【課題を解決するための手段】

【0004】

( 要旨 )

一実施例では、基準ゲノム配列に対してゲノムサンプルからバリエントを読み取るためのコンピュータ実装プロセスが、提供される。このプロセスは、リードのセットを収集することと、リードからk-merグラフを生成することを含む。例えば、k-merグラフは、収集されたリードの全ての可能なサブストリングを表すように構築されることができる。k-merグラフは、連続的グラフにまとめられ、可能なハプロタイプのセットが、連続的グラフから生成されてもよい。プロセスはさらに、エラーテーブルを生成してもよく（例えば、共通シーケンサエラーを識別するための多くの以前のサンプルから）、これは、共通シーケンサエラーのためのフィルタを提供する。プロセスは、次いで、ハプロタイプのセットおよびエラーテーブルに基づいて、ディプロタイプのセットを生成し、ディプロタイプのセットをスコアリングし、基準ゲノムからバリエントを識別してもよい。ディプロタイプをスコアリングすることは、ディプロタイプ毎に事後確率を判定することを含んでもよく、最高スコアリングディプロタイプが結果として報告される。

20

【0005】

別の実施例では、配列データのエラーテーブルを生成するためのコンピュータ実装プロセスが、提供される。例示的プロセスは、少なくとも1つのプロセッサおよびメモリを有する電子デバイスにおいて、ゲノムサンプルから収集されたリードのセットから可能なハプロタイプのセットを判定するステップと、収集されたリードのセットを基準サンプルに対してアライメントするステップと、基準サンプルから、収集されたリードのセットのリードがミスマッチを有する部位を判定するステップと、ミスマッチを有する部位をエラーテーブルに追加するステップとを含んでもよい。可能なハプロタイプのセットを判定するステップは、k-merグラフを、収集されたリードのセットから生成することと、生成されたk-merグラフを連続的グラフにまとめることと、連続的グラフから可能なハプロタイプのセットを判定することとを含んでもよい。

30

【0006】

加えて、バリエントコーラーのためと、エラーテーブルを生成するためのシステム、電子デバイス、グラフィカルユーザインターフェース、および非一過性コンピュータ可読記憶媒体（説明される1つまたはそれより多くのプロセスを実行するためのプログラムおよび命令を含む、記憶媒体）が、説明される。

40

【図面の簡単な説明】

【0007】

本願は、同一部分が同一数字によって参照され得る、付随の図面と関連して検討される以下の説明を参照することによって、最良に理解され得る。

【0008】

【図1】図1は、一実施形態による、例示的コーリングプロセスを図示する。

50

【 0 0 0 9 】

【図 2 A】図 2 A - 2 C は、図 1 のプロセスを参照して説明される例示的プロセスを図式的に図示する。

【図 2 B】図 2 A - 2 C は、図 1 のプロセスを参照して説明される例示的プロセスを図式的に図示する。

【図 2 C】図 2 A - 2 C は、図 1 のプロセスを参照して説明される例示的プロセスを図式的に図示する。

【 0 0 1 0 】

【図 3】図 3 A および 3 B は、異なるリードモデルのプロットを図示する。

【 0 0 1 1 】

【図 4】図 4 は、本発明の種々の実施形態が動作し得る、例示的システムおよび環境を図示する。

【 0 0 1 2 】

【図 5】図 5 は、例示的コンピューティングシステムを図示する。

【発明を実施するための形態】

【 0 0 1 3 】

以下の説明は、当業者が、種々の実施形態を作製および使用することが可能となるように提示される。具体的デバイス、技法、および用途の説明は、実施例のみとして提供される。本明細書に説明される実施例の種々の修正は、当業者に容易に明白となり、本明細書に定義される一般的原理は、本技術の精神および範囲から逸脱することなく、他の実施例および用途に適用されてもよい。したがって、開示される技術は、本明細書に説明および示される実施例に限定されず、請求項と一貫した範囲が与えられることが意図される。

【 0 0 1 4 】

本願は、概して、基準ゲノム配列からバリエントを識別するためのバリエントコーラーに関する。一実施例では、バリエントコーラーは、エラーテーブルを生成してエラーをハプロタイプから除去し、ディプロタイプを生成し、ディプロタイプをスコアリングし、基準ゲノム配列からバリエントを識別するためのプロセスを含む。バリエントコーラーの実施例は、P l a t y p u s、G A T K、F r e e b a y e s、およびその他等の公知のコーラーに優るいくつかの進歩を提供し得る。例えば、全ての実施形態または実施例に存在しないが、進歩は、リード内のアライメントの代わりに（例えば、アライメントのためにリードを蓄積し、全リードを使用して、1つのグラフを作成する代わりに）、局所化と、共通シーケンサエラーを防ぐために、エラーテーブルを介して、エラー較正とを含んでもよい。

【 0 0 1 5 】

一実施形態では、バリエントコーラーは、いくつかの処理段階に分割され、各段階は、次の段階への入力としてその出力を提供する。以下の実施例は、配列データを記憶するためのバイナリフォーマットである、B i n a r y A l i g n m e n t / M a p フォーマット「b a m」または「B A M」フォーマットの使用を仮定する。しかしながら、他のデータフォーマット（例えば、S e q u e n c e A l i g n m e n t / M A P フォーマットまたは「S A M」フォーマット）も、企図され、そして可能である。一実施例では、各 b a m ファイル内の各領域の処理は、全ての他の領域および b a m ファイルとは完全に別個である。

【 0 0 1 6 】

広義には、かつ一実施例では、ある領域のためのコールを生成するために、図 1 にプロセス 10 として図示される、以下のプロセスが、行われる。プロセス 10 の説明と併せて、プロセス 10 の種々の側面を図式的に図示する図 2 A - 2 C を基準されたい。

【 0 0 1 7 】

最初に、着目配列が、12において得られる。例えば、リードが、コールの領域と何らかの点で重複する b a m ファイルから収集されることができる。この処理は、B W A、B O W T I E、M A X 等のショートリードアライナを使用して、図 2 A に図式的に図示され

10

20

30

40

50

るように、リード 210 をゲノム領域 220 に対してアライメントすることを含んでもよい。収集されたリードは、次いで、その関連付けられたソフトクリッピング情報を使用して、クリッピングされることができる。アライナからの補助情報、例えば、塩基間 (base-to-base) アライメント情報は、次いで、破棄されることができ、リードは、単に、塩基の配列となる。(いくつかの実施例では、マッピング品質に基づくフィルタリングが、随意に、行われることができる。)

#### 【0018】

k-mer グラフが、次いで、14において、収集されたリードから構築され、k-mer グラフは、収集されたリードともに含まれる、長さkの全ての可能なサブストリングを表す。例示的 k-mer グラフは、図2Bに図示され、そこでは、k=3である(実際は、20~30のkが、k-mer が一意であって、例えば、1カ所にのみ生じることを確実にするために使用されてもよい)。例えば、各リードは、k-mer および k-mer 遷移を収集するために走査される。各エッジは、その関連付けられた遷移の確率でアノテーションされ、各 k-mer は、エッジの起点として認められた回数でアノテーションされる。k-mer AとBとの間の遷移の確率は、k-mer Aが認められた合計回数で除算された k-mer Aに続く k-mer Bが認められた回数である。

