

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7631330号  
(P7631330)

(45)発行日 令和7年2月18日(2025.2.18)

(24)登録日 令和7年2月7日(2025.2.7)

|                         |               |       |
|-------------------------|---------------|-------|
| (51)国際特許分類              | F I           |       |
| G 1 6 B 50/50 (2019.01) | G 1 6 B 50/50 |       |
| G 0 6 F 12/00 (2006.01) | G 0 6 F 12/00 |       |
| G 0 6 F 16/17 (2019.01) | G 0 6 F 16/17 | 1 0 0 |
| G 0 6 F 16/22 (2019.01) | G 0 6 F 16/22 |       |
| H 0 3 M 7/30 (2006.01)  | H 0 3 M 7/30  | Z     |
| 請求項の数 20 (全42頁)         |               |       |

|                   |                             |          |                          |
|-------------------|-----------------------------|----------|--------------------------|
| (21)出願番号          | 特願2022-522858(P2022-522858) | (73)特許権者 | 590000248                |
| (86)(22)出願日       | 令和2年10月17日(2020.10.17)      |          | コーニンクレッカ フィリップス エヌ       |
| (65)公表番号          | 特表2022-553199(P2022-553199  |          | ヴェ                       |
|                   | A)                          |          | Koninklijke Philips      |
| (43)公表日           | 令和4年12月22日(2022.12.22)      |          | N.V.                     |
| (86)国際出願番号        | PCT/EP2020/079298           |          | オランダ国 5 6 5 6 アーヘー アイン   |
| (87)国際公開番号        | WO2021/074440               |          | ドーフエン ハイテック キャンパス 5 2    |
| (87)国際公開日         | 令和3年4月22日(2021.4.22)        |          | High Tech Campus 52 ,    |
| 審査請求日             | 令和5年10月12日(2023.10.12)      |          | 5 6 5 6 AG Eindhoven , N |
| (31)優先権主張番号       | 62/923,141                  |          | etherlands               |
| (32)優先日           | 令和1年10月18日(2019.10.18)      | (74)代理人  | 110001690                |
| (33)優先権主張国・地域又は機関 | 米国(US)                      |          | 弁理士法人M&Sパートナーズ           |
| (31)優先権主張番号       | 62/956,952                  | (72)発明者  | チャンダク シュブハム              |
| (32)優先日           | 令和2年1月3日(2020.1.3)          |          | オランダ国 5 6 5 6 アーヘー アイン   |
|                   | 最終頁に続く                      |          | ドーフエン ハイ テック キャンパス 5     |
|                   |                             |          | 最終頁に続く                   |

(54)【発明の名称】 多様な表形式データの効果的な圧縮、表現、および展開のためのシステムおよび方法

(57)【特許請求の範囲】

【請求項1】

データの圧縮を制御するための方法であって、前記方法は、  
 複数の第1のファイル形式のうちの1つの第1のファイル形式の表形式データにアクセスするステップと、  
 前記表形式データから複数の属性を抽出するステップと、  
 前記表形式データを複数のチャンクに分割するステップと、  
 前記抽出された属性および前記チャンクを処理して関連付けられた情報にするステップと、

前記関連付けられた情報内で識別される前記属性および前記チャンクに対して、異なるコンプレッサおよび依存関係属性を選択するステップと、

チャンク間で共有されるコンプレッサのためのグローバルデータを保存するステップと、  
 前記関連付けられた情報と、前記関連付けられた情報内に示される前記チャンクおよび前記属性のための前記異なるコンプレッサを示す情報とを含む第2のファイル形式のファイル生成するステップとを含み、

前記複数の第1のファイル形式は、互いに互換性がなく、前記第2のファイル形式とは異なり、前記異なるコンプレッサを示す前記情報は、前記関連付けられた情報内に示される前記属性および前記チャンクの選択的展開を可能にするために、前記第2のファイル形式に処理される、方法。

【請求項2】

10

20

前記相関付けられた情報は、  
 前記属性のうちの一つまたは複数の属性を示す第 1 の情報、および  
 前記第 1 の情報内に示される前記一つまたは複数の属性に関連付けられたチャンクに対応する第 2 の情報を含む、  
 少なくとも一つのテーブルを含む、請求項 1 に記載の方法。

