

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7157141号  
(P7157141)

(45)発行日 令和4年10月19日(2022.10.19)

(24)登録日 令和4年10月11日(2022.10.11)

(51)国際特許分類 F I  
G 1 6 B 50/50 (2019.01) G 1 6 B 50/50

請求項の数 10 (全32頁)

(21)出願番号	特願2020-509515(P2020-509515)	(73)特許権者	390009531 インターナショナル・ビジネス・マシ ンズ・コーポレーション INTERNATIONAL BUSI NESS MACHINES CORPO RATION アメリカ合衆国10504 ニューヨ ーク州 アーモンク ニュー オーチャード ロード New Orchard Road, A rmonk, New York 105 04, United States of America
(86)(22)出願日	平成30年8月9日(2018.8.9)	(74)代理人	100112690 弁理士 太佐 種一
(65)公表番号	特表2020-533666(P2020-533666 A)		
(43)公表日	令和2年11月19日(2020.11.19)		
(86)国際出願番号	PCT/IB2018/056009		
(87)国際公開番号	WO2019/043481		
(87)国際公開日	平成31年3月7日(2019.3.7)		
審査請求日	令和3年1月22日(2021.1.22)		
(31)優先権主張番号	15/693,019		
(32)優先日	平成29年8月31日(2017.8.31)		
(33)優先権主張国・地域又は機関	米国(US)		

最終頁に続く

(54)【発明の名称】 ゲノム・ファイルのためのコンテキスト・アウェア差分アルゴリズム

(57)【特許請求の範囲】

【請求項1】

複数のゲノム・データ・ファイルに対する少なくとも1つの差分ファイルを圧縮する方法であって、前記方法は、

インプットとして第1のファイル及び第2のファイルを受信するステップであって、該第1のファイル及び該第2のファイルは、前記複数のゲノム・データ・ファイルに属し、該第1のファイルの各行が第1のファイルのフィールドからなり、該第2のファイルの各行が第2のファイルのフィールドからなるように、該第1のファイル及び該第2のファイルの両方がタブ区切りであり、ここで、該第1のファイルが、第1の標本に関連付けられた参照ゲノムからの第1のデータを含み、該第2のファイルが、第1の標本と異なる、第2の標本に關  
連付けられた別のゲノムからの第2のデータを含む、ステップと、

10

前記受信された前記第1のファイル及び受信された前記第2のファイルを精査することにより、受信された第1のファイルのソートされた行及び受信された第2のファイルの対応するソートされた行を判断し、該ソートされた行及び該対応するソートされた行は、ゲノムにおけるマッピング位置の昇順に配置されるステップと、

前記ソートされた行の第1のファイルのフィールドと、前記対応するソートされた行の対応する第2のファイルのフィールドとを対比するステップと、

前記対比に基づいて複数の派生差分ファイルを生成するステップと、

汎用ファイル・コンプレッサを用いて、前記生成された複数の派生差分ファイルを圧縮するステップであって、前記複数の派生差分ファイルは、第2のファイルを再構築するた

20

めに第1のファイルと共に利用されるようにフォーマットされる、圧縮するステップと、を含む、方法。

【請求項2】

前記圧縮された複数の派生差分ファイルをユーザのデバイスに格納するステップと、前記圧縮された複数の派生差分ファイルを前記ユーザに提示するステップと、をさらに含む、請求項1に記載の方法。

【請求項3】

前記第1のファイルは、少なくとも1つのソース・ファイルから構成され、前記第2のファイルは、少なくとも1つの目的ファイルから構成される、請求項1に記載の方法。

【請求項4】

前記受信された第1のファイルおよび前記受信された第2のファイルがソートされていないことを判断するステップと、

前記受信された第1および第2のファイル内の各行を、ソーティング・ツールを用いてソートするステップと、

をさらに含む、請求項3に記載の方法。

【請求項5】

前記受信された第1および第2のファイルが、対応可能なフォーマットであることを判断するステップ、

をさらに含む、請求項1に記載の方法。

【請求項6】

前記受信された第1のファイルのソートされた行の第1のファイルのフィールドと受信された第2のファイルのソートされた行の対応する第2のファイルのフィールドとを対比する準備として、受信された第2のファイルの対応するソートされた行の少なくとも2つの行をバッファに格納すること、

をさらに含む、請求項1に記載の方法。

【請求項7】

前記第1のファイルのフィールドがそれぞれの第1のコラムを含み、前記第2のファイルのフィールドがそれぞれの第2のコラムを含み、前記ソートされた行の第1のファイルのフィールドと、前記対応するソートされた行の対応する第2のファイルのフィールドとを対比するステップが、

前記第1のコラムと前記第2のコラムの対応するものを対比するステップと、

前記第1のコラムと前記第2のコラムの対応するものとの対比に基づいて、派生差分ファイル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。

【請求項8】

前記第1のコラムと前記第2のコラムの対応するものを対比するステップが、

前記第1のコラムの優先コラムが前記第2のコラムの特定のコラムと一致するかどうかを判断するステップと、

前記第1のコラムの優先コラムが前記第2のコラムの特定のコラムに一致するという判断に回答して、前記第1のコラムの追加のコラムを前記第2のコラムの対応する更なるコラムと対比するステップと、

をさらに含む、請求項7に記載の方法。

【請求項9】

請求項1～8のいずれか1項に記載の前記方法の各ステップをコンピュータ・ハードウェアによる手段として構成したシステム。

【請求項10】

請求項1～8のいずれか1項に記載の前記方法の各ステップをコンピュータに実行させる、コンピュータ・プログラム。

【発明の詳細な説明】

【技術分野】

10

20

30

40

50

## 【 0 0 0 1 】

本発明は、一般にコンピューティングの分野に関し、さらに具体的には計算生物学に関する。

## 【背景技術】

## 【 0 0 0 2 】

ゲノム解析パイプライン（パイプライン）は、生の配列リードから生物学的に意味のあるアウトプットを抽出するために、前処理、変異の発見、およびコール・セットのリファインメントなど多くのステップを伴う。パイプラインは、かかるステップ毎にアウトプット・ファイルを生成し、それらは、インプット配列リードのサイズの如何によって、メガバイト～テラバイトの間の範囲のサイズとなる。

10

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【 0 0 0 3 】

これらのファイルの見直しによって、あらゆるステップでこれらアウトプット・ファイルに保存された全ての情報が新規に生成されたものでないことが判明している。データのかなり大きなパーセントが、インプット・ファイルからアウトプットに単に移されているだけで、パイプラインの実行の過程で、さらにはその後も、ストレージに対し不必要な負担を生じさせている。パイプラインの各段階は多くの時間または日数を要し得るので、これら中間的なファイルが、パイプラインにおける今後の調査、変更、または分岐のために保存される。

20

## 【課題を解決するための手段】

## 【 0 0 0 4 】

本発明の諸実施形態は、複数のゲノム・データ・ファイルに対する少なくとも1つの差分ファイルを圧縮するための方法、コンピュータ・システム、およびコンピュータ・プログラム製品を開示する。本発明は、インプットとして複数のゲノム・データ・ファイルを受信するステップを含むことができる。また、本発明は、受信された複数のゲノム・データ・ファイルを精査することによって、複数の行を判断するステップを含むことが可能である。本発明は、次いで、これら精査された複数のゲノム・データ・ファイルに関連付けられた複数の行を対比するステップを含むことができる。本発明は、対比された複数の行に基づいて、複数の派生差分ファイルを生成するステップをさらに含むことが可能である。また、本発明は、汎用ファイル・コンプレッサを用いて、生成された複数の派生差分ファイルを圧縮するステップを含むことができる。

30

## 【 0 0 0 5 】

本発明のこれらのおよび他の目的、特徴、および利点は、以降の本発明の例示的な諸実施形態の詳細な説明を添付の図面と併せ読むことによって明らかとなる。これら図面中の様々な特徴は、当業者が詳細な説明と併せ本発明を理解するのを容易にする上での明瞭化のための例示である。

## 【図面の簡単な説明】

## 【 0 0 0 6 】

【図1】少なくとも1つの実施形態による、ネットワーク化コンピュータ環境を示す。

40

【図2】少なくとも1つの実施形態による、ゲノム・データ・ファイルに対する差分ファイルを圧縮するための或るプロセスを示すオペレーション・フローチャートである。

【図3】少なくとも1つの実施形態による、ゲノム・データ・ファイルに対する差分ファイルを圧縮するための典型的なプロセスを示すオペレーション・フローチャートである。

【図4】少なくとも1つの実施形態による、配列アラインメント/マップ・フォーマット（SAM：Sequence Alignment / Map format）のゲノム・データ・ファイルに対する差分ファイルを圧縮するための典型的なプロセスを示すオペレーション・フローチャートである。

【図5】少なくとも1つの実施形態による、変異コール・フォーマット（VCF：Variant Call Format）のゲノム・データ・ファイルに対する差分ファイルを

50

圧縮するための典型的なプロセスを示すオペレーション・フローチャートである。

【図 6】少なくとも 1 つの実施形態による、2 つのゲノム・データ・ファイルを対比するための典型的なプロセスを示すブロック図である。

【図 7】少なくとも 1 つの実施形態による、配列アラインメント/マップ・フォーマット (SAM) のゲノム・データ・ファイルの階層構造を識別するための典型的なプロセスを示すオペレーション・フローチャートである。

【図 8】少なくとも 1 つの実施形態による、変異コール・フォーマット (VCF) のゲノム・データ・ファイルの階層構造を識別するための典型的なプロセスを示すオペレーション・フローチャートである。

【図 9】少なくとも 1 つの実施形態による、図 1 に描かれたコンピュータおよびサーバの内部および外部のコンポーネントのブロック図である。

10

【図 10】本開示の或る実施形態による、図 1 に描かれたコンピュータ・システムを含む例示のクラウド・コンピューティング環境のブロック図である。

【図 11】本開示の或る実施形態による、図 8 の例示のクラウド・コンピューティング環境の機能層のブロック図である。

【発明を実施するための形態】

【0007】

請求対象の構造体および方法の詳細な実施形態が、本明細書で開示されるが、但し、当然のことながら、これら開示される実施形態は、様々な形態で具現化が可能な請求対象の構造体および方法の単なる例示である。しかしながら、本発明は多くの異なる形態で具現化することが可能で、本明細書中で述べる例示的な実施形態に限定されると解釈されるべきではない。むしろ、これらの例示的な実施形態は、本開示が、当業者によって徹底的且つ完全に理解され、本発明の範囲が十分に伝達されるようにするため提供される。本説明中で、提示された実施形態を、不必要に曖昧にするのを避けるために、周知の機能および技法の詳細は省略されることがある。

20

【0008】

本発明は、任意の可能な技術的詳細の集約度で、システム、方法、もしくはコンピュータ・プログラム製品またはこれらの組合せとすることができる。このコンピュータ・プログラム製品は、プロセッサに本発明の態様を実行させるためのコンピュータ可読プログラムを有するコンピュータ可読ストレージ媒体（または媒体群）を含むことが可能である。

30

【0009】

このコンピュータ可読ストレージ媒体は、命令実行デバイスが使用するための命令を保持し格納できる有形のデバイスとすることができる。該コンピュータ可読ストレージ媒体は、例えば、以下に限らないが、電子ストレージ・デバイス、磁気ストレージ・デバイス、光ストレージ・デバイス、電磁気ストレージ・デバイス、半導体ストレージ・デバイス、または前述のデバイスの任意の適切な組合せであってよい。コンピュータ可読ストレージ媒体のさらに具体的な例の非包括的リストには、携帯型コンピュータ・ディスク、ハード・ディスク、ランダム・アクセス・メモリ (RAM)、読み取り専用メモリ (ROM)、消去およびプログラム可能読み取り専用メモリ (EPROM: erasable programmable read-only memory またはフラッシュ・メモリ)、静的ランダム・アクセス・メモリ (SRAM: static random access memory)、携帯型コンパクト・ディスク読み取り専用メモリ (CD-ROM: compact disc read-only memory)、デジタル多用途ディスク (DVD: digital versatile disk)、メモリ・スティック、フレキシブル・ディスク、パンチカードまたは記録された命令を有する溝中の嵩上げ構造体などの機械的符号化デバイス、および前述の任意の適切な組合せが含まれる。本明細書で用いられるコンピュータ可読ストレージ媒体は、無線波または他の自由に伝播する電磁波、ウェーブガイドまたは他の送信媒体（例えば、光ファイバ・ケーブルを通過する光パルス）、またはワイヤを通して送信される電気信号などの、本質的に一時的な信号であると解釈されるものではない。

40

50

## 【 0 0 1 0 】

本明細書に述べられるコンピュータ可読プログラム命令は、コンピュータ可読ストレージ媒体から、それぞれのコンピューティング/処理デバイスに、または、例えばインターネット、ローカル・エリア・ネットワーク、広域ネットワークもしくはワイヤレス・ネットワークまたはこれらの組合せなどのネットワークを介して、外部のコンピュータもしくは外部のストレージ・デバイスにダウンロードすることが可能である。このネットワークには、銅送信ケーブル、光送信ファイバ、ワイヤレス通信、ルータ、ファイアウォール、スイッチ、ゲートウェイ・コンピュータ、もしくはエッジ・サーバまたはこれらの組合せが含まれてよい。それぞれのコンピューティング/処理デバイス中のネットワーク・アダプタ・カードまたはネットワーク・インターフェースは、ネットワークからコンピュータ可読プログラム命令を受信し、そのコンピュータ可読プログラム命令を、ストレージのため、それぞれのコンピューティング/処理デバイス内のコンピュータ可読ストレージ媒体中に転送する。