#### 【0019】

k-mer グラフは、次いで、16において、処理の単純化のために、連続的(「コンティグ」)グラフにまとめられることができる。コンティググラフは、概して、ゲノム情報の領域をとともに形成する、重複セグメントのセットを図示する。例えば、本ステップは、それらが常時、同一経路内で終端する場合、2つの k-mer を結合することができる。加えて、k-mer グラフは、閾値回数未満(例えば、4回未満)認められる任意の k-mer を破棄し、閾値を下回る(例えば、3%を下回る)確率を有する任意のエッジを破棄することによってフィルタリングされる。いったん k-mer グラフが作成されると、それは、サイクル、すなわち、それ自体に収束する経路に関してチェックされることができる。グラフがサイクルを有する場合、破棄され、kが増加され、グラフが再構築されることができる。したがって、本実施例では、k-mer グラフは、サイクルを伴わずに構築される。

#### 【0020】

ハプロタイプ生成が、次いで、18において、行われることができる。例えば、いったんコンティググラフが構築されると、ハプロタイプ候補のための開始点が、入エッジ(入次数0)を伴わない全コンティグを検索することによって見出されることができる。これらは、領域の開始時のコンティグであるはずであるが、領域の中央におけるコンティグもまた、雑音に起因して作成される場合、本特性を有し得る。次いで、それらのコンティグを開始点として見なし、コンティググラフを通して全ての可能な経路が、列挙され、各経路は、出エッジを伴わないコンティグ(終末)に到達すると終端する。次に進む前に、全経路は、それらのコンティグを結合することによって、ハプロタイプストリングに変換されることができる。簡略化された実施例は、図2Cに図示され、開始点は、「1」によって示され、「6」まで続く。各可能な経路は、可能なハプロタイプを生成し、そのうちの1つが、図に示される。

#### 【0021】

いったん可能なハプロタイプのセットが生成されると、例示的プロセスは、20において、十分に良好なコールを作成するために十分なデータを有することを検証する(1つまたはそれより多くの仮説形成法(huristics)を通して)。例えば、プロセスは、所望の領域内の各位置が、十分な k-mer によって網羅され、領域全体を網羅する、少なくとも1つのハプロタイプが存在することをチェックする。これらのチェックのいずれかが失敗する場合、領域全体に対するコールは、発行され得ない。仮説形成法は、コール内の所望の信頼性に関して調節されることができることを理解されたい。

#### 【0022】

可能なハプロタイプのセットはさらに、22において、任意のスコアリングプロセスの

10

20

30

40

50

前に、「精緻化」されることができる。コンティググラフから生成されたハプロタイプは、概して、出力またはスコアリングのために好適ではない。故に、一実施例では、スコアリングの前に、それらはいくつかの補正相を受ける。最初に、ハプロタイプは、コーラーが全重複リードを使用し、大部分のハプロタイプが、元々、着目領域のエッジを越えて延在し得るため、この領域にクリッピングされる。一実施例では、ハプロタイプをクリッピングするために、問題の領域に対してアライメントされ、アライメント外の任意の塩基は、破棄される。いったんハプロタイプがクリッピングされると、ハプロタイプ内のエラーが、補正されることができる。例えば、プロセスは、多くのサンプルから、共通シーケンサエラーをリスト化する、エラーテーブル（以下により詳細に説明される）を生成することができ、本エラーテーブルは、可能なハプロタイプのセットからそれらのエラーを除去するために使用されることができる。これらのステップは、複製を含むハプロタイプのセットをもたらし得、これらの複製は、ドロップされることができる。

10

20

30

40

#### 【0023】

ディプロタイプは、24において、ハプロタイプから生成され、そしてスコアリングされることができる。例えば、N個のハプロタイプのセットが、全ての可能なディプロタイプを生成するために、それ自体と組み合わせられることができる。N個のハプロタイプに関して、 $N(N+1)/2$ 個の一意のディプロタイプが存在するであろう。これらのディプロタイプは、次いで、スコアリングされることができ、ディプロタイプのスコアは、その事後確率  $P(\text{diplo type} | \text{reads})$  に等しい。最高スコアリングディプロタイプは、結果として報告されることができ、その信頼性は、最高確率と次に高い確率との間の比率の対数と等しい。ディプロタイプのスコアリングは、以下により詳細に説明される。

#### 【0024】

結果は、次いで、26において、フォーマットされ（必要に応じて）、要求に応じて書き出されることができる。例えば、フォーマットがJavaScript Object Notation（「json」または「JSON」）またはVariant Call Format（「vcf-full」である場合、さらなる処理は、本実施例では、必要とされず、コールは、単に、ディスクに書き出される。しかしながら、結果フォーマットが、Variant Call Format - Single Nucleotide Polymorphism（「vcf-snp」）である場合、結果は、より小さいコールに分割され、これは、領域をその個々のSNPおよび挿入欠失に分割する。vcf-snpフォーマットにおける単一コールは、異なるバリエーションが相互にある距離（例えば、10塩基）内にある、全バリエーションから成る。

#### 【0025】

##### ディプロタイプのスコアリング

#### 【0026】

一実施例では、前述のN個のハプロタイプのセットは、全ての可能なディプロタイプを生成するために、それ自体と組み合わせられることができる。N個のハプロタイプに関して、 $N(N+1)/2$ 個の一意のディプロタイプが存在し得る。これらのディプロタイプは、次いで、スコアリングされる；ディプロタイプのスコアは、その事後確率  $P(\text{diplo type} | \text{reads})$  に等しい。最高スコアリングディプロタイプは、結果として報告されることができ、その信頼性は、最高確率と次に高い確率との間の比率の対数と等しい。

#### 【0027】

最良ディプロタイプを候補のリストから判定するために使用される例示的確率スコアリングモデルが、ここで説明される。一実施例では、各ディプロタイプに割り当てられるスコアは、ディプロタイプの事後確率  $P(\text{diplo type} | \text{reads})$  である。スコアリングのために使用される確率は典型的には、小さいため、一実装では、対数確率が使用される。事後確率は、以下のように、尤度および事前確率に分解されることができる。



## 【数 1】

$$P(\text{diplotype} \mid \text{reads}) = (1/Z) P(\text{reads} \mid \text{diplotype}) P(\text{diplotype})$$

式中、 $Z = P(\text{reads})$  は、ある正規化定数であって、算出されない。 $Z$  は、ディプロタイプから独立するため、2つのディプロタイプを比較する目的のために無視されることができる。事前確率  $P(\text{diplotype})$  および尤度  $P(\text{reads} \mid \text{diplotype})$  が、次いで、別個に算出されることができる。

## 【0028】

事前確率を算出するために、本実施例では、大部分の領域は、基準に類似すると仮定され得る。ディプロタイプの確率は、したがって、ディプロタイプが基準から生物学的変異を介して生成された確率である。本実施例は、これが、単に、基準から生成されているハプロタイプの確率の積と仮定する（選択に起因して、完全に正確ではないが、概して、十分であることを理解されたい）。したがって、ディプロタイプの確率は、以下のように表されることができる。