【請求項 3】

前記第 1 の情報は二次元セル配列を含み、各セルは、前記二次元セル配列に対応する前記チャンク内に含まれる一つまたは複数の対応する属性を特定する、請求項 2 に記載の方法。

【請求項 4】

前記第 1 の情報は、前記チャンクに関連する次元固有属性の少なくとも一つの次元テーブルを含む、請求項 3 に記載の方法。

【請求項 5】

前記異なるコンプレッサのうちの一つのコンプレッサが前記ファイルに埋め込まれている、請求項 1 に記載の方法。

【請求項 6】

前記異なるコンプレッサのうちの一つは、前記表形式データの一つまたは複数の属性をスパース配列として符号化する、請求項 1 に記載の方法。

【請求項 7】

前記表形式データは異なるサイズのチャンクに分割される、請求項 1 に記載の方法。

【請求項 8】

前記方法はさらに、  
 一つまたは複数の属性固有インデックス生成するステップと、  
 前記一つまたは複数の属性固有インデックスを前記第 2 の ファイル形式の前記ファイルに組み込むステップとを含み、前記属性固有インデックスは、前記表形式データの特定の属性の値または範囲を含むチャンクを識別することを可能にする、請求項 1 に記載の方法。

【請求項 9】

前記方法はさらに、  
 前記第 2 のファイル形式の前記ファイルを M P E G - G ファイルに統合するステップを含む、請求項 1 に記載の方法。

【請求項 10】

前記方法はさらに、  
 前記相関付けられた情報のためのアクセス制御ポリシー情報を生成するステップと、  
 前記アクセス制御ポリシー情報を前記第 2 のファイル形式の前記ファイルに統合するステップとを含み、

前記アクセス制御ポリシー情報は、前記相関付けられた情報の第 1 の部分のための第 1 のアクセスレベルを示す第 1 の情報、および前記相関付けられた情報の第 2 の部分のための第 2 のアクセスレベルを示す第 2 の情報を含み、前記第 2 のアクセスレベルは前記第 1 のアクセスレベルとは異なる、請求項 1 に記載の方法。

【請求項 11】

データの圧縮を制御するための装置であって、前記装置は、  
 命令を保存するように構成されたメモリと、  
 前記命令を実行して動作を実行するように構成された少なくとも一つのプロセッサとを備え、前記動作は、  
 複数の第 1 のファイル形式のうちの一つの第 1 のファイル形式の表形式データにアクセスするステップと、  
 前記表形式データから複数の属性を抽出するステップと、  
 前記表形式データを複数のチャンクに分割するステップと、  
 前記抽出された属性および前記チャンクを処理して相関付けられた情報にするステップと、

10

20

30

40

50

前記相関付けられた情報内で識別される前記属性および前記チャンクに対して、異なるコンプレッサおよび依存関係属性を選択するステップと、

チャンク間で共有されるコンプレッサのためのグローバルデータを保存するステップと、前記相関付けられた情報と、前記相関付けられた情報内に示される前記チャンクおよび前記属性のための前記異なるコンプレッサを示す情報とを含む第2のファイル形式のファイル生成するステップとを含み、

前記複数の第1のファイル形式は、互いに互換性がなく、前記第2のファイル形式とは異なり、前記異なるコンプレッサを示す前記情報は、前記相関付けられた情報内に示される前記属性および前記チャンクの選択的展開を可能にするために、前記第2のファイル形式に処理される、装置。

10

【請求項12】

前記相関付けられた情報は、

前記属性のうちの1つまたは複数の属性を示す第1の情報、および

前記第1の情報内に示される前記1つまたは複数の属性に関連付けられたチャンクに対応する第2の情報を含む少なくとも1つのテーブルを含む、請求項11に記載の装置。

【請求項13】

前記第1の情報は二次元セル配列を含み、各セルは、前記二次元セル配列に対応する前記チャンク内に含まれる1つまたは複数の対応する属性を特定する、請求項12に記載の装置。