10

## 【 0 0 1 1 】

本発明のオペレーションを実行するためのコンピュータ可読プログラム命令は、アセンブラ命令、命令集合アーキテクチャ (ISA: instruction-set-architecture) 命令、マシン命令、マシン依存命令、マイクロコード、ファームウェア命令、状態設定データ、集積回路の構成データ、または、Smalltalk、C++などのオブジェクト指向プログラミング言語、および「C」プログラミング言語もしくは類似のプログラミング言語などの従来式の手続き型プログラミング言語を含む、1つ以上のプログラミング言語の任意の組合せで記述されたソース・コードもしくはオブジェクト・コードであってよい。このコンピュータ可読プログラム命令は、スタンドアロン・ソフトウェア・パッケージとしてユーザのコンピュータで専ら実行することも、ユーザのコンピュータで部分的に実行することもでき、一部をユーザのコンピュータで一部を遠隔コンピュータで実行することもでき、あるいは遠隔のコンピュータまたはサーバで専ら実行することもできる。後者の場合は、ローカル・エリア・ネットワーク (LAN: local area network) または広域ネットワーク (WAN: wide area network) を含む任意の種類ネットワークを介して、遠隔コンピュータをユーザのコンピュータに接続することもでき、あるいは (例えばインターネット・サービス・プロバイダを使いインターネットを介し) 外部のコンピュータへの接続を行うことも可能である。いくつかの実施形態において、例えば、プログラム可能論理回路、フィールドプログラム可能ゲート・アレイ (FPGA: field-programmable gate array)、またはプログラム可能論理アレイ (PLA: programmable logic array) を含む電子回路は、本発明の諸態様を実行すべく、該電子回路をカスタマイズするためコンピュータ可読プログラム命令の状態情報を利用することによって、該コンピュータ可読プログラム命令を実行することができる。

20

30

## 【 0 0 1 2 】

本発明の諸態様は、本発明の諸実施形態による方法、装置 (システム)、およびコンピュータ・プログラム製品のフローチャート図もしくはブロック図またはその両方を参照しながら本明細書で説明される。当然のことながら、フローチャート図もしくはブロック図またはその両方の各ブロック、およびフローチャート図もしくはブロック図またはその両方のブロックの組合せは、コンピュータ可読プログラム命令によって実装することが可能である。

40

## 【 0 0 1 3 】

これらのコンピュータ可読プログラム命令を、汎用コンピュータ、特殊用途コンピュータ、またはマシンを形成する他のプログラム可能データ処理装置のプロセッサに提供し、そのコンピュータまたは他のプログラム可能データ処理装置のプロセッサを介して実行されるこれらの命令が、フローチャートもしくはブロック図またはその両方のブロックもしくはブロック群中に特定されている機能群/動作群を実装するための手段を生成することができる。また、コンピュータ、プログラム可能データ処理装置、もしくは他

50

のデバイスまたはこれらの組合せに対し特定の仕方でも機能するように命令することが可能なこれらのコンピュータ可読プログラム命令を、コンピュータ可読ストレージ媒体に格納し、格納された命令を有するコンピュータ可読ストレージ媒体が、フローチャートもしくはブロック図またはその両方のブロックまたはブロック群中に特定されている機能/動作の諸態様を実装する命令群を包含する製造品を構成するようにすることができる。

【0014】

さらに、これらコンピュータ可読プログラム命令を、コンピュータ、他のプログラム可能データ処理装置、または他のデバイスにロードし、そのコンピュータ上で、他のプログラム可能装置上で、または他のデバイス上で一連のオペレーション・ステップを実施させて、コンピュータ実装のプロセスを作り出し、当該コンピュータ上で、他のプログラム可能装置上でもしくは他のデバイス上で実行される命令が、フローチャートもしくはブロック図またはその両方のブロックもしくはブロック群中に特定されている機能群/動作群を実装するようにすることも可能である。

10

【0015】

諸図面中のフローチャートおよびブロック図は、本発明の様々な実施形態による、システム、方法、およびコンピュータ・プログラム製品から可能となる実装のアーキテクチャ、機能性、およびオペレーションを表している。この点に関し、フローチャートまたはブロック図中の各ブロックは、特定の論理機能(群)を実装するための一つ以上の実行可能命令を含む、モジュール、セグメント、または命令の部分を表し得る。一部の別の実装においては、ブロック中に記載された機能が、図面に記載された順序から外れて行われてよい。例えば、連続して示された2つのブロックが、関与する機能性に応じ、実際にはほぼ同時に実行されることがあり、時にはこれらのブロックが逆の順序で実行されることもあり得る。さらに、ブロック図もしくはフローチャート図またはその両方の各ブロック、およびブロック図もしくはフローチャート図またはその両方中のブロック群の組合せは、特定の機能または動作を実施する特殊用途ハードウェア・ベースのシステムによって実装でき、または特殊用途ハードウェアとコンピュータ命令との組合せによって実行できることにも留意すべきである。

20

【0016】

以降で説明する例示的な諸実施形態は、ゲノム・データ・ファイルに対する差分ファイルを圧縮するためのシステム、方法、およびプログラム製品を提供する。しかして、本諸実施形態は、ゲノム解析パイプラインにより生成される中間的なファイルのストレージ必要容量を低減することによって、計算生物学の技術分野を改良する効果を有する。さらに具体的には、少なくとも2つのゲノム・データ・ファイル(例えば、少なくとも1つのソース・ファイルおよび1つの目的ファイル)を、インプットとしてゲノムベース差分圧縮プログラムに投入することができ、次いで、ゲノムベース差分圧縮プログラムが、ソート・アルゴリズムを含み得る既存のオープン・ソース・ツールを利用してゲノム・データ・ファイルの行をソートすることができる。ゲノムベース差分圧縮プログラムは、次いで、ソース・ファイルおよび目的ファイルの各行を逐次的に精査することが可能で、ソース・ファイルと目的ファイルと間に対応する行を対比し、派生差分ファイルを生成することが可能である。次いで、ゲノムベース差分圧縮プログラムは、既存の汎用ファイル・コンプレッサを用いて、生成された派生差分ファイルを圧縮することができる。

30

40

【0017】

前述したように、ゲノム解析パイプライン(パイプライン)は、生の配列リードから生物学的に意味の或るアウトプットを抽出するために、前処理、変異の発見、およびコール・セットのリファインメントなど多くのステップを伴う。該パイプラインは、かかるステップ毎にアウトプット・ファイルを生成し、それらは、インプット配列リードのサイズの如何によって、メガバイト~テラバイトの間の範囲のサイズになる。

【0018】

これらのファイルの見直しによって、あらゆるステップでこれらアウトプット・ファイルに保存された全ての情報が新規に生成されたものでないことが判明している。データの

50

かなり大きなパーセントが、インプット・ファイルからアウトプットに単に移されているだけで、パイプラインの実行の過程で、さらにはその後も、ストレージに対し不必要な負担を生じさせている。パイプラインの各段階は多くの時間または日数を要し得るので、これら中間的なファイルが、パイプラインにおける今後の調査、変更、または分岐のために保存される。

**【 0 0 1 9 】**

デルタまたは差分圧縮アルゴリズムが存在し得るが（例えば、`x d e l t a 3`または`v c d i f f`）、但し、これらの圧縮アルゴリズムは、ソース・ファイルと目的ファイルとの間に、比較的に大きく連続した共通の部分列がある場合により効果的であり得る。配列アラインメント/マップ・フォーマット（`S A M`）および変異コール・フォーマット（`V C F`）などの或る種のフォーマットによるゲノム・データは、フィールドから成る比較的小さな記録を含むことがあり、一部のフィールドおよびサブフィールドが1つのバージョンと次のバージョンとで異なることがある。

10

**【 0 0 2 0 】**

したがって、小さな記録群に対してゲノム解析パイプラインによって生成された中間的ファイルのストレージ必要容量、または繰り返し情報の量を最小化するためのデータ低減方法を考案することが、なかんずく有益であり得る。

**【 0 0 2 1 】**

少なくとも1つの実施形態によれば、本ゲノムベース差分圧縮プログラムは、インプット・ファイルとアウトプット・ファイルの対（ついで）の間の差分ファイル（すなわち、差分）を、一致、置き換え、挿入、および削除など、単純なコラム操作のセットに低減することができる。ソースおよび目的ファイルは、それぞれのコラムとともに全行を通貫して構文解析することが可能である。差分ファイルは、ソース・ファイルを目的ファイルに変換するために、ソース・ファイルに対し行うことが可能な操作のセットの連続的な記述であってよい。派生差分ファイルは、`. z i p`（すなわち、ファイル圧縮および解凍に使用される或るファイル・フォーマットおよびソフトウェア・アプリケーション）など、既存の汎用ファイル・コンプレッサを用いて圧縮することができる。

20

**【 0 0 2 2 】**

少なくとも1つの実施形態によれば、本ゲノムベース差分圧縮プログラムは、医療産業において、ストレージおよびファイル・システム、製品、またはクラウド提供物が、より小さな必要容量によるパイプラインをサポートし作動させ、様々な通信ネットワークを介した伝送を加速化する差別化を可能にすることができる。また、本ゲノムベース差分圧縮プログラムは、様々な中間ステップから、再スタートさせることによって別の分岐の実行を可能にすることができ、仮説および変異のより迅速な確認を可能にすることができる。

30

**【 0 0 2 3 】**

少なくとも1つの実施形態によれば、ゲノム・データ・ファイルは、テキストの行（すなわち、記録）のセットとして表されてよく、行の各々は、（例えば、`S A M`または`V C F`フォーマットのようにタブで区切られた）固定数のコラムまたはフィールドを含む。ソース・ファイルと目的ファイルとを対比することによって、差分を計算することができる。ソースおよび目的ファイルの行の各々は、それらのマッピング位置の昇順にソートすることが可能である。次いで、ソースおよび目的ファイルは行ごとに同調して精査することができ、ソースおよび目的ファイル中の同じ位置の対応する行が対比され、差分ファイルを得ることができる。この派生差分ファイルは、次いで、汎用コンプレッサを使って圧縮されてよい。

40

**【 0 0 2 4 】**

少なくとも1つの実施形態によれば、本ゲノムベース差分圧縮プログラムは、ソースおよび目的ファイルをソートし、リファレンス・ゲノム・ファイルにマップすることが可能であるという事実を考慮に入れることが可能である。しかして、本ゲノムベース差分圧縮プログラムは、ソース・ファイルと目的ファイルとの間の2つの類似の記録に対し、階層識別システムを用いてよい。正確な一致が見出された後、次いで、両方の記録を、個別の

50

フィールドに構文解析し、差分についてフィールド毎ベースで対比することが可能である。

【0025】

少なくとも1つの実施形態によれば、本ゲノムベース差分圧縮プログラムは、同じ位置にマップされた目的ファイル中の行をバッファ保存し、同じ位置のソース・ファイル中の次の行に対して目的ファイルを再スキャンすることによって、行の照合を改良することができる。ゲノム・ファイル中の各行は、リファレンス・ゲノムにマップされた「リード」に対応してよい。一般に、リファレンス・ゲノム中の同じ位置にマップされた多くのリードが存在し得る。事実上、ゲノムの同じ領域にオーバーラップしているリードの平均数をカバレッジ係数とすることができる。カバレッジは、同じ位置にマップされ得る数十、さらには数百のリードと高くなり得る。

10

【0026】

少なくとも1つの実施形態によれば、より小さな所定粒度を有し得るコラム中の差分を見出すために、より細かい粒度（例えば、サブコラム）を用いることができる。ゲノム・フォーマット中の各記録（すなわち、行）の最後のコラムは、キー - 値ペアのセットを含んでよい。パイプラインの諸ステップで、いくつかの新規のキー - 値ペアが追加され得、いくつかは削除され得る。しかして、より小さな粒度を用いることによって、より小さい差分を得ることができる。

【0027】

少なくとも1つの実施形態によれば、差分の符号化（すなわち、圧縮方法）は、ファイルの対の間の差の形でデータを格納することができる。たとえ、配列データの間の差を調べるいくつかの既存の差分圧縮アルゴリズム（例えば、vcdiffおよびxdelta）が存在していても、これらの既存の差分アルゴリズムは、ゲノム・データ・ファイルに適用されたとき、実際の差よりも大きな派生差分ファイルを生成し得る。これら既存の差分アルゴリズムは、ゲノム・データ・ファイルのソートされた順序、および行のタブ区切り構造の利点を利用できていない可能性がある。SAMおよびVCFフォーマットのファイルを含め、これらファイルは、効率的な差分圧縮のためゲノムベース差分圧縮プログラムから取り除かれてよい、ヌクレオチド配列および品質値に加え、過度の補助情報を含み得る。

20

【0028】

少なくとも1つの実施形態によれば、このアルゴリズムの設計は、パイプラインとデータ・ファイルとの主要特徴に基づくことができる。パイプラインの最初のステップは、配列のアラインメントであり、これはSAMファイルをもたらし、このアラインメントでは、あらゆる生の配列リードが、リファレンス・ゲノム中の特定の位置にマップされるか、またはマップされずに遺棄される。アラインメントの位置は、再アラインメントの過程で変更されない限り、残りのパイプラインの間ずっと同じのままに留まり得る。ほとんどのパイプライン・ツールは、アラインメント・ステップの後でだけファイルをソートすることができる。本ゲノムベース差分圧縮プログラムは、ファイルの対の間の差を同調させて計算するため、このソートを利用することが可能である。2つのファイルの間でソート順序が大きく変化すると、差分の大きさもまた増大し得る。本ゲノムベース差分圧縮プログラムは、差分が受容可能な閾値を超えて増加し得るかどうかを検出することができる。かかる場合に、本ゲノムベース差分圧縮プログラムは、差分圧縮を中止し目的ファイルを変更しないでおくことが可能である。

30

40

【0029】

少なくとも1つの実施形態によれば、本ゲノムベース差分圧縮プログラムは、ソースおよび目的ファイルのヘッダを迅速に検索し、第一記録に到達することができる。ソース・ファイル中のあらゆる個別記録に対し、差分ファイル中に対応する処置があり得る。あらゆるソース記録に対し、本ゲノムベース差分圧縮プログラムは、最初に、目的ファイル中に一致を見出す試みをすることができる。この一致を見出すための検索領域は、SAMの場合は、同じ値を持つリファレンス配列名（REFまたはRNAME）および1から始まる左端のマッピング位置（POS: position）に、VCFの場合は、リファレン