## 【数 2】

$$P(\text{diplotype}) = P(\text{haplotype}_1) P(\text{haplotype}_2)$$

## 【0029】

生成されているハプロタイプの確率は、それが全ての可能な方法において生成されている確率の総和であって、ハプロタイプと基準の各可能なアライメントは、ハプロタイプを生成する異なる方法に対応する。しかしながら、全アライメントにわたって総和を行うことは、算出上、扱いにくくあり得、したがって、本実施例は、確率質量の大部分が単一アライメント内に含有され、最高確率を有するものと仮定する。したがって、 $P(\text{haplotype})$  を算出するために、プロセスは、ハプロタイプを基準に対してアライメントする。アライメントの間に使用されるマッチ、ミスマッチ、ギャップ開放、およびギャップ伸長パラメータは、生物学的変異に起因して生じるような事象の対数確率に対応する。アライメントは、スコアを最大限にするため、対数確率を最大限にし、したがって、最高確率アライメントをもたらす。例えば、1塩基変化は、約1,000塩基毎に生じ、したがって、ミスマッチパラメータは、 $\log(1/1000)$  となり得る。

## 【0030】

尤度  $P(\text{reads} \mid \text{diplotype})$  の算出は、類似プロセスを使用する。最初に、実施例は、全リードが独立であると仮定し、これは、尤度が以下のように書き直されることを可能にする。

## 【数 3】

$$P(\text{reads} \mid \text{diplotype}) = \text{product}_i \{ P(\text{read}_i \mid \text{diplotype}) \}$$

## 【0031】

次いで、実施例は、リードが、2つのハプロタイプのディプロタイプから生じ得る（等確率を伴う）か、またはゲノム内の別の場所から無作為に生成され得る（非常に低確率を伴う）かのいずれかであると仮定する。後者の場合は、アライナエラーおよび稀な外れ値を効果的にモデル化する。したがって、リードの確率は、以下のように表されることができる。

## 【数 4】

$$P(\text{read} \mid \text{diplotype}) = \text{epsilon} P(\text{read is random}) + (0.5 - \text{epsilon}) P(\text{read} \mid \text{haplotype}_1) + (0.5 - \text{epsilon}) P(\text{read} \mid \text{haplotype}_2)$$

## 【0032】

無作為に生成されたリードの確率は、4つの等しい可能性の塩基が存在するため、生成された各塩基に等しい。

10

20

30

40

50

## 【数 5】

$$P(\text{read is random}) \sim 0.25^{\text{len}(\text{read})}$$

## 【0033】

ハプロタイプが与えられたリードの確率は、アライメントを使用して見出されることができる。本実施例は、ハプロタイプが、基となるゲノムの真の配列であって、リードが、エラーを含んだ配列決定プロセスを使用して、この配列から生成されると仮定する。したがって、アライメントパラメータは、シーケンサエラーの率であるはずである。例えば、ミスマッチパラメータは、任意の塩基においてシーケンサが1つの塩基変化を作成する確率の対数であるはずである。事前確率と同様に、プロセスは、最良アライメントを算出し、スコアを確率として使用する。

10

## 【0034】

他のスコアリングプロセスが、ここで説明されるものの代わりに、またはそれに加えて、使用されてもよく、例えば、他のパラメータ、値、仮定、および算出プロセスを含むことが、当業者によって理解されるはずである。

## 【0035】

エラーテーブル生成

## 【0036】

一般に、そして一実施例では、エラーテーブルは、共通シーケンサエラーを防ぐためのフィルタのように作用し、これは、いくつかの領域を別様にコーリングすることを非常に困難にし得る。一実施例では、エラーテーブルを生成するために、同一領域に対してデータを含む、数百（例えば、100～300またはそれを上回る）サンプルが、使用される。本実施例では、所与の領域に対するエラーテーブル生成は、以下のステップを受ける。

20

1. サンプル毎に、リードを基準に対してアライメントする。基準内の塩基毎に、異なるバリエーションがそこで認められる回数をカウントする（バリエーションは、4つの塩基、異なる長さ欠失、および異なる挿入である）。本プロセスは、フォワードリードおよびバックワードリードに関して別個に行われることができる。

2. ある閾値を上回るバリエーションが存在する、すなわち、ある閾値%のリードが非基準対立遺伝子を有する、部位を見出す。例えば、閾値は、1%であることができる。これらの部位は、エラーテーブルに入る候補部位である。

30

3. 次に、エラーテーブル部位が、フィルタリングされる。フィルタリングにおける例示的ステップは、以下の次の節においてより詳細に説明される。

4. フィルタは、部位のいくつかをエラーテーブルから除去する。フィルタリング後、部位は、Single Nucleotide Polymorphism Database「dbSNP」（および潜在的に、複数のdbSNP Variant Caller Formats「VCF」）と比較される。dbSNP内に生じ、共通である、任意の部位は、エラーテーブルから除去されることができる。

5. エラーテーブルは、大容量JSONファイルとしてディスクに書き込まれ、部位毎の記録は、基準塩基および各代替塩基の頻度を示す。例えば、1%を上回る頻度を伴う任意の代替塩基は、フィルタリングされてもよい。フィルタリングのためのカットオフは、システム自体内で構成可能であることができ、したがって、エラーテーブル内であっても、フィルタリングを保証するために十分ではない。しかしながら、カットオフは、非常に類似する。例えば、プロセスは、エラーテーブル内にある1.5%を上回る頻度を伴う任意のものをフィルタリングすることができる。

40

## 【0037】

エラーテーブルは、着目領域毎に1回生成され、次いで、後の使用のために記憶されることができる。

## 【0038】

エラーテーブルフィルタリング統計

## 【0039】

50

エラーテーブル生成プロセスのステップ3（前述）に記載のように、高相違部位は全て、エラーテーブルに対する候補である。候補部位は、一連の統計的試験を通して（ならびにdbSNPとの比較を通して）フィルタリング除去されることができる。以下は、2つの例示的試験を含む、候補エラーテーブル部位をフィルタリングするために使用される例示的手順を説明する。

#### 【0040】

最初に、各部位毎に、Hardy-Weinberg試験統計が、算出されることができる。これは、非常にネイティブなジェノタイピングによって行われることができる。例えば、塩基が、リードの20%未満のサンプル内に認められる場合、ホモ接合型基準（「HOM REF」）と見なされ、リードの20%～75%に認められる場合、ヘテロ接合型（「HET」）と見なされ、リードの75%を上回って認められる場合、ホモ接合型代替（「HOM ALT」）と見なされる。次いで、サンプルは、これらの3つのカテゴリ（HOM REF、HET、およびHOM ALT）にピン化され、Hardy-Weinberg試験は、0.5%のアルファに対して標準的カイ二乗統計を使用して行われる。したがって、エラーテーブル内の本部位が実在のSNPから由来し得る可能性が存在する場合、エラーテーブルからの除去が検討される。

10

#### 【0041】

しかしながら、これらの部位は、本実施例では、エラーテーブルから直ちに除去されない。エラーテーブルから除去されるためにはまた、Bayes Factor試験にも合格しなければならない。Bayes Factor試験は、以下のように、2つの異なるモデル、すなわち、SNPモデルおよび雑音モデルを前提として、データの確率の比率を算出する。

20

#### 【数6】

$$B = P(\text{data} | \text{SNP model}) / P(\text{data} | \text{noise model})$$

#### 【0042】

Bayes Factorが高い（例えば、10を上回る）場合、データは、SNPモデルに由来するより高い確率を有し、したがって、部位は、エラーテーブルから除去される。

#### 【0043】

2つのモデルは、リードフラクション分布のモデルである。対立遺伝子の頻度が20%である場合、対立遺伝子は、雑音であり得、サンプル内の頻度の分布は全て、約20%となるであろう。すなわち、各サンプル内において、リードの約20%は、本対立遺伝子を有するであろう。代替として、対立遺伝子は、実在し得、その場合、いくつかのサンプルは、100%に近い対立遺伝子を有し、いくつかのサンプルは、0%を有し、いくつかのサンプルは、50%を有するであろう（HOM ALT、HOM REF、およびHETに対応する）。

30

#### 【0044】

これらの2つのモデルは、異なる数のパラメータを有する。概して、雑音モデルでは、リード内で雑音を観察する確率（観察される対立遺伝子頻度に対応する）が、必要とされ、SNPモデルでは、HOM ALT、HOM REF、およびHETサンプルの確率（これらの2つは、1つに総和されなければならないため、2つのみのパラメータ）が、必要とされる。モデルと異なる数のパラメータを比較するために、パラメータは、積分されることができる。したがって、 $P(\text{data} | \text{noise model})$ を算出するために、プロセスは、雑音確率の全ての可能な値（0から1）にわたって $P(\text{data} | \text{noise model}, \text{noise probability})$ を積分することができる。同様に、 $P(\text{data} | \text{SNP model})$ を算出するために、プロセスは、hom refおよびhet proportionsの全ての可能な値にわたって $P(\text{data} | \text{SNP model}, \text{hom ref proportion}, \text{het proportion})$ を積分することができる（hom alt proportionは、1マイナ

40

50

スそれらの2つである)。(積分の面積は、それらの3つの総和がちょうど1であって、それらのいずれも[0, 1]範囲外にないように制約される。)本積分は、Scientific Python「SciPy」数値積分関数(または均等物)を使用して実装されることができる。

#### 【0045】

これらモデル(雑音およびSNPモデル)は両方とも、リードがある種類のBernoulli分布から求められている、すなわち、プロセスは、ある確率 $p$ を用いて、問題の対立遺伝子を認めるかどうかという仮定に基づく。雑音モデルに関して、 $p$ は、パラメータ(雑音確率)であって、プロセスは、その $p$ にわたって積分する。確率 $P(\text{data} | \text{noise model}, p)$ は、二項分布確率質量関数を使用することによって算出されることができ、 $p$ は、プロセスが当該対立遺伝子を認める確率である。PMFに対する $x$ および $n$ パラメータは、単に、対立遺伝子が認められた回数およびサンプル内の合計リード数である。これは、所与のサンプルの確率を算出することを可能にし、データセット内の全サンプルにわたって全それらの確率をとともに乗算することは、パラメータ $p$ を前提としたモデルの全体的確率を提供する。(注記:例示的計算におけるアンダーフローを回避するために、プロセスは、各確率を10で乗算してもよい。したがって、算出される確率は、 $10^N$ でスケールされ、 $N$ は、データセット内のサンプルの数である。)

10

#### 【0046】

SNPモデルに関して、例示的プロセスは、サンプルがHOM REFである可能性に関するものと、HETに関するものと、HOM ALTに関するものの3つの二項分布を含む。しかしながら、各場合において、プロセスは、サンプルがHOM REFまたはHOM ALTである場合でも、汚染が、依然として、ある基準をもたらし得るため、確率 $p$ を把握しない。同様に、HETの場合も、汚染および他の影響(マッピング品質等)が、正確に50%ではない $p$ をもたらし得る。これに対処するために、プロセスは、 $p$ をベータ分布を伴う無作為変数であるようにし得る。すなわち、 $p$ の全ての可能な値にわたって積分することは、ベータ二項分布を与え、これは、SNPモデル内のこれらの3つの場合における単純二項の代わりに使用され得る。事前確率情報(HOM REF、HET、またはHOM ALTである)をモデル化するために、プロセスは、ベータ事前確率のために、アルファおよびベータパラメータを使用して、分布を適切に我々の歪ませることができる。HOM REFおよびHOM ALTの場合、プロセスは、アルファ=20およびベータ=1(またはその逆)を使用して、図3Aに示されるもののようなプロットをもたらしすることができる。HETの場合、プロセスは、アルファ=20およびベータ=20を使用して、図3Bに示されるもののようなプロットをもたらしすることができる。

20

30

#### 【0047】

Bayes Factor試験に不合格の任意の部位は、Hardy-Weinberg比例計算において生じた雑音であると仮定され、したがって、エラーテーブル内に保たれる。

#### 【0048】

Bayes Factor試験に加え、一実施例では、部位がエラーテーブルから外されるためには、Strand Bias試験に合格しなければならない。Strand Bias試験は、非常に単純である。すなわち、基準に関するリードおよび対立遺伝子に関するリードが、どのストランドのトラックがカウントされているのかを維持しながら、全サンプルにわたって集めされる。全体的対立遺伝子頻度 $p$ もまた、算出される。次いで、フォワードリードの確率を算出し(それらが確率 $p$ を用いた二項分布に由来すると仮定する)、バックワードリードに関する同一確率を算出する。それらの確率の比が、非常に高いか、または非常に低い場合、対立遺伝子の分布が一方のストランドまたは他方のストランドに向かって非常にバイアスされていることを示す。したがって、その比の対数がある閾値を上回る(例えば、10を上回る)大きさを有する場合、部位は、ストランドバイアスされていると見なされ、エラーテーブル内に含まれる。

40

#### 【0049】

50

故に、一実施例では、部位が、Hardy - Weinberg 試験、Bayes Factor 試験、および Strand Bias 試験に合格する場合、エラーテーブル候補部位から除去される。

#### 【0050】

種々の他の試験または試験の組み合わせも、エラーテーブルを生成（またはフィルタリング）するために採用されてもよいことを認識されたい。さらに、他の変数または閾値も、本明細書に説明される実施例とともに採用され、シーケンサエラーと実在バリエーションとの間の差異を判定してもよい。

#### 【0051】

コマンドラインインターフェース：

10

#### 【0052】

以下の節は、例示的バリエーションコーラーと、それとともに提供され得るツールの実践的インストールおよび使用を説明する。本明細書に説明される例示的バリエーションコーラーは、標準的 Python パッケージとして実装されることができ（一実施例では、唯一の依存物は、配列アライメントのための C++ ライブラリ `seqan` である）。当然ながら、当業者は、他のプログラミング言語、データフォーマット、および同等物も、可能性として考えられ、検討されることを認識するであろう。

#### 【0053】

一実施例では、例示的バリエーションコーラーは、エラー補正のために、事前に構築されたエラーテーブル（例えば、本明細書に説明されるように）に依拠する。エラーテーブルを生成するために、プロセスは、コーリングのための領域に関するデータを伴う複数のサンプル（例えば、数百サンプルまたはそれを上回る）を収集する。エラーテーブルが、次いで、以下の例示的コマンドを介して、具体的領域に関して生成されることができ（`chr1:100-200` 等）。