50

ス・ゲノムからの染色体識別子またはアセンブリ・ファイル（#CHROM）中のコンテ  
 イグ（例えば、オーバーラップしている配列データまたはリード）と位置（POS）を指  
 している角括弧内の識別文字列を有する記録に限定され得る。本ゲノムベース差分圧縮プ  
 ログラムは、目的ファイルを、一致を見出すことができるまでまたは検索領域が尽きるま  
 で、逐次的に読み取ることが可能である。目的ファイルからの不一致記録は、1つのソー  
 ス記録に対する一致の検索の過程で遭遇することができ、それらは、本ゲノムベース差分  
 圧縮プログラムが、ソース・ファイルの次の記録に進む際に不一致記録を改めて処理す  
 るために、バッファ（例えば、ファイルが一時的に格納されるメモリの小部分）にしまわれ  
 てよい。バッファにしまわれた記録の数が受容可能な閾値を超えた場合、本ゲノムベース  
 差分圧縮プログラムは検索を中止することができる。一致が見出せない場合、当該ソー  
 ス記録は、差分ファイル中に「記録を削除」と標識されてよい。本ゲノムベース差分圧縮プ  
 ログラムが、目的ファイルを逐次に読取っている間に、ソース記録のものよりも小さい可  
 能性のあるREF、POSまたは#CHROM、REFを有する記録に遭遇した場合、目  
 的ファイルからのその特定の記録は、「記録を挿入」と表記されてよく、差分ファイル中  
 に加えられてよい。

10

#### 【0030】

少なくとも1つの実施形態によれば、本ゲノムベース差分圧縮プログラムは、特定の請  
 求項に基づき、ソース・ファイルと目的ファイルとを照合するため階層識別システムを利用  
 することができる。例えば、SAMファイルでは、本ゲノムベース差分圧縮プログラム  
 は、クエリ・テンプレート名（QNAME）、ビット単位フラグ（FLAG）、REF、  
 POS（すなわち、第一～第四コラム）に関連付けられたコラムを利用することが可能で  
 、VCFファイルでは、本ゲノムベース差分圧縮プログラムは、#CHROM、POS、  
 REF、および代替塩基（群）（ALT）（すなわち、第一～第二コラム、および第四～  
 第五コラム）を利用することが可能である。

20

#### 【0031】

少なくとも1つの実施形態によれば、ソース記録と目的記録との間の一致が判断された  
 後、特定のフィールド（例えば、SAMファイルでは、任意もしくは非必須のフィールド  
 （OPTIONAL）または第十二コラム、VCFファイルでは、付加情報（INFO）  
 または第八コラム）に対してフィールド毎の対比を行うことができる。これら対比された  
 フィールドに対し、フィールド中に不一致があった場合、その特定のフィールドは「置き  
 換え」と標識されてよく、差分ファイルにその新規の値を加えることができる。対比から  
 除外されたファイル（例えば、OPTIONALおよびINFO）は、いくつかのキー・  
 値ペア、またはキー単体および値単体を含むことが可能で、これらは、セミコロン、空白  
 スペース、または別のマークで区切られてよい。対比から除外されたこれらのフィールド  
 は、次いで、目的およびソース・ファイルの両方に対して、それらそれぞれのサブフィー  
 ルドに構文解析することができる。ソースに存在し、目的に不在のキーまたは値は、差分  
 ファイル中に「削除」と標識されてよい。目的に存在し、ソースに不在のキーまたは値は  
 、「挿入」と標識されてよく、差分ファイル中に加えられてよい。ソース記録に対するこ  
 れら操作が決められ差分ファイルに加えられた後、本ゲノムベース差分圧縮プログラムは  
 、ソース・ファイル中の次の記録の読み取りに進むことができ、ソースおよび目的に対し  
 、ファイルの終了点（EOF：end-of-file）に達するまで、このプロセスを  
 繰り返すことができる。生成された差分ファイルは、高度に圧縮可能であり得、アーカイ  
 ブのため圧縮されてよい。

30

40

#### 【0032】

少なくとも1つの実施形態によれば、あらゆるソース記録に対し、本ゲノムベース差分  
 圧縮プログラムは、目的ファイルから新しい記録を読み取る前に、まずバッファ中の一致  
 を検索することができる。バッファは、検索領域が変更される都度、毎回クリアするこ  
 とが可能である。全ての不一致記録は、「挿入」とタグされ、差分に加えられてよい。

#### 【0033】

図1を参照すると、少なくとも1つの実施形態による、典型的なネットワーク化コンピ

50

ユーザ環境 100 が描かれている。このネットワーク化コンピュータ環境 100 は、プロセッサ 104 と、ソフトウェア・プログラム 108 の実行が可能にされたデータ・ストレージ・デバイス 106 と、ゲノムベース差分圧縮プログラム 110 a とを備えたコンピュータ 102 を含んでよい。また、ネットワーク化コンピュータ環境 100 は、データベース 114 および通信ネットワーク 116 と相互作用することが可能なゲノムベース差分圧縮プログラム 110 b の実行を可能にされたサーバ 112 も含んでよい。ネットワーク化コンピュータ環境 100 は、複数のコンピュータ 102 およびサーバ 112 を含むことができ、そのうちの 1 つのみが示されている。通信ネットワーク 116 は、広域ネットワーク (WAN)、ローカル・エリア・ネットワーク (LAN)、電気通信ネットワーク、ワイヤレス・ネットワーク、公衆交換ネットワーク、もしくは衛星ネットワーク、またはこれらの組合せなど、様々な種類の通信ネットワークを含んでよい。当然のことながら、図 1 は、1 つの実装の単なる例を提示しており、各種の実施形態が実装可能な環境に関し、いかなる限定をも意味するものではない。設計および実装上の要件に基づいて、描かれた環境に多くの修改を加えることが可能である。

#### 【0034】

クライアント・コンピュータ 102 は、通信ネットワーク 116 を介してサーバ・コンピュータ 112 と通信することが可能である。通信ネットワーク 116 は、有線、ワイヤレス通信リンク、または光ファイバ・ケーブルなどの接続手段を含んでよい。図 9 を参照しながら説明するように、サーバ・コンピュータ 112 は、内部コンポーネント 902 a および外部コンポーネント 904 a をそれぞれ含んでよく、クライアント・コンピュータ 102 は、内部コンポーネント 902 b および外部コンポーネント 904 b をそれぞれ含んでよい。また、サーバ・コンピュータ 112 は、サービスとしてのソフトウェア (SaaS: Software as a Service)、サービスとしてのプラットフォーム (PaaS: Platform as a Service)、またはサービスとしてのインフラストラクチャ (IaaS: Infrastructure as a Service) など、クラウド・コンピューティング・サービス・モデル中で運用することも可能である。また、サーバ 112 は、プライベート・クラウド、コミュニティ・クラウド、パブリック・クラウド、またはハイブリッド・クラウドなど、クラウド・コンピューティング展開モデル中に配置されてもよい。クライアント・コンピュータ 102 は、例えば、携帯デバイス、電話、携帯情報端末、ネットブック、ラップトップ・コンピュータ、タブレット・コンピュータ、デスクトップ・コンピュータ、または、プログラムを実行し、ネットワークにアクセスし、データベース 114 にアクセスする能力のある任意の種類のコンピューティング・デバイスであってよい。本実施形態の様々な実装によれば、ゲノムベース差分圧縮プログラム 110 a、110 b は、以下に限らないが、コンピュータ/携帯デバイス 102、ネットワーク・サーバ 112、またはクラウド・ストレージ・サービスなどの様々なストレージ・デバイス中に内蔵することが可能なデータベース 114 とやり取りすることができる。

#### 【0035】

本実施形態によれば、クライアント・コンピュータ 102 またはサーバ・コンピュータ 112 を用いるユーザは、ゲノム・データ・ファイルに対する差分ファイルを圧縮するためにゲノムベース差分圧縮プログラム 110 a、110 b を (それぞれ) 使用することが可能である。ゲノムベース差分圧縮の方法については、以降で図 2 ~ 8 に関連してさらに詳しく説明する。

#### 【0036】

ここで図 2 を参照すると、少なくとも 1 つの実施形態による、ゲノムベース差分圧縮プログラム 110 a および 110 b によって用いられる、ゲノム・データ・ファイルに対する典型的な圧縮プロセス 200 を示すオペレーション・フローチャートが描かれている。

#### 【0037】

202 で、ゲノム・データ・ファイルが、ゲノムベース差分圧縮プログラム 110 a、110 b へのインプットとして受信される。ゲノム・データ・ファイルは、ユーザのデバ

10

20

30

40

50

イス（例えば、ユーザのコンピュータ102）上のソフトウェア・プログラム108を用い、通信ネットワーク116を介してゲノムベース差分圧縮プログラム110a、110b中に伝送することができる。ゲノム・データは、或る生体に関連するゲノム・データを含んでよく、このゲノム・データのファイルは、デオキシリボ核酸（DNA：deoxyribonucleic acid）配列のリード（すなわち、DNA分子内のヌクレオチドの正確な順序）、リファレンス・ゲノムへのマッピング、およびリファレンス・ゲノムに対して検出された変異を含んでよい。また、ゲノム・データ・ファイルは、少なくとも2つのファイル、ソース（すなわち、リファレンス）および目的ファイルに対し、各々が固定数のコラム（すなわち、フィールド）を含む行（すなわち、記録）のセットを含んでよい。ソース・ファイルは、生体のリファレンス・ゲノムから直接導出することが可能なデータである。各ファイルは、そのファイルが目的ファイルまたはソース・ファイルのどちらかによって標識されてよい（例えば、「target\_\_name.format.type（目的\_\_名前.フォーマット・タイプ）」または「source\_\_name.format.type（ソース\_\_名前.フォーマット・タイプ）」）。ゲノムベース差分圧縮プログラム110a、110bは、目的ファイルを再構築するためにソース・ファイルを利用することが可能である。

10

#### 【0038】

例えばゲノムベース差分圧縮プログラム110a、110bは、ゲノム解析ツールキット（GATK：genome analysis toolkit）から伝送されてくるソース・ファイル（すなわち、source\_\_liverHS1567.sam）および目的ファイル（すなわち、target\_\_liverHS1567.sam）を受信する。ソース・ファイルは、或る健康な人の肝細胞標本に関連付けられたリファレンス・ゲノムからのデータを含み、目的ファイルは、肝がんを患っている人からの肝細胞標本に関連付けられたデータを含む。ゲノムベース差分圧縮プログラム110a、110bが、302で、インプットとしてソースおよび目的ファイルを受信するステップを含め、ゲノム・データ・ファイルに対する典型的な圧縮プロセス300については、図3に関連して後記でもっと詳しく説明することとする。

20

#### 【0039】

本実施形態において、ゲノム・データ・ファイルは、GATKなどのゲノム解析パイプライン実行からのアウトプットとして生成されてよい。パイプラインのスク립トは、目的およびソース・ファイルにゲノムベース差分圧縮プログラム110a、110bを実行するための特定のコマンドを使って充実させることが可能である。さらに、圧縮されたアウトプット・ファイル（例えば、SAM（BAM）およびVCF（BCF）ファイルのバイナリ・バージョン）の一部は、ゲノムベース差分圧縮プログラム110a、110bを開始する前に、外部のエンジンを利用して解凍することができる。

30

#### 【0040】

本実施形態において、ソースおよび目的ファイルは、アウトプット・ファイルが生成され格納される場所であるフォルダまたはディレクトリ（すなわち、データベース114）中の、パイプラインの実行中に生成された諸ファイルから読み出すことが可能である。

#### 【0041】

本実施形態において、ゲノムベース差分圧縮プログラム110a、110b中にインプットとして1つのゲノム・データ・ファイルしか投入されなかった場合、ゲノムベース差分圧縮プログラム110a、110bは、派生差分ファイルを計算することはできない。代わりに、ゲノムベース差分圧縮プログラム110a、110bは、ユーザにエラー・メッセージを返すことができる。ユーザは、同じインプット・ファイルを再投入することが可能であるが、但し、そのファイルは、ゲノムベース差分圧縮プログラム110a、110bが派生差分ファイルを計算するための、対応するソースまたは目的ファイルと共に投入される必要がある。

40

#### 【0042】

次に204で、ゲノムベース差分圧縮プログラム110a、110bによって、ゲノム

50

・データ・ファイルの行がソートされる。ソースおよび目的ファイルの行は、リファレンス・ゲノム中のそれらのマッピング位置（POS）に基づいて昇順に配置することができ、これにより、ソース・ファイルと目的ファイルとの間での行の照合のための検索領域を絞りこみ、バッファ占拠域を最小化し、ゲノムベース差分圧縮プログラム 110 a、110 b の実効時間を最適化することが可能になる。ゲノム・データ・ファイルの行に対する行群をソートするために、既存のソート・アルゴリズムを含む既存のオープン・ソース・ソート・ツールを利用することができる。

【0043】

目的およびソース・ファイルのフォーマットの型によっては、これらファイルは、パイプラインにおいて既にソートされていることがある。ソースおよび目的ファイルがパイプラインにおいてソート済みの場合、それら目的およびソース・ファイルは、202で、インプットとしてゲノムベース差分圧縮プログラム 110 a、110 b 中に投入されるとき、それらそれぞれの行のソートされた順序を保持し得る。

10

【0044】

前述の例を続けると、ゲノムベース差分圧縮プログラム 110 a、110 b は、GATK から得られたソースおよび目的ファイルが既にソートされており、したがって、ゲノムベース差分圧縮プログラム 110 a、110 b がマッピング位置（POS）に基づいて行をソートする必要はない、と判断する。

【0045】

本例の目的の上で、ゲノムベース差分圧縮プログラム 110 a、110 b は、SAM または VCF フォーマットのファイルにだけ対応可能である。ファイルがなんらかの他のフォーマットである場合、ユーザには、エラー・メッセージが提示されることになり、ゲノムベース差分圧縮プログラム 110 a、110 b は、対応可能なソースおよび目的ファイルがインプットとして投入されるまで、実行を停止する。本例では、ゲノムベース差分圧縮プログラム 110 a、110 b は、ソースおよび目的ファイルは SAM フォーマットであると判断し、これは、ゲノムベース差分圧縮プログラム 110 a、110 b が対応可能である。しかして、ゲノムベース差分圧縮プログラム 110 a、110 b はバッファ（すなわち、目的ファイルがソース・ファイルに一致しなかった場合に、その目的ファイルを格納できる、メモリの小部分）を起動する。ゲノムベース差分圧縮プログラム 110 a、110 b がフォーマットを識別するためファイルをチェックし、バッファを初期化する（INIT:initializing）ステップを含め、ゲノム・データ・ファイルに対する差分ファイルの圧縮のための典型的なプロセスについては、図3と関連して後記でもっと詳しく説明することとする。

20

30

【0046】

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、インプットとしてファイルを受信する前にソースおよび目的ファイルがソートされていなかった場合、ゲノムベース差分圧縮プログラム 110 a、110 b は、それらのPOS（すなわち、各ファイルの行の第四コラム）に基づいてオープン・ソース・ソート・ツールを利用することによって、ソースおよび目的ファイルの各々の行を昇順でソートすることができる。

【0047】

40

本実施形態において、ゲノムベース差分圧縮プログラム 110 a、110 b は、ファイルのフォーマット型（例えば、SAM、BAM、VCF、BCF）の対応可能性をチェックすることが可能である。ファイルのフォーマット型が、ゲノムベース差分圧縮プログラム 110 a、110 b によって対応可能である場合、ゲノムベース差分圧縮プログラム 110 a、110 b は先に進んでよい。また一方、ファイルのフォーマット型が、ゲノムベース差分圧縮プログラム 110 a、110 b にとって対応可能でない場合、ゲノムベース差分圧縮プログラム 110 a、110 b は、実行を停止してよく、ユーザにエラー・メッセージを提示することができる。

【0048】

本実施形態において、ゲノム・データ・ファイルの行群を通貫してソートするために用

50

いることが可能なソート・アルゴリズムは、ゲノム・データ・ファイルの複雑さの如何によつては、部分的に変更をする必要があるかもしれない。

【0049】

次いで206で、ソースおよび目的ファイルが精査される。ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルを逐次的に精査する（すなわち、読み取る）ことが可能で、このとき、ゲノム・データ・ファイル中の各ラインを一行として判断してよい。インプット・ファイルは、改行コード（すなわち、ラインの終了と新ラインの開始を表す特殊な文字または文字列）に遭遇するまで、ゲノムベース差分圧縮プログラム110a、110bによって読み取られてよい。さらに、これらの行は、コラム（すなわち、フィールド）に分割することができる。2つのファイルが同調されていることを確認するために、これら行の値が対比されてよい。両方のファイルは、これらファイルのコンテンツの間の差を、時間およびメモリ効率的な仕方で計算するために同調させて精査することができる。

10

【0050】

前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルの各行を精査する。ソースおよび目的ファイルは、各行に3つのデータ片を備える4つの行を含む。2つのゲノム・データ・ファイルを対比するブロック図について、図6に関連して後記でもっと詳しく説明することとする。ソース・ファイルの各行を精査した後、ゲノムベース差分圧縮プログラム110a、110bは、バッファが空かどうかを判断する（すなわち、バッファが空の場合、目的ファイルは、バッファの代わりにファイルから読み取ることが可能である）。本例では、バッファは空である。したがって、ゲノムベース差分圧縮プログラム110a、110bは、ファイルからの各目的行を精査すればよい。ゲノムベース差分圧縮プログラム110a、110bがソースおよび目的ファイルの各行を精査した後、ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルが同調されていることを確認する。ゲノムベース差分圧縮プログラム110a、110bがソースおよび目的ファイルの行を読み取るステップを含め、ゲノム・データ・ファイルに対する差分ファイルを圧縮するための典型的なプロセスについて、図3に関連して後記でもっと詳しく説明することとする。

20

【0051】

次いで208で、目的ファイルとソース・ファイルとの行が対比される。ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルと目的ファイルとの間の一切の差を判断するために、ソースおよび目的ファイル上の対応する行（すなわち、同じマッピング位置を有する行）を対比し、派生差分ファイル（すなわち、派生Dファイル）を生成することができる。

30

【0052】

ゲノムベース差分圧縮プログラム110a、110bは、部分一致している行の階層的フィールド単位の対比を（例えば、各ファイルに対し、特定のコラムが対比される）、および一致している行に対しては行の全面的対比を行う（例えば、両方のファイル中のさらなるコラムが対比される）ことが可能である。ゲノムベース差分圧縮プログラム110a、110bを用いて、対比される行が一致しているかどうかによりどちらの対比を使えばよいかを判断することができる。目的およびソース・ファイルからのそれぞれの行の対比は、特定のコラムに対するフィールド単位の対比によって開始されてよい。階層的フィールド単位の対比の過程での各対比によって、中間的な一致または不一致を記録することができる。階層的フィールド単位の対比の終了点で、行の間的一致または不一致のいずれかが生じ得る。この一致または不一致は、対比された行が同調しているかしていないかを判断し、さらに、不一致があるかどうかを判断するための手っ取り早い方法であり得る。これら2つの行の間に一致がある場合、行全体の対比（すなわち、全コラムが対比される）を行うことができる。この対比の結果は、派生差分ファイルに取り込むことができ、このファイルはゲノムベース差分圧縮プログラム110a、110bのアウトプットである。

40

【0053】

50

前述の例を続けると、ゲノムベース差分圧縮プログラム 110 a、110 bは、最初に、目的およびソース・ファイル内のそれぞれの行の階層的フィールド単位の対比を行う。このソースおよび目的ファイルはSAMフォーマットなので、11の固定コラムがある。しかしながら、階層的フィールド単位の対比は、4つのコラム（すなわち、POS、REF、QNAME、およびFLAG）だけを対比する。ゲノムベース差分圧縮プログラム 110 a、110 bは、目的およびソース・ファイル中の第一行のPOSおよびREFコラムを最初に対比し、それらPOSおよびREFコラムが一致するかどうかを判断する。目的およびソース・ファイル中の第一行のPOSおよびREFコラムが一致したので、ゲノムベース差分圧縮プログラム 110 a、110 bは中間的な一致があると判断し、次いで、同じ行のQNAMEコラムが対比される。ゲノムベース差分圧縮プログラム 110 a、110 bは、同じ行のQNAMEコラムが一致しており、中間的な一致があると判断する。次いで、ゲノムベース差分圧縮プログラム 110 a、110 bは、同じ行のFLAG（ビット単位）コラムを対比する。これらの行のFLAG（ビット単位）コラムが一致しないので、ゲノムベース差分圧縮プログラム 110 a、110 bは、これらの行のFLAG（ビット単位）の間の差が $0 \times 400$ （すなわち、光学的デュプリケート距離）以下かどうかを判断する。ゲノムベース差分圧縮プログラム 110 a、110 bは、この差が $0 \times 400$ に等しかったと判断する。したがって、目的ファイルとソース・ファイルとの第一行は一致する。ソース・ファイルと目的ファイルとの対比によって生成されたDファイルの第一行に「m」が含まれる。次いで、ゲノムベース差分圧縮プログラム 110 a、110 bは、ソースおよび目的ファイル中の次の行の階層的フィールド単位の対比に進んでよい。SAMフォーマットのゲノム・データ・ファイルの階層的構造を識別するための典型的なプロセスについて、図7に関連して後記でもっと詳しく説明することとする。

#### 【0054】

部分一致の行のフィールド単位の対比が行われたとき、以下の結果が生成され、派生Dファイルを生成するために用いられた。

`field`（フィールド）`_i`（S）が`field__i`（T）と異なる場合：D（`field__i`（T），i）

#### 【0055】

階層的フィールド単位の対比で対比された行が一致する場合、ゲノムベース差分圧縮プログラム 110 a、110 bは、行全体の対比を行うことになり、この対比は、目的およびソース・ファイルの行中の、RNAME（すなわち、第三コラムであり、アラインメントの名前を表す）と、MAPQ（すなわち、第五コラムであり、品質をマッピングしている）と、CIGAR（すなわち、第六コラムであり、配列が一致しているかどうかに関連する文字列演算を表す）と、RNEXT（すなわち、第七コラムであり、次のリードのリファレンス名を表す）と、PNEXT（すなわち、第八コラムであり、次のリードの位置を表す）と、TLEN（すなわち、第九コラムであり、観測されたテンプレートの長さを表す）と、SEQ（すなわち、第十コラムであり、部分配列を表す）と、QUAL（すなわち、第十一コラムであり、塩基のエラー確率によって塩基の品質を表す）とを含む、OPTIONALコラム（すなわち、残りの任意のまたは非必須の7つのコラム、または第十二コラム）の対比、およびソース・ファイルと目的ファイルとの行の間の差を計算するためにOPTIONALコラムのサブフィールド全体を通貫して構文解析することを含む。次いで、これらコラムの差が派生ファイルに加えられ、ゲノムベース差分圧縮プログラム 110 a、110 bは、ソース・ファイルから次の行を読み取る。SAMフォーマットのファイルに対する階層的識別のためのプロセスを含むゲノム・データ・ファイルの典型的な圧縮プロセス 300について、図4に関連して後記でもっと詳しく説明することとする。

#### 【0056】

行が不一致なので、行の全面的対比は行わなくてもよい。また一方、行の全面的対比が行われた場合は、以下の結果を生成することができる。

一致（D「m」）、削除（D「d」）、挿入（D T（i））

#### 【0057】

10

20

30

40

50

また一方、ソースおよび目的ファイルがVCFフォーマットであった場合は、各行内に8つの固定で必須の列がある。この階層的フィールド単位の対比では、4つの列だけ(すなわち、#CHROM、POS、REF、およびALT)を対比する。ゲノムベース差分圧縮プログラム110a、110bは、最初に目的およびソース・ファイル中の第一行の#CHROMおよびPOS列を対比し、それら#CHROMおよびPOS列が一致するかどうかを判断する。ゲノムベース差分圧縮プログラム110a、110bが#CHROMおよびPOS列が中間的な一致または不一致のどちらであると判断したかによって、同じ行のREFおよびALT列が対比されることになる。VCFフォーマットのゲノム・データ・ファイルの階層構造を識別するための典型的なプロセスについて、図8に関連して後記でもっと詳しく説明することとする。

10

#### 【0058】

階層的フィールド単位の対比で対比された行が一致すれば、ゲノムベース差分圧縮プログラム110a、110bは、行の全面的対比を行うことができ、この対比は、目的およびソース・ファイルの当該行中のINFO列(すなわち、付加情報を表す第八列)の対比と、ソース・ファイルと目的ファイルとの行の間の差を計算して派生ファイルを作成するためにINFO列のサブフィールドを通貫して構文解析することを含む。次いで、ゲノムベース差分圧縮プログラム110a、110bはソース・ファイルから次の行を読み取る。階層的識別のためのプロセスを含め、VCFフォーマットのファイルのゲノム・データ・ファイルに対する典型的な圧縮プロセス300について、図5に関連して後記でもっと詳しく説明することとする。

20

#### 【0059】

次いで210で、派生差分ファイルは、既存の汎用ファイル・コンプレッサを使って圧縮される。派生差分ファイル(すなわち、派生Dファイル)の汎用ファイル・コンプレッサへの伝送は、派生Dファイルの一時的変更に対するコードを使って自動化することができる。派生Dファイルの圧縮が完了した後、ゲノムベース差分圧縮プログラム110a、110bは、ユーザのデバイス(例えば、ユーザのコンピュータ102)のメモリに保存されてよい。ゲノムベース差分圧縮プログラム110a、110bは、ユーザに、派生Dファイルの圧縮が完了したことを(例えば、ダイアログ・ボックスを介して)プロンプトすることができる。このダイアログ・ボックスには、例えば、圧縮派生Dファイルが完了したことのメッセージ、およびダイアログ・ボックスの下方に「詳細を見る」ボタンを含めてよい。ユーザが「詳細を見る」ボタンをクリックすると、このダイアログ・ボックスは消え、圧縮派生Dファイルを含む別のダイアログ・ボックスをユーザに提示することができる。この派生Dファイルは、目的ファイルを再生するために、ソース・ファイルと併せ用いることが可能である。

30

#### 【0060】

前述の例を続けると、ソース・ファイルと目的ファイルとの対比によって生成された(すなわち、ソース・ファイルが、目的ファイルを派生Dファイルとして再構成するために用いられた)派生Dファイルは、.zipによって圧縮される。派生差分ファイルは、当初のサイズの52%に低減されており、圧縮派生Dファイルは、ユーザのコンピュータ(すなわち、ユーザ・デバイス102)のメモリに保存される。ユーザは、ダイアログ・ボックスによって、圧縮派生Dファイルが完了したことをプロンプトされ、ユーザはメッセージの下方の「詳細を見る」ボタンをクリックする。ユーザが「詳細を見る」ボタンをクリックすると、このダイアログ・ボックスは消え、圧縮派生Dファイルをユーザに提示する別のダイアログ・ボックスが現れる。

40

#### 【0061】

ここで図3を参照すると、少なくとも1つの実施形態による、ゲノムベース差分圧縮プログラム110aおよび110bによって用いられるゲノム・データ・ファイルに対する、典型的な圧縮プロセス300を示すオペレーション・フローチャートが描かれている。図示のように、ゲノムベース差分圧縮プログラム110a、110bは、Dファイルを作成するためにソースおよび目的ファイルを用い、派生Dファイルを圧縮するために汎用フ

50

ファイル・コンプレッサを用いる。

【0062】

302で、ソースおよび目的ファイルが、インプットとしてゲノムベース差分圧縮プログラム110a、110b中に投入される。ソースおよび目的ファイルは、ユーザのデバイス（例えば、ユーザのコンピュータ102）上のソフトウェア・プログラム108を使い、通信ネットワーク116を介して、ゲノムベース差分圧縮プログラム110a、110b中に、インプットとして伝送することが可能である。各ファイルは、そのファイルがソース・ファイルかまたは目的ファイルなのかを示すために標識されてよい。しかして、ゲノムベース差分圧縮プログラム110a、110bが、伝送されてきたファイルの種類（すなわち、そのファイルがソース・ファイルかまたは目的ファイルのいずれか）を識別

10

【0063】

例えば、ゲノムベース差分圧縮プログラム110a、110bは、ゲノム解析ツールキット（GATK）から、ソース・ファイルおよび目的ファイルを受信する。このソース・ファイルは、健全なヒトスジシマ蚊に関連付けられたリファレンス・ゲノムからのデータを含み、目的ファイルは、西ナイル・ウイルス（すなわち、フラビウイルス）に感染したヒトスジシマ蚊のゲノムに関連付けられたデータを含む。