20

#### 【数 7】

```
python -m kcall gen-table
  --reference /path/to/hg19.fa
  --output my_error-table.err
  --from /directory/with/bam/files
  --threads $NTHREADS
  --region chr1:100-200
--dbsnp dbsnp.vcf
```

30

#### 【0054】

代替として、プロセスは、以下の \* .bed ファイルを提供することができる。

#### 【数 8】

```
python -m kcall gen-table
  --reference /path/to/hg19.fa
  --output my_error-table.err
  --from /directory/with/bam/files
  --threads $NTHREADS
  --bed /path/to/my/bedfile.bed
--dbsnp dbsnp.vcf
```

40

#### 【0055】

最後に、ディレクトリの代わりに、\* .bam ファイルのリストを用いて、プロセスは、--from の代わりに、そのリストを提供することができる。

## 【数 9】

```
python -m kcall gen-table
  --reference /path/to/hg19.fa
  --output my_error-table.err
  --from /path/to/list-of-bam-files.txt
  --threads $NTHREADS
  --bed /path/to/my/bedfile.bed
--dbnp dbnp.vcf
```

## 【0056】

10

ユーザが、クラスタ内のいくつかのノードにわたってエラーテーブル生成を並列処理することを所望する場合、プロセスは、`*.bed`ファイル内の領域毎に別個のジョブを引き起こすことができる。プロセスは、次いで、生成された断片の全てを単一テーブルに組み合わせることができる。エラーテーブルが単純 `json` フォーマットであるため、プロセスは、`jq` ツールを使用して、これを行うことができる。

# 全エラーテーブル断片は、`pieces/as json files.cat pieces/* .json | jq -s add > combined__table.json` 内に記憶されると仮定する。

## 【0057】

エラーテーブルが生成されると、プロセスは、以下のコマンドを用いて、`Kcall` バリエーションコーラーを起動することができる。

20

## 【数 10】

```
python -m kcall call
  --reference /path/to/hg19.fa
  --errors my_error-table.json
  --bam /path/to/sample.bam
  --threads $NTHREADS
  --bed /path/to/bed/file.bed
  --output-json output.json
  --output-vcf-full full.vcf
--output-vcf-snp snp.vcf
```

30

## 【0058】

例示的バリエーションコーラーは、前述に示される対応するフラグの下で、少なくとも3つのフォーマット、例えば、`json`、`vcf-snp`、および `vcf-full` において出力を提供することができる。プロセスは、これらのフラグの任意のサブセットを有してもよい。すなわち、いずれも提供されない場合、プロセスは、`vcf-snp` フォーマットを出力し、標準化する。`json` フォーマットは、概して、最も単純であって、単に、ディクショナリを伴う `JSON` ファイルをもたらす、各キーは、領域を記述するストリング（「chr1:100-200」等）であって、値は、無コール理由を記述するストリング（領域がコーリングされない場合）、またはディプロタイプおよび領域に関する配列を提供する信頼キーを伴うディクショナリのいずれかである。`vcf-full` フォーマットは、`VCF` と同一情報を出力し、各領域は、ちょうど1行に対応する。無コールについての情報は、`VCF` から利用可能である（ジェノタイプ `GT` フィールドが、`.` となるであろうため）が、無コール理由は、`JSON` 出力フォーマットから利用可能であることに留意されたい。最後に、`vcf-snp` フォーマットは、個々のハプロタイプコールを介して、出力 `VCF` を分割し、それらが数塩基の分離より近い場合、`SNPS` をともに結合する。これは、`GATK` および `FreeBayes` に類似するコールを生成する。

40

## 【0059】

いったん例示的バリエーションコーラーがコールを生成すると、プロセスは、それらを別の

50

コールのセットと比較することができる。例えば、バリエーションコーラーは、本目的のために、基準ゲノム内のそれらの場所によってインデックス化された塩基毎の差異を見出す、積分比較ツールを含んでもよい。これは、プロセスが、VCFと異なる出力フォーマットを比較することを可能にし、したがって、コールセットは、容易に、FreeBayes、GATK1、またはGATK2コールセットと比較されることができる。2つのVCFを比較するために、以下のコマンドが、使用されることができる。

【数 1 1】

```
python -m kall compare first_vcf.vcf second_vcf.vcf
--reference /path/to/hg19.fa
--output output.diff
--stats output.stats
--name $SAMPLE_NAME
--bed /path/to/bed/file.bed
```

10

【0060】

生成された出力は、前述の2つのタブ分離テーブル(output.diffおよびoutput.stats)内に含有される。これらの2つのTSVファイルは、それぞれ、2つのコールセット間の差異および差異の頻度についてのいくつかの統計を含有する。

【0061】

例示的アーキテクチャおよび処理環境：

20

【0062】

本明細書に説明されるシステムおよびプロセスのある側面および実施例が動作し得る、例示的環境およびシステムが、ここで説明される。図4に示されるように、いくつかの実施例では、本システムは、クライアント-サーバモデルに従って実装されることができる。本システムは、ユーザデバイス102上で実行されるクライアント側部分と、サーバシステム110上で実行されるサーバ側部分とを含むことができる。ユーザデバイス102は、デスクトップコンピュータ、ラップトップコンピュータ、タブレットコンピュータ、PDA、携帯電話（例えば、スマートフォン）、または同等物等の任意の電子デバイスを含むことができる。

【0063】

ユーザデバイス102は、インターネット、イントラネット、または任意の他の有線もしくは無線公共もしくはプライベートネットワークを含み得る、1つまたはそれより多くのネットワーク108を通して、サーバシステム110と通信することができる。ユーザデバイス102上の例示的システムのクライアント側部分は、ユーザ対応入力および出力処理ならびにサーバシステム110との通信等、クライアント側機能性を提供することができる。サーバシステム110は、個別のユーザデバイス102上に常駐する任意の数のクライアントのために、サーバ側機能性を提供することができる。さらに、サーバシステム110は、クライアント対応I/Oインターフェース122と、1つまたはそれより多くの処理モジュール118と、データおよびモデル記憶120と、外部サービスとのI/Oインターフェース116とを含み得る、1つまたはそれより多くのコーラサーバ114を含むことができる。クライアント対応I/Oインターフェース122は、コーラサーバ114のためのクライアント対応入力および出力処理を促進することができる。1つまたはそれより多くの処理モジュール118は、本明細書に説明されるように、種々の問題および候補スコアリングモデルを含むことができる。いくつかの実施例では、コーラサーバ114は、タスク完了または情報取得のために、ネットワーク108を通して、テキストデータベース、サブスクリプションサービス、政府記録サービス、および同等物の外部サービス124と通信することができる。外部サービスとのI/Oインターフェース116は、そのような通信を促進することができる。

30

40

【0064】

サーバシステム110は、コンピュータの1つまたはそれより多くの独立型データ処理

50

デバイスまたは分散型ネットワーク上に実装されることができる。いくつかの実施例では、サーバシステム 110 は、第三者サービスプロバイダ（例えば、第三者クラウドサービスプロバイダ）の種々の仮想デバイスおよび／またはサービスを採用し、サーバシステム 110 の下層コンピューティングリソースおよび／またはインフラストラクチャリソースを提供することができる。