【0064】

次に304で、ゲノムベース差分圧縮プログラム110a、110bは、フォーマットを識別するためにファイルをチェックし、バッファ（INITバッファ）を初期化する。ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルが対応可能なフォーマット（例えば、SAMまたはVCF）であることをチェックすることができる。次いで、ゲノムベース差分圧縮プログラム110a、110bは、バッファを初期化することができ、該バッファは、ゲノムベース差分圧縮プログラム110a、110bによってソース・ファイルから一致する行が見出されるまで、目的ファイルからの行を格納することが可能なメモリの小部分であってよい。前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、受信されたソースおよび目的ファイルがVCFフォーマットであることを確認した。加えて、ゲノムベース差分圧縮プログラム110a、110bは、バッファを起動する。

20

【0065】

次いで306で、ソース・ファイルから一行が読み取られる。フォーマットの型、およびファイルがパイプラインから読み出されているかどうかによって、それぞれのファイルの行が既にソートされている可能性があり、その場合、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルの精査（または読み取り）に進んでよい。前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルが、ゲノムベース差分圧縮プログラム110a、110bに伝送される前にパイプラインにおいて既にソートされていたと判断する。しかして、受信されたソースおよび目的ファイルは、ゲノムベース差分圧縮プログラム110a、110bが受信ソース・ファイルの行を読み取る前に、ソートを行う必要はない。しかして、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルの第一行の読み取りに進むこと

30

40

【0066】

308で、ゲノムベース差分圧縮プログラム110a、110bが、バッファが空である（例えば、バッファ内に目的ファイルがない）と判断した場合、312で、ゲノムベース差分圧縮プログラム110a、110bはファイルから目的行を読み取ってよい。ゲノムベース差分圧縮プログラム110a、110bは、バッファを検索し、バッファ中に何らかの目的行があるかどうかを判断する。バッファ中にいかなる目的行もない場合、ゲノムベース差分圧縮プログラム110a、110bは、バッファが空であると判断する。

【0067】

また一方、308で、バッファが空でない（例えば、バッファが少なくとも1つの目的

50

行を含む)場合、310で、その目的行がバッファから読み取られる。バッファ中に目的行が存在する場合、ゲノムベース差分圧縮プログラム110a、110bは、バッファが空でないと判断する。

**【0068】**

前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、バッファが空でなく3つの目的行を含むと判断する。しかして、それら目的行の1つがバッファから読み取られる。

**【0069】**

ゲノムベース差分圧縮プログラム110a、110bが、310でバッファから目的行を読み取るか、または312でファイルから目的行を読み取ったとき、314で、ゲノムベース差分圧縮プログラム110a、110bは、当該ソースおよび目的ファイルに対するフォーマットがSAMかまたはVCFかを判断する。ゲノムベース差分圧縮プログラム110a、110bがソースおよび目的ファイルの両方を読み取った後、ゲノムベース差分圧縮プログラム110a、110bは、派生Dファイルの生成および圧縮に進み、これらはフォーマットの型によって異なる。314で、ゲノムベース差分圧縮プログラム110a、110bが、ファイルはSAMフォーマットであると判断した場合、ゲノムベース差分圧縮プログラム110a、110bは図4に進む。また一方、314で、ゲノムベース差分圧縮プログラム110a、110bが、ファイルはVCFフォーマットであると判断した場合、ゲノムベース差分圧縮プログラム110a、110bは図5に進む。

**【0070】**

前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルの両方がVCFフォーマットであることを確認し、次いでゲノムベース差分圧縮プログラム110a、110bは図5に進む。

**【0071】**

ここで図4を参照すると、少なくとも1つの実施形態による、ゲノムベース差分圧縮プログラム110aおよび110bによって用いられる、SAMフォーマットのゲノム・データ・ファイルに対する典型的な圧縮プロセス300を示すオペレーション・フローチャートが描かれている。SAMフォーマットのソース・ファイルと目的ファイルとの行内の諸フィールドの対比が示されている。図4では、ソースおよび目的ファイルはSAMフォーマットである。

**【0072】**

ゲノムベース差分圧縮プログラム110a、110bが、316で、POSおよびREFは同じでないと判断した場合、318で、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルのPOSおよびREFが目的ファイルのものより小さいかどうかを判断する。POSおよびREFが同じかどうかを判断するために、ゲノムベース差分圧縮プログラム110a、110bによって、SAMフォーマットのソースおよび目的ファイルのそれぞれの行を評価することができる。

**【0073】**

前述の例を続けると、ソースおよび目的ファイルがSAMフォーマットであった場合、ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルの第一行を評価し、それぞれのファイルのPOSおよびREFが同じかどうかを判断することが可能である。ソースおよび目的ファイルが同じである場合、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルと目的ファイルとのQNAMEが同じであるかどうかを評価してよい。また一方、ソース・ファイルと目的ファイルとのPOSおよびREFが同じでない場合、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルのPOSおよびREFが、目的ファイルより小さいかどうかを判断することができる。

**【0074】**

ゲノムベース差分圧縮プログラム110a、110bが、318で、ソース・ファイルのPOSおよびREFは、目的ファイルより大きいと判断した場合、320で、ゲノムベ

10

20

30

40

50

ソース差分圧縮プログラム 110 a、110 b は、目的の記録を差分に挿入し、306 のソースから行を読み取る、に戻る。したがって、ソース・ファイル中の行の P O S および R E F が目的ファイル中のものより大きい場合、その目的の記録を派生差分ファイルに含めることができる。前述の例を続けると、ソース・ファイルが目的ファイルよりも大きい場合、その目的の記録は差分中に挿入されてよい。したがって、派生差分ファイルの当該行中に、その目的の記録を挿入することが可能である。

**【0075】**

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、318 で、ソース・ファイルの P O S および R E F が目的ファイル中の対応する行よりも小さいと判断した場合、322 で、差分中の対応する行は削除され、ゲノムベース差分圧縮プログラム 110 a、110 b は、306 のソース・ファイルから次の行を読み取る、に戻る。前述の例を続けると、ソース・ファイルが目的ファイルよりも小さい場合、そのソースの記録は、派生差分から除去されてよい。しかして、その差分は、派生差分ファイルの当該行に、ソースの記録の代わりに「d」を有してよい。

10

**【0076】**

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、316 で、ソース・ファイル中の行の P O S および R E F は目的ファイル中の対応する行と同じであると判断した場合、324 で、その行は、それぞれの Q N A M E が同じかどうかを判断するため評価される。前述の例を続けると、受信されたソースおよび目的ファイルの行が同じ P O S および R E F を有する場合、ゲノムベース差分圧縮プログラム 110 a、110 b は、受信されたソース・ファイルと目的ファイルとの対応する行の Q N A M E が同じかどうかを判断することができる。

20

**【0077】**

ゲノムベース差分圧縮プログラム 110 a、110 b が、324 で、ソースおよび目的ファイルの行の Q N A M E は同じでないと判断した場合、326 で、目的ファイルからの当該行はバッファに保存され、ゲノムベース差分圧縮プログラム 110 a、110 b は、312 のファイルから次の目的行を読み取る、に戻る。前述の例を続けると、目的およびソース・ファイルの対応する行の Q N A M E が同じでない（すなわち、異なる）場合、その目的行は、該目的行に対するより良好な一致が判断されるまで、バッファ中に保存される。

30

**【0078】**

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、324 で、それぞれの行の Q N A M E が同じであると判断した場合、328 で、目的およびソース・ファイルの当該行は、F L A G が同じであるかどうかを判断するために評価される。前述の例を続けると、ソースおよび目的ファイルの対応する行の Q N A M E が同じである場合、ゲノムベース差分圧縮プログラム 110 a、110 b は、受信されたソースおよび目的ファイルの対応する行の F L A G（ビット単位）を対比する。

**【0079】**

ゲノムベース差分圧縮プログラム 110 a、110 b が、328 で、それぞれの行の F L A G が同じでないと判断した場合、330 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、その差が閾値（例えば、0 x 4 0 0）以下かどうかを判断する。330 で、差が閾値（例えば、0 x 4 0 0）より大きい場合、326 で、目的ファイルからの行がバッファに保存され、ゲノムベース差分圧縮プログラム 110 a、110 b は、ファイルから目的行を読み取るため 312 に戻る。前述の例を続けると、受信されたソースおよび目的ファイルの対応する行のそれぞれの F L A G（ビット単位）が相異なる場合、ゲノムベース差分圧縮プログラム 110 a、110 b は、その差が閾値以下かどうかを判断することができる。この閾値は、ユーザが前もって決めておくことが可能で、本例では、0 x 4 0 0（すなわち、光学的デュープリケート距離）の閾値がユーザによって事前に選択されていた。しかして、ゲノムベース差分圧縮プログラム 110 a、110 b は、F L A G（ビット単位）の差が 0 x 4 0 0 以下かどうかを判断すればよい。

40

50

## 【 0 0 8 0 】

また一方、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b が、3 3 0 で、この差は閾値（例えば、 $0 \times 4 0 0$ ）以下であると判断するか、または、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b が、3 2 8 で、目的およびソース・ファイルの行の F L A G が同じであると判断した場合、3 3 2 で、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、O P T I O N A L フィールドを除き、コラムの差を計算する。ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、コラムの差を調べ、次いでどのコラムがかかる差を含むかを示すことができる。差のあるコラムの位置およびコラムの差の程度によって、派生差分ファイル内の記録が影響され得る。前述の例を続けると、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b が、受信されたソースおよび目的ファイルの対応する行の間の差が  $0 \times 4 0 0$  以下であると判断するか、または受信されたソースおよび目的ファイルの対応する行の F L A G（ビット単位）が同じであると判断した場合、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、O P T I O N A L フィールドを除いたコラムの差の計算に進んでよい。

10

## 【 0 0 8 1 】

次いで 3 3 4 で、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、O P T I O N A L フィールドをサブフィールドに構文解析する。O P T I O N A L のサブフィールドは、S A M フォーマット・ファイルのアラインメント記録中に含めることが可能で、これらには、ソース・ファイルと目的ファイルとの対応する行の O P T I O N A L のサブフィールドの間の相似点または差を判断するため、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b によって構文解析することが可能な、2 0 より多い所定の標準タグを含めることが可能である。前述の例を続けると、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、受信されたソースおよび目的ファイルの対応する行の間に存在する一切の相似点および差があるかどうかを判断するために、O P T I O N A L フィールドを通貫して、O P T I O N A L フィールド内のサブフィールド中の所定の標準タグを検索することができる。

20

## 【 0 0 8 2 】

次いで 3 3 6 で、目的ファイルとソース・ファイルとのそれぞれの行のこれらサブフィールドの差が計算される。このとき、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、目的およびソース・ファイルのそれぞれの行の O P T I O N A L のサブフィールドを通貫して差を調べることができ、目的およびソース・ファイルそれぞれの行に対し、どのサブフィールドに差があるかを示すことができる。この差の程度によって、派生差分ファイル内の記録が影響され得る。前述の例を続けると、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、受信されたソース・ファイルと目的ファイルとの対応する行の O P T I O N A L サブフィールドの差を計算することができる。

30

## 【 0 0 8 3 】

次いで 3 3 8 で、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、コラムの差を差分ファイルに加え、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、ソース・ファイルから次の行を読み取るため 3 0 6 に戻る。ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、最短ではソースおよび目的ファイルが終わるまで、ゲノム・データ・ファイルに対するこの典型的な圧縮プロセス 3 0 0 を通してのループを続ける。前述の例を続けると、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b は、コラムの差を派生差分ファイルに加え、受信されたソース・ファイルの次の対応行を読み取るために 3 0 6 に戻る。

40

## 【 0 0 8 4 】

ここで図 5 を参照すると、少なくとも 1 つの実施形態による、ゲノムベース差分圧縮プログラム 1 1 0 a、1 1 0 b によって用いられる、V C F フォーマットのゲノム・データ・ファイルに対する典型的な圧縮プロセスを示すオペレーション・フローチャートが描かれている。V C F フォーマットのソース・ファイルと目的ファイルとの行内のフィールドの対比が示されている。図 5 では、ソースおよび目的ファイルは V C F フォーマットであ

50

る。

【0085】

ゲノムベース差分圧縮プログラム110a、110bが、340で、#CHROMおよびPOSが同じでないと判断した場合、342で、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルの当該行の#CHROMおよびPOSが目的ファイルの行より小さいかどうかを判断する。ゲノムベース差分圧縮プログラム110a、110bが、342で、ソース・ファイルの行の#CHROMおよびPOSは目的ファイルの行より大きいと判断した場合、320で、ゲノムベース差分圧縮プログラム110a、110bは、目的の記録を差分中に挿入し、ソース・ファイルから行を読み取るために306に戻る。

10

【0086】

また一方、ゲノムベース差分圧縮プログラム110a、110bが、342で、ソース・ファイルの行の#CHROMおよびPOSは目的ファイル中の対応する行より小さいと判断した場合、322で、差分中の対応する行は除去され、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルから次の行を読み取るために306に戻る。

【0087】

また一方、ゲノムベース差分圧縮プログラム110a、110bが、340で、ソースおよび目的ファイルのそれぞれの行の#CHROMおよびPOSは同じであると判断した場合、344で、ゲノムベース差分圧縮プログラム110a、110bは、ソースおよび目的ファイルのそれぞれの行のREFおよびALTが同じかどうかを判断する。前述の例を続けると、受信されたソースおよび目的ファイルは、事前にVCFフォーマットであると判断されていた。しかして、ゲノムベース差分圧縮プログラム110a、110bは、図5に進んでおり、本図では、受信されたソースおよび目的ファイルの第一行が、#CHROMおよびPOSが同じかどうかを判断するために評価される。受信されたソース・ファイルと目的ファイルとの対応する行の#CHROMは21であり、それぞれの行に対するPOSは3である。したがって、ソースおよび目的ファイルの対応する行の#CHROMおよびPOSは同じであり、ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイルの行中のREF、ALTが目的ファイル中の対応する行と同じかどうかを判断するために344に進む。

20

【0088】

ゲノムベース差分圧縮プログラム110a、110bが、344で、ソース・ファイルと目的ファイルとの行のREFおよびALTは同じでないと判断した場合、326で、目的ファイルからの当該行はバッファに保存され、ゲノムベース差分圧縮プログラム110a、110bは、ファイルから目的行を読み取るために312に戻る。

30

【0089】

また一方、ゲノムベース差分圧縮プログラム110a、110bが、344で、ソースおよび目的ファイルの行のREFおよびALTは同じであると判断した場合、346で、ゲノムベース差分圧縮プログラム110a、110bは、INFOフィールドを除いたコラムの差を計算する。ゲノムベース差分圧縮プログラム110a、110bは、コラムの差を調べ、どのコラムがかかる差を含むかを示すことができる。差のあるコラムの位置およびコラムの差の程度の如何によって、派生差分ファイル内の記録が影響され得る。前述の例を続けると、両方の行のREFはGであり、両方の行のALTはCである。しかして、ゲノムベース差分圧縮プログラム110a、110bは、受信されたソースおよび目的ファイルの両方の行のALTおよびREFは同じであると判断する。ソースおよび目的ファイルの行のREFおよびALTが同じなので、ゲノムベース差分圧縮プログラム110a、110bは、一切のコラムの差を計算するために、VCFフォーマットのソースおよび目的ファイル中に含まれる8つの必須の固定コラム全てを通貫して検索する。ゲノムベース差分圧縮プログラム110a、110bは、IDおよびQUALの2つのコラム中にいくつかの差を生成する。

40

【0090】

50

次いで348で、ゲノムベース差分圧縮プログラム110a、110bは、INFOフィールドをサブフィールドに構文解析する。INFOフィールド（すなわち、VCFファイルの第八コラム）は、ソースまたは目的ファイルに関する付加情報を含んでよく、ソース・ファイルと目的ファイルとのそれぞれの行の間の相似点および差を判断するために検索することができる。前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、VCFフォーマットのソースおよび目的ファイルの諸行のINFOまたは第八コラムを通貫して検索する。ソース・ファイルの行中の各INFOサブフィールドのフォーマットは、メタ情報中にDP = 154 ; MQ = 52 ; H2として指定されている。また一方、目的ファイルの行中の各INFOサブフィールドのフォーマットは、DP = 159 ; MQ = 79 ; H2である。

10

【0091】

次いで350で、ゲノムベース差分圧縮プログラム110a、110bは、目的ファイルとソース・ファイルとのそれぞれの行に対するこれらサブフィールドを対比し、ソース・ファイルから次の行を読み取るため306に戻る。ゲノムベース差分圧縮プログラム110a、110bは、最短ではソースおよび目的ファイルが終わるまで、300を通してのループを続ける。前述の例を続けると、ゲノムベース差分圧縮プログラム110a、110bは、これらサブフィールドを対比して差を計算し、それらは派生差分ファイル中に挿入される。次いで、ゲノムベース差分圧縮プログラム110a、110bは、受信されたソース・ファイルの次の行を読み取るため306に戻る。

【0092】

20

ここで図6を参照すると、少なくとも1つの実施形態による、ゲノムベース差分圧縮プログラム110aおよび110bによって用いられる、ゲノム・データ・ファイルを対比する典型的なプロセス400を示すブロック図が描かれている。図示のように、ゲノム・データに対するファイル402および404は、行およびコラムのセットに編成することが可能である。

【0093】

402のファイル1はソース・ファイルであり、404のファイル2は目的ファイルである。ゲノムベース差分圧縮プログラム110a、110bがこれら2つのファイルを対比するのに、404のマッピング位置に対比される402のマッピング位置に基づく行毎のインジケータが利用される。ゲノムベース差分圧縮プログラム110a、110bは、ソース・ファイル406の記録の1つ中に含まれるデータを、目的ファイル408中の別の記録と対比したとき、これら2つのファイルの間の差があることを判断することができる。図4は色の差を示しているが、この差は記録内の相異なるデータを含んでよい。しかして、ゲノムベース差分圧縮プログラム110a、110bは、派生Dファイルに対し以下の結果を書込むことができる。

30

差分： <一致行>

<第二コラム - 置き換え> <新規の値>

<一致行>

<一致行>

【0094】

40

これによれば、408中の新規の値を除いて、ソース・ファイル402と目的ファイル404とは、対比された行4つのうちの3つが同じであり、相互に一致する。

【0095】

ここで図7を参照すると、少なくとも1つの実施形態による、ゲノムベース差分圧縮プログラム110a、110bによって用いられる、SAMフォーマットのゲノム・データ・ファイルの階層構造を識別するための典型的なプロセス500を示すオペレーション・フローチャートが描かれている。図示のように、ゲノムベース差分圧縮プログラム110a、110bは、SAMの記録に対し、目的ファイルとソース・ファイルとの行を対比する際に、階層的な識別のシステムを利用する。

【0096】

50

ゲノムベース差分圧縮プログラム 110 a、110 b が、502 で、当該行の POS および REF が一致しない（すなわち、不一致）と判断した場合、504 で、目的ファイルとソース・ファイルとは不一致であると判定され、目的ファイルとソース・ファイルとの対比は、ソース・ファイル中の次の行に移る。ゲノムベース差分圧縮プログラム 110 a、110 b は、目的ファイルとソース・ファイルとのそれぞれの行の POS および REF を対比し、POS、REF が同じかどうかを判断する。

【0097】

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、502 で、POS および REF が一致すると判断した場合、506 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、ソースおよび目的ファイルの行は中間的な一致であると判定し、ゲノムベース差分圧縮プログラム 110 a、110 b は、508 でソースおよび目的ファイルの次の行の QNAME を対比するために行を下降する。

10

【0098】

ゲノムベース差分圧縮プログラム 110 a、110 b が、508 で、これらの行の QNAME が一致しないと判断した場合、510 で、目的ファイルとソース・ファイルとは不一致であると判定され、目的ファイルとソース・ファイルとの行の対比を終了する。また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、508 で、QNAME が一致すると判断した場合、512 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、目的およびソース・ファイルの行は中間的な一致であると判定し、ゲノムベース差分圧縮プログラム 110 a、110 b は、514 でソースおよび目的ファイルの FLAG（ビット単位）を対比するために行を下降する。

20

【0099】

ゲノムベース差分圧縮プログラム 110 a、110 b が、514 で、当該行の FLAG（ビット単位）が一致すると判断した場合、516 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、目的およびソース・ファイルは一致であると判定し、図 4 中のより複雑なオペレーション・フローチャートに進む。

【0100】

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、514 で、FLAG（ビット単位）が一致しないと判断した場合、518 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、ソース・ファイルと目的ファイルとの行は中間的不一致であると判定し、ゲノムベース差分圧縮プログラム 110 a、110 b は、520 で、その差が閾値（例えば、 $0 \times 400$ ）以下かどうかを判断することができる。520 で、差が閾値以下であった場合、522 で、目的およびソース・ファイルの当該行は一致と見なされ、一致に対する「m」が派生差分ファイルに書込まれる。また一方、520 で、差が閾値よりも大きかった場合、524 で、ソース・ファイルと目的ファイルとの行は不一致であると判断され、このソースおよび目的ファイルの行の階層的な識別は終了する。

30

【0101】

ここで図 8 を参照すると、少なくとも 1 つの実施形態による、ゲノムベース差分圧縮プログラム 110 a および 110 b によって用いられる、VCF フォーマットのゲノム・データ・ファイルの階層構造を識別するための簡単で典型的なプロセス 600 を示すオペレーション・フローチャートが描かれている。図示のように、ゲノムベース差分圧縮プログラム 110 a、110 b は、VCF の記録に対する目的およびソース・ファイルの行を対比する際に、階層的な識別のシステムを利用する。

40

【0102】

ゲノムベース差分圧縮プログラム 110 a、110 b が、602 で、当該行の #CHROM および POS が一致しない（すなわち、不一致）と判断した場合、604 で、目的ファイルとソース・ファイルとは不一致であるとして判定され、目的およびソース・ファイルの対比は終了する。ゲノムベース差分圧縮プログラム 110 a、110 b は、目的およびソース・ファイル中の各行の #CHROM および POS を対比する。

【0103】

50

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、602 で、#CHROM および POS は一致すると判断した場合、606 で、目的およびソース・ファイルの行は中間的な一致として判定され、ゲノムベース差分圧縮プログラム 110 a、110 b は、608 でソースおよび目的ファイルの次の行の REF および ALT を対比するために行を下降する。

【0104】

ゲノムベース差分圧縮プログラム 110 a、110 b が、608 で、当該行の REF および ALT は一致すると判断した場合、612 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、目的およびソース・ファイルの行は一致すると判定し、目的ファイルとソース・ファイルとの行の対比は終了する。一致に対する「m」が派生差分ファイル中に書込まれる。

10

【0105】

また一方、ゲノムベース差分圧縮プログラム 110 a、110 b が、608 で、当該行の REF および ALT は一致しないと判断した場合、610 で、ゲノムベース差分圧縮プログラム 110 a、110 b は、目的ファイルとソース・ファイルとの行は不一致であると判定し、ソースおよび目的ファイルの階層的識別は終了する。

【0106】

当然のことながら、図 2 ~ 8 は、1 つの実施形態の単なる例示を提供しており、各種の実施形態がどのように実装できるかに関して、いかなる限定をも意味するものではない。設計および実装上の要件に基づいて、描かれた（諸）実施形態に多くの修改を加えることが可能である。

20

【0107】

図 9 は本発明の或る例示的な実施形態による、図 1 に描かれたコンピュータの内部および外部の諸コンポーネントのブロック図 900 である。当然のことながら、図 9 は、単に一実装の例示を提供するものであり、環境に関するいかなる限定を意味するものでもなく、この環境には種々の実施形態を実装することができる。設計および実装上の要件に基づいて、描かれた環境に多くの修改を加えることが可能である。

【0108】

データ処理装置 902、904 は、マシン可読のプログラム命令を実行する能力のある任意の電子デバイスを表している。データ処理装置 902、904 は、スマート・フォン、コンピュータ・システム、PDA、または他の電子デバイスを表してよい。データ処理装置 902、904 によって表されてよい、コンピューティング・システム、環境、もしくは構成体またはこれらの組合せには、以下に限らないが、パーソナル・コンピュータ・システム、サーバ・コンピュータ・システム、シン・クライアント、シック・クライアント、ハンドヘルドまたはラップトップ・デバイス、マルチプロセッサ・システム、マイクロプロセッサベースのシステム、ネットワーク PC、ミニコンピュータ・システム、および前述のシステムまたはデバイスのいずれかを含む分散型クラウド・コンピューティング環境が含まれる。

30

【0109】

ユーザ・クライアント・コンピュータ 102 およびネットワーク・サーバ 112 は、図 9 に示された内部コンポーネント 902 a、b および外部コンポーネント 904 a、b のそれぞれのセットを含んでよい。内部コンポーネント 902 a、b のセットの各々は、1 つ以上のバス 912 上の、1 つ以上のプロセッサ 906、1 つ以上のコンピュータ可読 RAM 908、および 1 つ以上のコンピュータ可読 ROM 910、ならびに 1 つ以上のオペレーティング・システム 914 および 1 つ以上のコンピュータ可読有形ストレージ・デバイス 916 を含む。1 つ以上のオペレーティング・システム 914、ソフトウェア・プログラム 108、ならびにクライアント・コンピュータ 102 中のゲノムベース差分圧縮プログラム 110 a およびネットワーク・サーバ 112 中のゲノムベース差分圧縮プログラム 110 b は、1 つ以上の RAM 908（これは、通常、キャッシュ・メモリを含む）を介しての、1 つ以上のプロセッサ 906 による実行のために、1 つ以上のコンピュータ可

40

50

読有形ストレージ・デバイス 916 に格納することができる。図 9 に示された実施形態では、コンピュータ可読有形ストレージ・デバイス 916 の各々は内部ハード・ドライブの磁気ディスク・ストレージ・デバイスである。あるいは、コンピュータ可読有形ストレージ・デバイス 916 の各々は、ROM 910、EPROM、フラッシュ・メモリなどの半導体ストレージ・デバイス、またはコンピュータ・プログラムおよびデジタル情報を格納できる任意の他のコンピュータ可読有形ストレージ・デバイスである。

【0110】

また、内部コンポーネント 902 a、b の各セットは、CD-ROM、DVD、メモリ・スティック、磁気テープ、磁気ディスク、光ディスク、または半導体ストレージ・デバイスなどの、1つ以上の携帯型コンピュータ可読有形ストレージ・デバイス 920 から読み取りおよび書込みをするための R/W ドライブまたはインターフェース 918 を含む。ソフトウェア・プログラム 108、およびゲノムベース差分圧縮プログラム 110 a および 110 b などのソフトウェア・プログラムは、それぞれの携帯型コンピュータ可読有形ストレージ・デバイス 920 の 1つ以上に格納し、それぞれの R/W ドライブまたはインターフェース 918 を介して読み取り、それぞれのハード・ドライブ 916 中にロードすることができる。

10

【0111】

また、内部コンポーネント 902 a、b の各セットは、TCP/IP アダプタ・カード、ワイヤレス Wi-Fi インターフェース・カード、または 3G もしくは 4G ワイヤレス・インターフェース・カード、または他の有線またはワイヤレス通信リンクなどの、ネットワーク・アダプタ（もしくはスイッチ・ポート・カード）またはインターフェース 922 も含んでよい。クライアント・コンピュータ 102 中のソフトウェア・プログラム 108 およびゲノムベース差分圧縮プログラム 110 a、ならびにネットワーク・サーバ・コンピュータ 112 中のゲノムベース差分圧縮プログラム 110 b は、ネットワーク（例えば、インターネット、ローカル・エリア・ネットワーク、またはその他の広域ネットワーク）と、それぞれのネットワーク・アダプタまたはインターフェース 922 とを介して、外部のコンピュータ（例えば、サーバ）からダウンロードすることが可能である。ソフトウェア・プログラム 108、ならびにクライアント・コンピュータ 102 中のゲノムベース差分圧縮プログラム 110 a、およびネットワーク・サーバ 112 中のゲノムベース差分圧縮プログラム 110 b は、ネットワーク・アダプタ（もしくはスイッチ・ポート・アダプタ）またはインターフェース 922 から、それぞれのハード・ドライブ 916 の中にロードすることができる。このネットワークは、銅線、光ファイバ、ワイヤレス送信、ルータ、ファイアウォール、交換機、ゲートウェイ・コンピュータ、もしくはエッジ・サーバ、またはこれらの組合せを含んでよい。