#### 【0065】

コーラサーバ 114 の機能性は、クライアント側部分およびサーバ側部分の両方を含むように図 4 に示されるが、いくつかの実施例では、本明細書に説明されるある機能（例えば、ユーザインターフェース特徴およびグラフィカル要素に関して）は、ユーザデバイス上にインストールされた独立型アプリケーションとして実装されることができる。加えて、システムのクライアント部分とサーバ部分との間における機能性の分割は、異なる実施例において変動することができる。例えば、いくつかの実施例では、クライアント上で実行されるユーザデバイス 102 は、ユーザ対応入力および出力処理機能のみを提供し、システムの全ての他の機能性をバックエンドサーバに委譲する、シンククライアントであることができる。

#### 【0066】

サーバシステム 110 およびクライアント 102 はさらに、例えば、処理ユニット、メモリ（本明細書に説明される機能の一部または全部を実施するための論理またはソフトウェアを含み得る）、および通信インターフェース、ならびに他の従来のコンピュータ構成要素（例えば、キーボード／タッチスクリーン等の入力デバイスおよびディスプレイ等の出力デバイス）を有する、種々のタイプのコンピュータデバイスの任意の 1 つを含んでもよいことに留意されたい。さらに、サーバシステム 110 およびクライアント 102 の一方または両方は、概して、論理（例えば、http ウェブサーバ論理）を含む、またはローカルもしくは遠隔データベースもしくは他のデータおよびコンテンツのソースからアクセスされる、データをフォーマットするようにプログラムされる。この目的を達成するために、サーバシステム 110 は、共通ゲートウェイインターフェース（CGI）プロトコルおよび関連付けられたアプリケーション（または「スクリプト」）、Java（登録商標）「servelets」、すなわち、サーバシステム 110 上で起動する Java（登録商標）アプリケーション、または情報を提示し、クライアント 102 からの入力を受信するための同等物等の種々のウェブデータインターフェース技法を利用してもよい。サーバシステム 110 は、単数形で本明細書に説明されるが、実際には、複数のコンピュータ、デバイス、データベース、関連付けられたバックエンドデバイス、および同等物を備え、（有線および／または無線で）通信し、本明細書に説明される機能の一部または全部を行うように協働してもよい。サーバシステム 110 はさらに、アカウントサーバ（例えば、電子メールサーバ）、モバイルサーバ、メディアサーバ、および同等物を含むか、またはそれと通信してもよい。

#### 【0067】

さらに、本明細書に説明される例示的方法およびシステムは、種々の機能を行うために別個のサーバおよびデータベースシステムの使用を説明するが、他の実施形態も、説明される機能性が行われる限り、設計選択肢上、説明される機能を単一デバイスまたは複数のデバイスの任意の組み合わせ上で生じさせるように動作する、ソフトウェアまたはプログラミングを記憶することによって実装され得ることに留意されたい。同様に、説明されるデータベースシステムも、単一データベース、分散型データベース、分散型データベースの集合、冗長オンラインもしくはオフラインバックアップもしくは他の冗長性を伴うデータベース、または同等物として実装されることができる、分散型データベースまたは記憶ネットワークおよび関連付けられた処理インテリジェンスを含むことができる。図に描写されないが、サーバシステム 110（ならびに本明細書に説明される他のサーバおよびサービス）は、概して、限定ではないが、プロセッサ、RAM、ROM、クロック、ハードウェアドライバ、関連付けられた記憶、および同等物を含む、サーバシステム内で通常見出されるような当該技術で認識される構成要素を含む（例えば、以下に論じられる図 5 参

10

20

30

40

50



照)。さらに、説明される機能および論理は、ソフトウェア、ハードウェア、ファームウェア、またはそれらの組み合わせ内に含まれてもよい。

【0068】

図5は、種々のコーリングおよびスコアリングモデルを含む前述の説明されたプロセスのうちの任意の1つを実施するために構成された例示的コンピューティングシステム1400を図示する。本文脈では、コンピューティングシステム1400は、例えば、プロセッサ、メモリ、記憶装置、および入力/出力デバイス(例えば、モニタ、キーボード、ディスクドライブ、インターネット接続等)を含んでもよい。しかしながら、コンピューティングシステム1400は、プロセスのいくつかまたは全ての側面を実行するための電気回路または他の特殊ハードウェアを含んでもよい。いくつかの動作設定において、コンピューティングシステム1400は、1つまたはそれを上回るユニットを含むシステムとして構成されてよく、それぞれは、ソフトウェア、ハードウェア、またはそのいくつかの組み合わせにおいて、プロセスのいくつかの側面を実行するために構成される。

10

【0069】

図5は、前述の説明されたプロセスを実施するために使用され得る、いくつかの構成要素を有するコンピューティングシステム1400を図示する。メインシステム1402は、入力/出力(「I/O」)セクション1406と、1つまたはそれを上回る中央処理装置(「CPU」)1408と、メモリセクション1410とを有する主回路基板1404を含み、それに関連されるフラッシュメモリカード1412を有してよい。I/Oセクション1406は、ディスプレイ1424、キーボード1414、ディスク記憶装置1416、およびメディアドライブ装置1418に接続される。メディアドライブ装置1418は、コンピュータ可読媒体1420の読み取り/書き込みが可能で、プログラム1422および/またはデータを含むことができる。

20

【0070】

前述の説明されたプロセスの結果に基づく少なくともいくつかの値は、その後の使用のために保存されることができる。加えて、非一過性コンピュータ可読媒体は、コンピュータを用いて、前述の説明されたプロセスのうちの任意の1つを実施するための1つまたはそれを上回るコンピュータプログラムを記憶する(例えば、有形に具現化する)ために使用されることができる。コンピュータプログラムは、例えば、汎用プログラミング言語(例えば、Pascal、C、C++、Python、Java(登録商標))またはある特殊用途専用言語で書き込まれてもよい。

30

【0071】

種々の例示的实施形態が、本明細書に記載される。これらの実施例は、非限定的意味で参照される。それらは、公開された本技術のより広く適用できる側面を例証するために提供される。種々の実施形態の厳密な精神および範囲から逸脱することなく、種々の変更がなされ、また、均等物が代用されてよい。加えて、多くの修正が、特定の状況、材料、組成物、プロセス、プロセス行為、またはステップを、種々の実施形態の目的、精神、または範囲に適合させるためになされてよい。さらに、当業者によって理解されるであろうように、本明細書に記載および例証される個々の変形例はそれぞれ、種々の実施形態の範囲または精神から逸脱することなく、任意の他のいくつかの実施形態の特徴から容易に分離されてよい、またはそれらと併用されてよい個別の構成要素および特徴を有する。全てのそのような修正は、本開示と関連付けられる請求項の範囲内であることが意図される。