20

30

【0112】

外部コンポーネント 904 a、b の各セットには、コンピュータ・ディスプレイ・モニター 924、キーボード 926、およびコンピュータ・マウス 928 を含めることができる。また、外部コンポーネント 904 a、b は、タッチ・スクリーン、仮想キーボード、タッチ・パッド、ポインティング・デバイス、および他の人間用のインターフェース・デバイスも含んでよい。内部コンポーネント 902 a、b の各セットは、コンピュータ・ディスプレイ・モニター 924、キーボード 926、およびコンピュータ・マウス 928 にインターフェース接続するためのデバイス・ドライバ 930 をさらに含む。デバイス・ドライバ 930、R/W ドライブもしくはインターフェース 918、およびネットワーク・アダプタもしくはインターフェース 922 は、ハードウェア、および（ストレージ・デバイス 916 もしくは ROM 910 またはその両方に格納された）ソフトウェアを含む。

40

【0113】

本開示は、クラウド・コンピューティングの詳細な説明を含むが、前もって当然のことながら、本明細書で述べる教示の実装は、クラウド・コンピューティング環境に限定されない。それどころか、本開示の諸実施形態は、現在知られた、または今後開発される任意の他の種類のコンピューティング環境に合わせて実装することができる。

50

## 【0114】

クラウド・コンピューティングは、最小の管理作業またはサービスのプロバイダとのやり取りで、迅速に立ち上げてリリースすることができる、構成可能なコンピューティング・リソース（例えば、ネットワーク、ネットワーク帯域幅、サーバ、処理、メモリ、ストレージ、アプリケーション、仮想マシン、およびサービス）の共用のプールへの、便利でオンデマンドのネットワーク・アクセスを可能にするためのサービス・デリバリのモデルである。このクラウド・モデルは、少なくとも5つの特徴と、少なくとも3つのサービス・モデルと、少なくとも4つの展開モデルとを含んでよい。

## 【0115】

特徴は次の通りである。

オンデマンド・セルフサービス：クラウド・コンシューマは、サービスのプロバイダと人間のやり取りなしに、必要な場合は自動的に、サーバ時間およびネットワーク・ストレージなどのコンピューティング能力を一方向的にセットアップすることができる。

広範なネットワーク・アクセス：諸能力はネットワークを介して利用可能で、異機種から成るシンまたはシック・クライアント・プラットフォーム（例えば、携帯電話、ラップトップ、およびPDA）による使用を推進する標準的なメカニズムを通してアクセスされる。

リソースのプール化：プロバイダのコンピューティング・リソースは、マルチテナント・モデルを用いる複数のコンシューマにサービスするために、デマンドに従って動的に割り当ておよび再割り当てされる各種の物理的および仮想のリソースとしてプール化される。一般にコンシューマは提供されたリソースの正確な場所の制御または知識を持たない、という点で、場所無依存性の感覚があるが、抽象化のより高位レベルでは場所（例えば、国、州、またはデータセンタ）を特定することを可能にする。

敏速な伸縮性：諸能力は、迅速にスケール・アウトするため、場合によっては自動的に、敏速且つ伸縮自在にセットアップでき、迅速にスケール・インするために敏速にリリースされる。コンシューマには、セットアップのため利用可能な諸能力が多くの場合無制限に見え、いつでもどのような量でも購入が可能である。

計量されるサービス：クラウド・システムは、サービスの種類（例えば、ストレージ、処理、帯域幅、および有効なユーザ・アカウント）に適した抽象化のいずれかのレベルで、計量機能を利用することによって、リソース使用を自動的に管理し、最適化する。リソース利用は、プロバイダおよび利用されるサービスのコンシューマの双方に透明性を提供しながら、モニタし、管理し、報告することができる。

## 【0116】

サービス・モデルは次のとおりである。

サービスとしてのソフトウェア（SaaS：Software as a Service）：コンシューマに提供される能力は、クラウド・インフラストラクチャ上で実行されているプロバイダのアプリケーションを使うことである。これらのアプリケーションは、様々なクライアント・デバイスから、ウェブ・ブラウザ（例えば、ウェブベースのeメール）などのシン・クライアント・インターフェースを介してアクセス可能である。コンシューマは、限られたユーザ固有のアプリケーション構成設定のあり得る例外を除いて、ネットワーク、サーバ、オペレーティング・システム、ストレージ、または個別のアプリケーション機能でさえも含め、根底にあるクラウド・インフラストラクチャを管理または制御はしない。

サービスとしてのプラットフォーム（PaaS：Platform as a Service）：コンシューマに提供される能力は、クラウド・インフラストラクチャ上に、プロバイダによってサポートされるプログラミング言語およびツールを使って生成された、コンシューマ生成の、またはコンシューマ取得のアプリケーションを展開することである。コンシューマは、ネットワーク、サーバ、オペレーティング・システム、またはストレージを含め、根底にあるクラウド・インフラストラクチャを管理または制御することはないが、これら展開されるアプリケーション、およびおそらくはアプリケーションのホステ

10

20

30

40

50

イング環境設定に対する制御を有する。

サービスとしてのインフラストラクチャ ( I a a S : I n f r a s t r u c t u r e a s a S e r v i c e ) : コンシューマに提供される能力は、コンシューマが、処理、ストレージ、ネットワーク、およびオペレーティング・システムおよびアプリケーションを含み得る、任意のソフトウェアを展開し実行することが可能な他の基本的コンピューティング・リソースをセットアップすることである。コンシューマは、根底にあるクラウド・インフラストラクチャを管理または制御することはないが、オペレーティング・システム、ストレージ、展開されるアプリケーションに対する制御、およびおそらくは選択したネットワーク・コンポーネント (例えば、ホストのファイアウォール) の限定された制御を有する。

10

【 0 1 1 7 】

展開モデルは次のとおりである。

プライベート・クラウド : このクラウド・インフラストラクチャは、一組織のためだけに運営される。これは、その組織または第三者によって管理されてよく、自組織の構内に所在しても自組織の構外に所在してもよい。

コミュニティ・クラウド : このクラウド・インフラストラクチャは、いくつかの組織によって共有され、共有の利害関係 (例えば、任務、安全要件、指針、およびコンプライアンス配慮事項) を有する特定のコミュニティをサポートする。これは、これらの組織または第三者によって管理されてよく、これら組織の構内に所在してもこれら組織の構外に所在してもよい。

20

パブリック・クラウド : このクラウド・インフラストラクチャは、一般公衆または大きな産業グループが利用可能なようにされており、クラウド・サービスを販売する組織によって所有されている。

ハイブリッド・クラウド : このクラウド・インフラストラクチャは、独自のエンティティに留まりながら、データおよびアプリケーションの可搬性 (例えば、クラウド間の負荷バランスのためのクラウド・パースティング) を可能にする標準化されたまたは独自の技術によって一緒に結ばれた2つ以上のクラウド (プライベート、コミュニティ、またはパブリック) の合成体である。

【 0 1 1 8 】

クラウド・コンピューティング環境は、ステートレスネス、弱連結、モジュール性、および意味相互運用性に焦点を合わせて方向付けられたサービスである。クラウド・コンピューティングの中心には、相互接続されたノードのネットワークを含むインフラストラクチャがある。

30

【 0 1 1 9 】

図 1 0 を参照すると、例示的なクラウド・コンピューティング環境 1 0 0 0 が描かれている。図示のように、クラウド・コンピューティング環境 1 0 0 0 は、例えば、携帯情報端末 ( P D A : p e r s o n a l d i g i t a l a s s i s t a n t ) もしくはセルラ電話 1 0 0 0 A、デスクトップ・コンピュータ 1 0 0 0 B、ラップトップ・コンピュータ 1 0 0 0 C、もしくは車載コンピュータ・システム 1 0 0 0 N、またはこれらの組合せなど、クラウド・コンシューマによって使用されるローカルのコンピューティング・デバイスが通信可能な1つ以上のクラウド・コンピューティング・ノード 1 0 0 を含む。ノード 1 0 0 は、相互に通信することができる。これらは、前述したプライベート、コミュニティ、パブリック、またはハイブリッド・クラウド、またはこれらの組合せなど、1つ以上のネットワークに物理的にまたは仮想的にグループ化する (図示せず) ことが可能である。これは、クラウド・コンピューティング環境 1 0 0 0 が、クラウド・コンシューマがローカル・コンピューティング・デバイス上にリソースを維持する必要のない、サービスとしてのインフラストラクチャ、プラットフォーム、もしくはソフトウェア、またはこれらの組合せを提供することを可能にする。当然のことながら、図 1 0 に示されたコンピューティング・デバイス 1 0 0 0 A ~ N の種類は例示だけを意図されたものであり、コンピューティング・ノード 1 0 0 およびクラウド・コンピューティング環境 1 0 0 0 は、任意の

40

50

種類のネットワークもしくは(例えば、ウェブ・ブラウザを使って)ネットワーク・アドレス指定が可能な接続またはその両方を介して、任意の種類のコピュータ化デバイスと通信することができる。

#### 【0120】

ここで図11を参照すると、クラウド・コンピューティング環境1000によって設けられる機能的抽象化層のセット1100が示されている。前もって当然のことながら、図11に示されたコンポーネント、層、および機能は例示だけを意図されたものであり、本発明の諸実施形態はこれに限定されない。図示のように、以下の層および対応する機能が設けられている。

#### 【0121】

ハードウェアおよびソフトウェア層1102は、ハードウェアおよびソフトウェア・コンポーネントを含む。ハードウェア・コンポーネントの例には、メインフレーム1104、RISC(Reduced Instruction Set Computer(縮小命令セット・コンピュータ))アーキテクチャ・ベースのサーバ1106、サーバ1108、ブレード・サーバ1110、ストレージ・デバイス1112、ならびにネットワークおよびネットワークング・コンポーネント1114が含まれる。いくつかの実施形態において、ソフトウェア・コンポーネントは、ネットワーク・アプリケーション・サーバ・ソフトウェア1116およびデータベース・ソフトウェア1118を含む。

#### 【0122】

仮想化層1120は、仮想サーバ1122、仮想ストレージ1124、仮想プライベート・ネットワークを含む仮想ネットワーク1126、仮想アプリケーションおよびオペレーティング・システム1128、ならびに仮想クライアント1130、の仮想エンティティの諸例を設けることが可能な抽象化層を提供する。

#### 【0123】

一例において、管理層1132には以下に記載の機能を備えることが可能である。リソース・セットアップ1134は、クラウド・コンピューティング環境内でタスクを実行するために利用されるコンピューティング・リソースおよび他のリソースの動的な調達を提供する。計量および価格計算1136は、クラウド・コンピューティング環境内でリソースが利用されるのに応じて、これらリソースの消費に対するコストの追跡、および請求書または納入書の作成を提供する。一例において、これらのリソースには、アプリケーション・ソフトウェア・ライセンスを含めてよい。セキュリティは、クラウド・コンシューマおよびタスクに対する身元検証、ならびにデータおよび他のリソースに対する保護を提供する。ユーザ・ポータル1138は、コンシューマおよびシステム管理者に対し、クラウド・コンピューティング環境へのアクセスを提供する。サービス・レベル管理1140は、要求されるサービス・レベルが満たされるように、クラウド・コンピューティング・リソースの割り当ておよび管理を提供する。サービス・レベル合意(SLA: Service Level Agreement)計画および達成1142は、SLAに沿って将来の必要性が予期されるクラウド・コンピューティング・リソースに対する事前配置およびそれらリソースの調達を提供する。

#### 【0124】

作業負荷層1144は、クラウド・コンピューティング環境で使用できる機能の諸例を提供する。この層から提供できる作業負荷および機能の例には、マッピングおよびナビゲーション1146、ソフトウェア開発およびライフサイクル管理1148、仮想教室教育配信1150、データ解析処理1152、トランザクション処理1154、およびゲノムベース差分圧縮1156が含まれる。ゲノムベース差分圧縮プログラム110a、110bは、ゲノム・データ・ファイルに対する差分ファイルを圧縮する方法を提供する。

#### 【0125】

本発明の様々な実施形態の説明は、例示目的で提示されたもので、網羅的であることも、または開示された態様に限定することも意図されていない。当業者には、記載された諸態様の範囲から逸脱することのない多くの修改および別形が明白であろう。本明細書で用

10

20

30

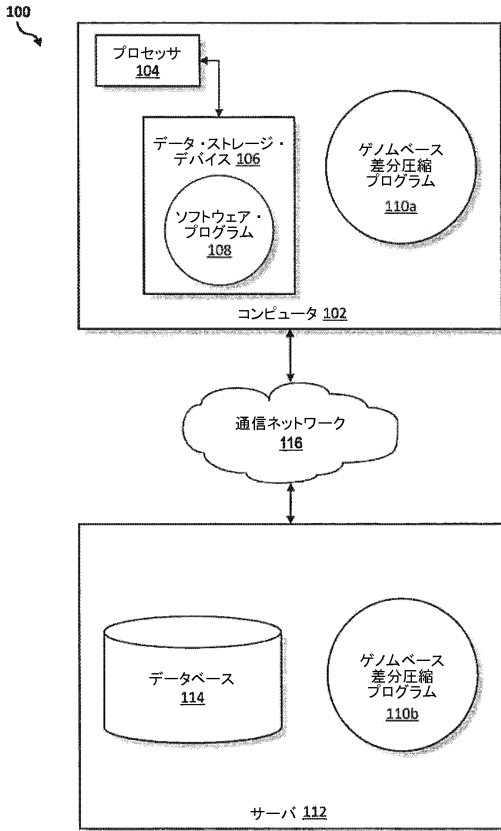
40

50

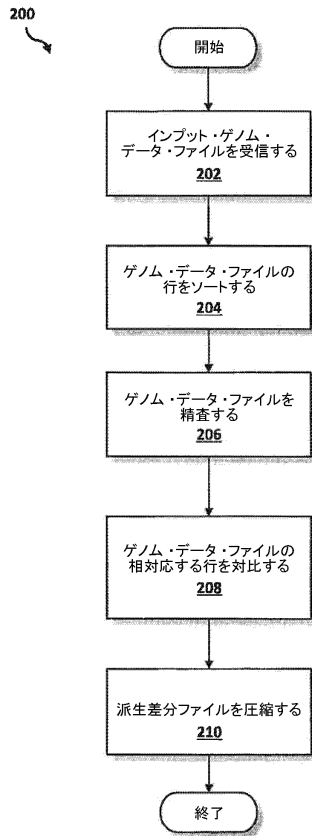
いられた用語は、諸実施形態の原理、実際上の応用、または市販の技術の技術的な改良を最善に説明し、または他の当業者が本明細書に開示された諸実施形態を理解できるように選択されたものである。

【図面】

【図 1】



【図 2】



10

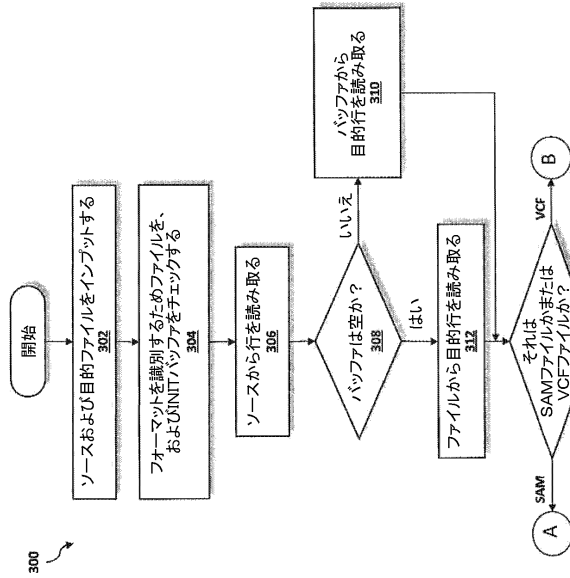
20

30

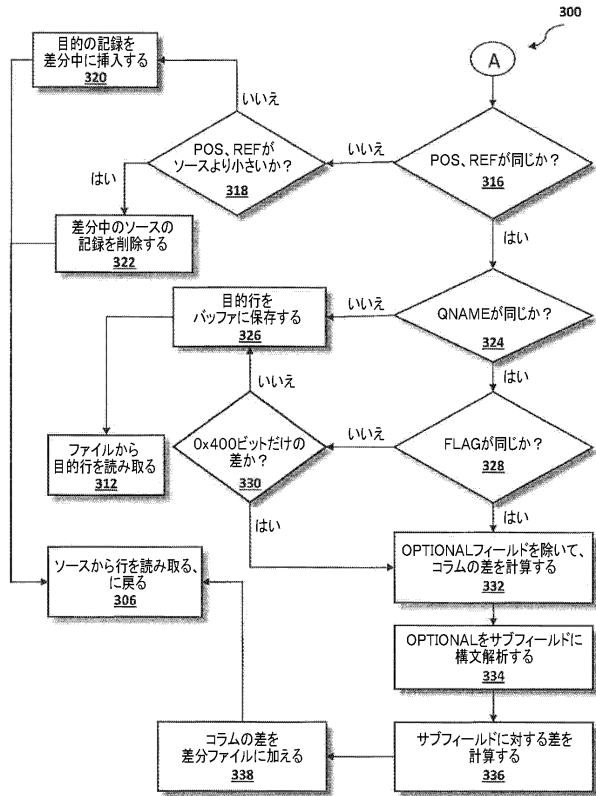
40

50

【図3】



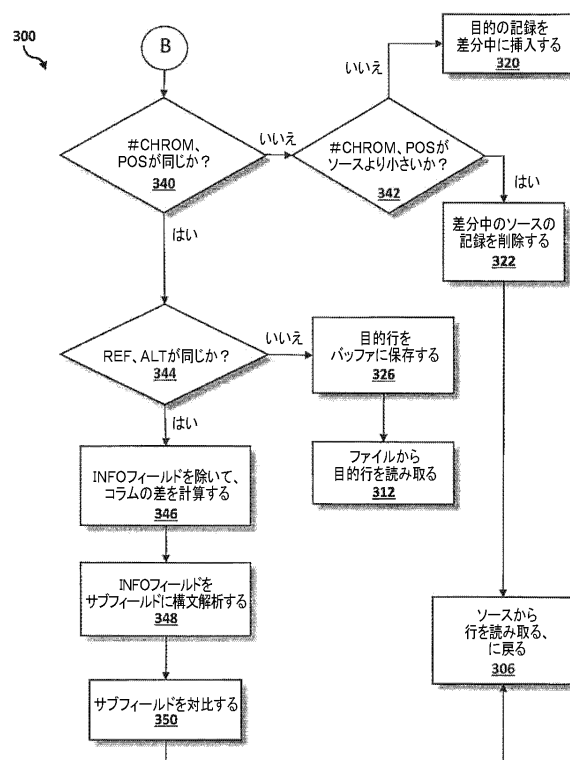
【図4】



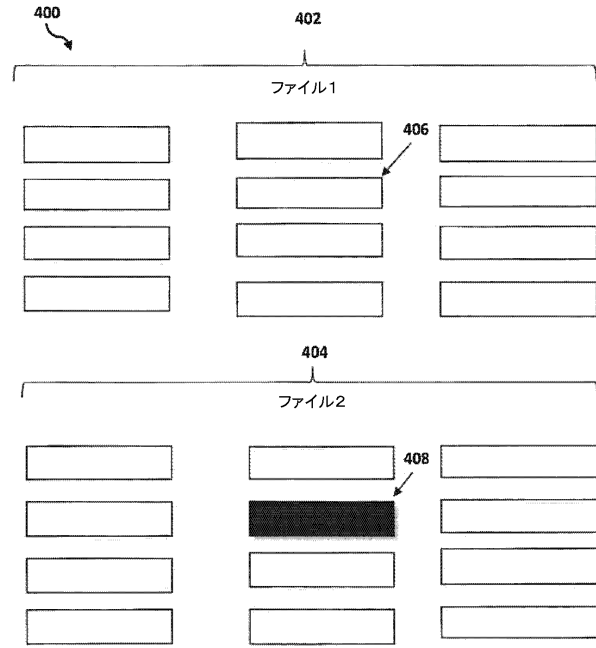
10

20

【図5】



【図6】

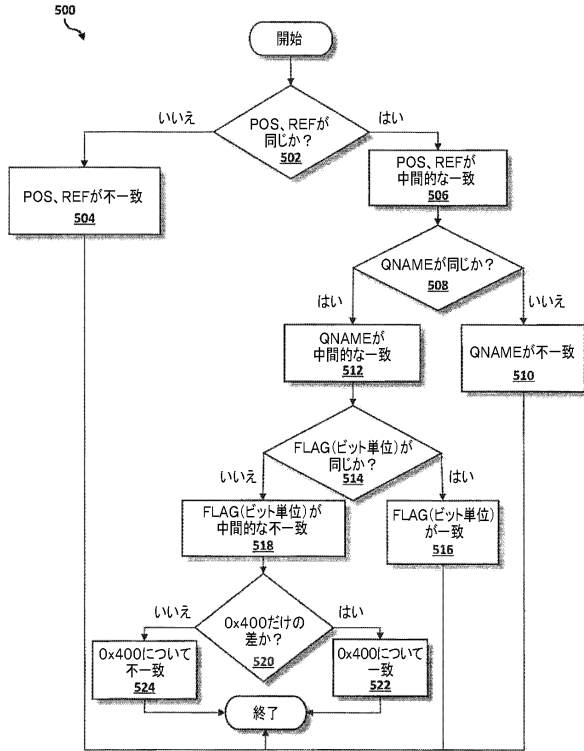


30

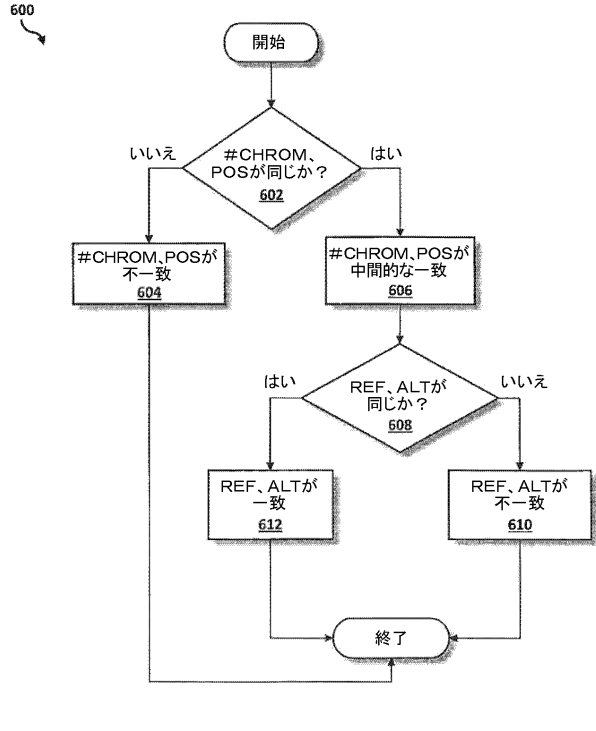
40

50

【図7】



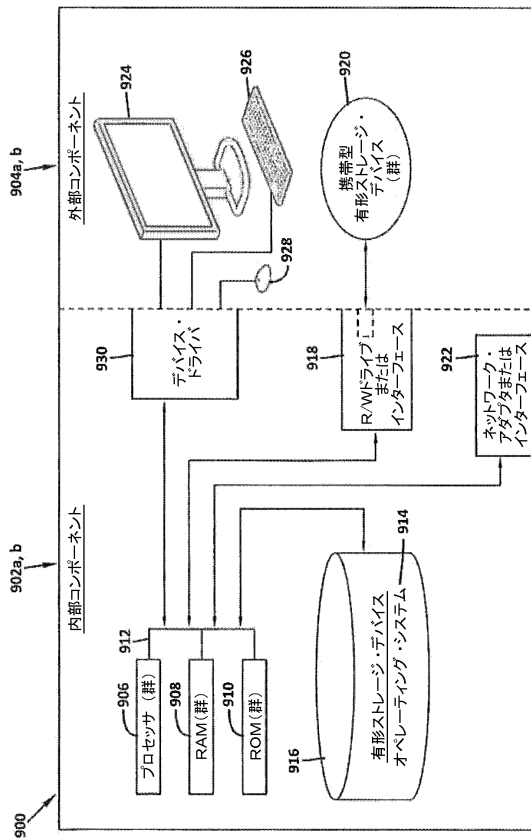
【図8】



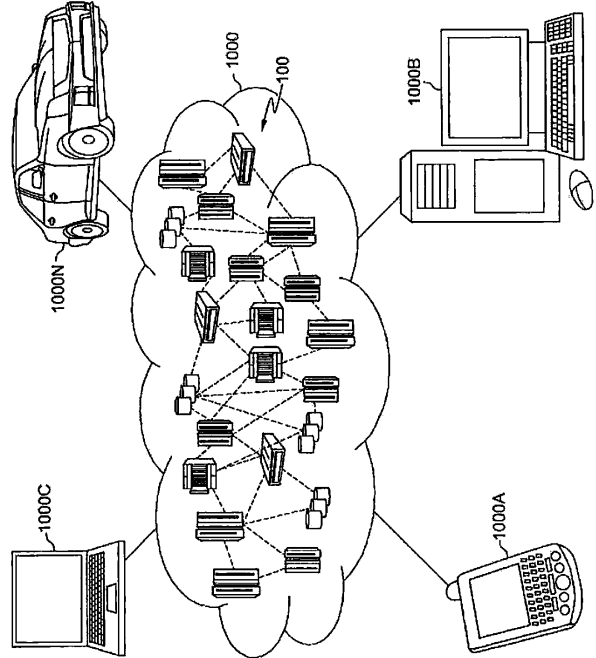
10

20

【図9】



【図10】

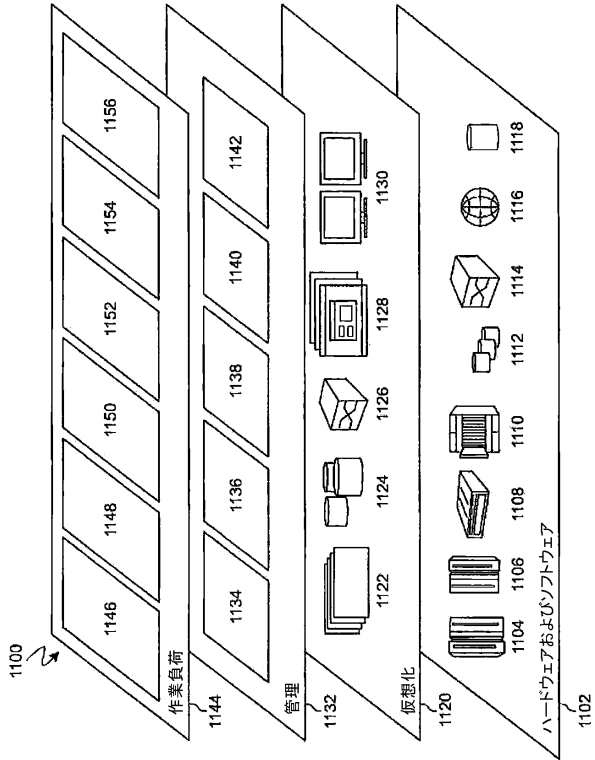


30

40

50

【 1 1 1 】



10

20

30

40

50

---

フロントページの続き

- (72)発明者 マハラナ、アジャシャ  
アメリカ合衆国 9 5 2 1 - 6 0 9 9 カリフォルニア州 サンノゼ ハリー・ロード 6 5 0
- (72)発明者 コンスタンチネスキュ、ミハイル、コルネリウ  
アメリカ合衆国 9 5 2 1 - 6 0 9 9 カリフォルニア州 サンノゼ ハリー・ロード 6 5 0
- 審査官 山内 裕史
- (56)参考文献 米国特許出願公開第 2 0 1 3 / 0 1 3 2 3 5 3 ( U S , A 1 )  
米国特許出願公開第 2 0 0 8 / 0 0 7 7 6 0 7 ( U S , A 1 )  
米国特許出願公開第 2 0 1 1 / 0 1 1 9 2 4 0 ( U S , A 1 )
- (58)調査した分野 (Int.Cl. , D B 名)  
G 1 6 B 5 / 0 0 - 9 9 / 0 0