40

【図 1】

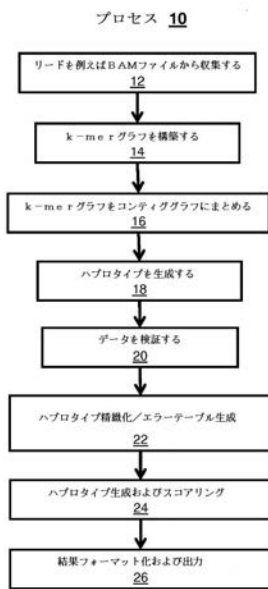
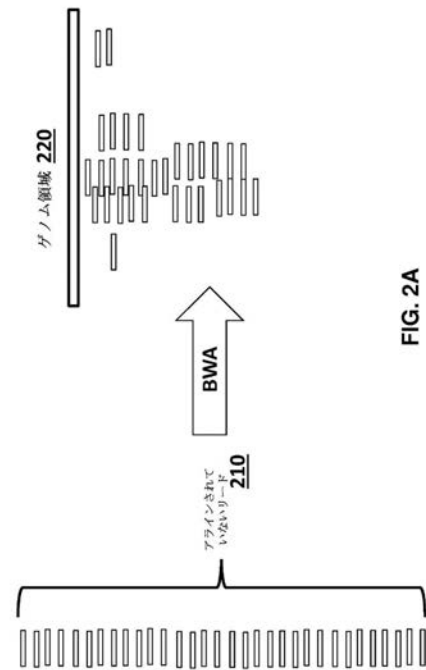
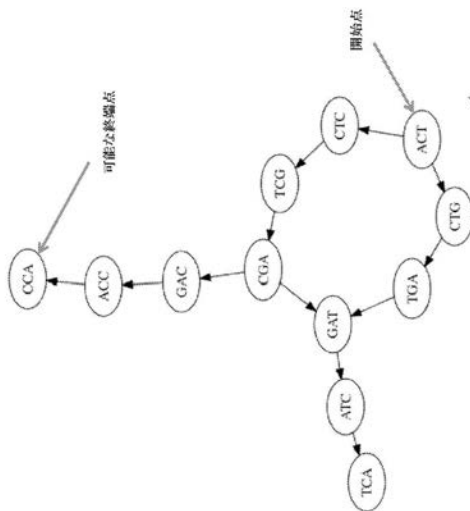


FIG. 1

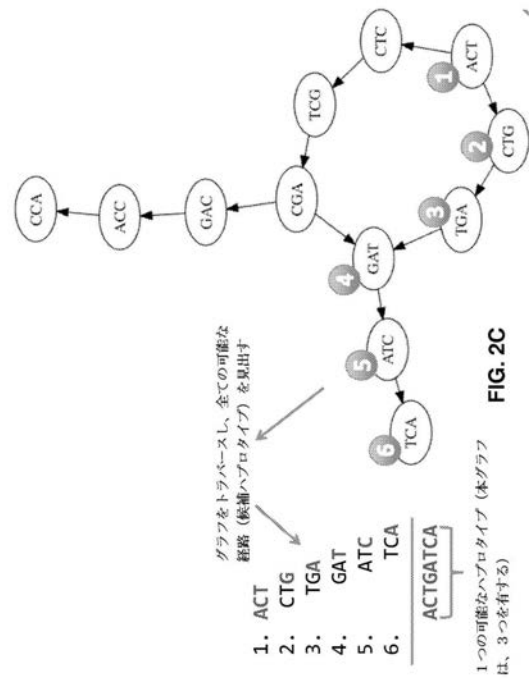
【図 2 A】



【図 2 B】



【図 2 C】



【図 3 A】

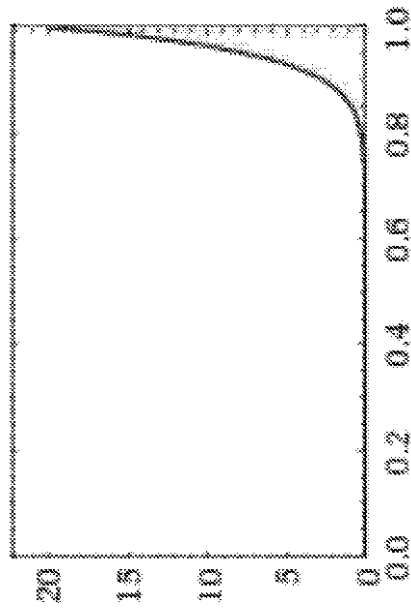


FIG. 3A

【図 3 B】

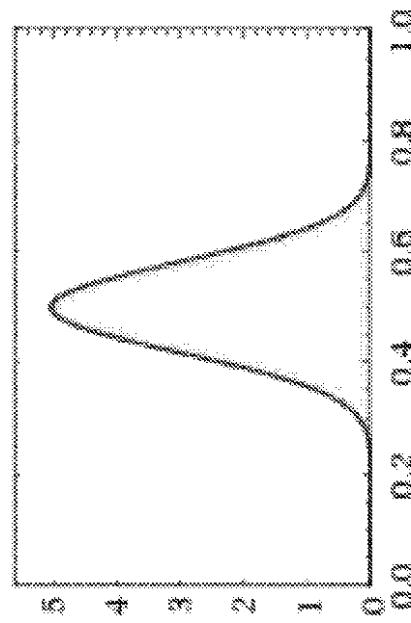


FIG. 3B

【図 4】

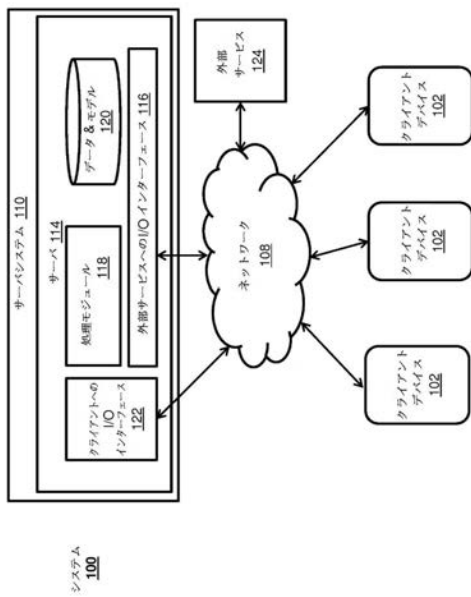


FIG. 4

【図 5】

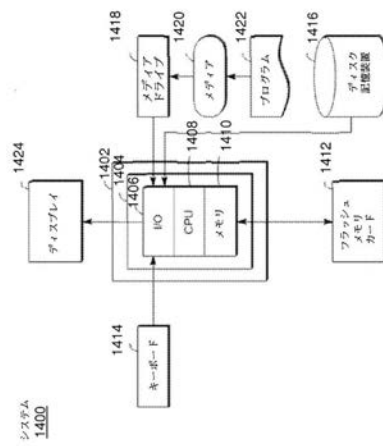


FIG. 5

## 【国際調査報告】

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/US 15/55807

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> IPC(8) - G01N 33/48, C40B 20/00, G06F 19/18 (2016.01) CPC - G06F19/18, C07H 21/00, G06F19/24 According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) IPC(8): G01N 33/48, C40B 20/00, G06F 19/18 (2016.01) CPC: G06F19/18, C07H 21/00, G06F19/24 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched USPC: 702/20, 506/2, 702/20 Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) PatBase, Google Patents, Google Scholar, Google Web, search terms: computer-implemented, genome sample, variants, reference genomic sequence, reads, error table, diplotypes, haplotypes, posterior probability, k-mer graph, Hardy-Weinberg test, Bayes Factor test, Strand Bias Test, sequencer error, output variants, scoring		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	US 2012/0053845 A1 (BRUSTLE et al.) 01 March 2012 (01.03.2012) para [0009], [0010], [0015], [0019], [0024], [0025], [0026], [0028], Fig. 1	1-6, 9-20
Y	US 2013/0054508 A1 (KERMANI et al.) 28 February 2013 (28.02.2013) para [0004], [0048], [0060], [0069], [0217]	1-6, 9-20
Y	ZUO et al. Use of diplotypes-matched haplotype pairs from homologous chromosomes in gene-disease association studies. Shanghai Arch Psychiatry. (June 2014) vol 26, no 3, pp 165-170, abstract, pg 166, col 1, para 2 3, pg 168 col 1 para 2	1-6, 9-20
Y	US 2004/0265816 A1 (TANAKA et al.) 30 December 2004 (30.12.2004) para [0012], [0049], [0071]	3, 11, 17
Y	US 2005/0214811 A1 (MARGULIES et al.) 29 September 2005 (29.09.2016) abstract, para [0073]	6, 14, 20
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/>		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 27 January 2016		Date of mailing of the international search report <b>16 FEB 2016</b>
Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US, Commissioner for Patents P.O. Box 1450, Alexandria, Virginia 22313-1450 Facsimile No. 571-273-8300		Authorized officer: Lee W. Young PCT Helpdesk: 571-272-4300 PCT OSP: 571-272-7774

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/US 15/55807

**Box No. II Observations where certain claims were found unsearchable (Continuation of item 2 of first sheet)**

This international search report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:  
because they relate to subject matter not required to be searched by this Authority, namely:
2. ☐ Claims Nos.:  
because they relate to parts of the international application that do not comply with the prescribed requirements to such an extent that no meaningful international search can be carried out, specifically:
3. ☐ Claims Nos.:  
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

**Box No. III Observations where unity of invention is lacking (Continuation of item 3 of first sheet)**

This International Searching Authority found multiple inventions in this international application, as follows:

\*\*\*\*\* See Supplemental Sheet to continue \*\*\*\*\*

1. ☐ As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims.
2. ☐ As all searchable claims could be searched without effort justifying additional fees, this Authority did not invite payment of additional fees.
3. ☐ As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:
4. ☒ No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:  
1-6, 9-20

**Remark on Protest**

- ☐ The additional search fees were accompanied by the applicant's protest and, where applicable, the payment of a protest fee.
- ☐ The additional search fees were accompanied by the applicant's protest but the applicable protest fee was not paid within the time limit specified in the invitation.
- ☐ No protest accompanied the payment of additional search fees.

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/US 15/55807

Continuation of Box No. III, Observations where unity of invention is lacking:

This application contains the following inventions or groups of inventions which are not so linked as to form a single general inventive concept under PCT Rule 13.1. In order for all inventions to be examined, the appropriate additional examination fees must be paid.

Group I: Claims 1-6, 9-20 drawn to a computer-implemented method or system for determining variants from a genome sample relative to a reference genomic sequence

Group II: Claims 7-8, drawn to a computer-implemented method for generating an error table of sequence data.

The inventions listed as Groups I and II do not relate to a single general inventive concept under PCT Rule 13.1 because, under PCT Rule 13.2, they lack the same or corresponding special technical features for the following reasons:

**Special Technical Features:**

Group I has the special technical feature of ACCESSING AN ERROR TABLE of sequence data from previously sequenced samples, not required by Group II.

Group I has the special technical feature of determining a set of possible haplotypes from a set of collected reads, where the set of possible haplotypes is FILTERED BY THE ERROR TABLE, not required by Group II.

Group II has the special technical feature of aligning a set of collected reads to a reference sample, determining mismatches with the reference table, and ADDING SITES THAT HAVE A MISMATCH TO AN ERROR TABLE, not required by Group I.

**Common Technical Feature:**

Groups I and II have the common technical feature of an error table, where Group I accesses and uses a preconstructed error table, and Group II constructs and generates an error table for future use.

However, said common technical feature does not represent a contribution over the prior art, and is anticipated by US 2013/0332081 A1 to REESE et al. (hereinafter "Reese").

As to the common technical feature, Reese teaches (para [0116-0117]: "Variant Masking. In one aspect, disclosed herein are methods of identifying and/or prioritizing phenotype causing variants by comparing a target cohort to a background cohort, wherein variants within certain regions of the reference sequence are excluded or masked from a CLRT [i.e. composite likelihood ratio test]. Variant calling can be error-prone in repetitive and/or homologous sequences. Variant calling errors in repetitive or homologous sequences can be caused by short sequence reads aligning to multiple sites within the reference sequence. Variant masking can decrease the effect of sequencing platform bias (FIG. 3B). A CLRT with variant masking can be performed with AAS information and/or variant frequency estimation"; para [0153]: "The variant masking option allows the user to exclude a list of nucleotide sites from the likelihood calculations based on information obtained prior to the genome analysis. The masking files used in these analyses excludes sites where short reads would map to more than one position in the reference genome. This procedure mitigates the effects introduced by cross-platform biases by excluding sites that are likely to produce spurious variant calls due to improper alignment of short reads to the reference sequence. The three masking schemes employed are a) 60-bp single-end reads, b) 35-bp single-end reads, and c) 35-bp paired-end reads separated by 400-bp. These three masking files are included with the VAAST distribution")

As the common technical feature was known in the art at the time of the invention, this cannot be considered a common special technical feature that would otherwise unify the groups. The inventions lack unity with one another.

Therefore, Groups I and II lack unity of invention under PCT Rule 13 because they do not share a same or corresponding special technical feature

NOTE, Claims 16-18 depends from "the system of claim 9", as drafted, is objected to, because claim 9 does not teach "a system". For this International Search and Opinion, claims 16-18 are re-constructed as dependent claims of claim 15.

## フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(特許庁注：以下のものは登録商標)

1. J A V A S C R I P T

- (72)発明者 ギピアンスキー, アンドリュー レオニドヴィッチ  
アメリカ合衆国 カリフォルニア 94080, サウス サンフランシスコ, キンボール ウ  
エイ 180, カウンシル, インコーポレイテッド 気付
- (72)発明者 ハケ, イムラン サイーダル  
アメリカ合衆国 カリフォルニア 94080, サウス サンフランシスコ, キンボール ウ  
エイ 180, カウンシル, インコーポレイテッド 気付
- (72)発明者 マグワイア, ジャレッド ロバート  
アメリカ合衆国 カリフォルニア 94080, サウス サンフランシスコ, キンボール ウ  
エイ 180, カウンシル, インコーポレイテッド 気付
- (72)発明者 ロバートソン, アレクサンダー デ ヨング  
アメリカ合衆国 カリフォルニア 94080, サウス サンフランシスコ, キンボール ウ  
エイ 180, カウンシル, インコーポレイテッド 気付