【請求項14】

20

前記第1の情報は、前記チャンクに関連する次元固有属性の少なくとも1つの次元テーブルを含む、請求項13に記載の装置。

【請求項15】

前記異なるコンプレッサのうちの少なくとも1つのコンプレッサが前記ファイルに埋め込まれている、請求項1に記載の方法。

【請求項16】

前記異なるコンプレッサのうちの少なくとも1つは、前記表形式データの1つまたは複数の属性をスパース配列として符号化する、請求項1に記載の方法。

【請求項17】

前記表形式データは異なるサイズのチャンクに分割される、請求項11に記載の装置。

30

【請求項18】

前記少なくとも1つのプロセッサは、前記命令を実行して、

1つまたは複数の属性固有インデックス生成し、

前記1つまたは複数の属性固有インデックスを前記第2のファイル形式の前記ファイルに組み込むように構成され、前記属性固有インデックスは、前記表形式データの特定の属性の値または範囲を含むチャンクを識別することを可能にする、請求項11に記載の装置。

【請求項19】

前記少なくとも1つのプロセッサは、前記命令を実行して、前記第2のファイル形式の前記ファイルをMPEG-Gファイルに統合するように構成される、請求項11に記載の装置。

40

【請求項20】

前記少なくとも1つのプロセッサは、前記命令を実行して、

前記相関付けられた情報のためのアクセス制御ポリシー情報を生成し、

前記アクセス制御ポリシー情報を前記第2のファイル形式の前記ファイルに統合するように構成され、

前記アクセス制御ポリシー情報は、前記相関付けられた情報の第1の部分のための第1のアクセスレベルを示す第1の情報、および前記相関付けられた情報の第2の部分のための第2のアクセスレベルを示す第2の情報を含み、前記第2のアクセスレベルは前記第1のアクセスレベルとは異なる、請求項11に記載の装置。

【発明の詳細な説明】

50

**【関連出願の相互参照】****【0001】**

【0001】本出願は、2019年10月18日に出願された米国仮特許出願第62/923,141号に関連し、その全内容があらゆる目的のために参照により本明細書に援用される。

**【0002】**

【0002】本出願は、2019年10月18日に出願された米国仮特許出願第62/923,113号に関連し、その全内容があらゆる目的のために参照により本明細書に援用される。

**【0003】**

【0003】本出願は、本出願と同時に提出された米国仮特許出願第62/956,941号(代理人整理番号2019P00831US01)、“Customizable Delimited Text Compression Framework”に関連し、その全内容があらゆる目的のために参照により本明細書に援用される。

**【技術分野】****【0004】**

【0004】本明細書の1つまたは複数の実施形態は、これらに限定はされないが、表形式データおよび区切られたテキストファイルなどの情報の圧縮および展開を管理することに関する。

**【背景技術】****【0005】**

【0005】ゲノムデータは多くの研究用途にとって重要であることが知られている。RAW形式のゲノムデータは、多数の配列決定されたゲノムを含む可能性がある。場合によっては、管理できないほどに配列の数が多い。したがって、RAWデータから関連性のある情報を抽出する、またはより適切に解釈するために様々な試みがなされてきた。

**【0006】**

【0006】1つの処理技術はアノテーションツールの使用を伴う。そのようなツールは、ゲノムデータ内のコーディング領域、およびそれらの対応する位置を定めるために使用され得る。アノテーションはまた、疾患や遺伝的異常に関する情報を伝達する可能性のある反復領域の数および空間分布の指標を提供し得る。ゲノムアノテーションデータの例として、マッピング統計量、定量的ブラウザトラック、バリエーション、ゲノム機能的アノテーション、遺伝子発現データ、およびHi-Cコンタクト行列などが挙げられる。

**【0007】**

【0007】アノテーションデータ(およびその他のゲノム関連情報)は、現在、様々なファイル形式、例えば、バリエーションコール形式(VCF)、browser extensible data(BED)形式、およびwiggle(WIG)形式などで表されている。これらの形式は互いに互換性を有さない。結果として、これらの形式の使用は、相互運用性に関連する問題を伴い、データを視覚化できるようにするために頻繁に形式を変換する必要がある。また、単一の統一および標準化された形式がないため、圧縮アルゴリズムへの取り組みが妨げられ、次善の圧縮アルゴリズム(例えば、gzip)が広く使用されることとなった。

**【0008】**

【0008】ゲノム情報を圧縮するための既存のアルゴリズムに関連する欠点は、とりわけ、選択性の欠如に関連する。例えば、ゲノムアノテーションデータは通常、異なる統計的特性を有する複数のフィールド(属性)を含む。既存のアルゴリズムは、すべてのフィールドをまとめて圧縮するため、個別のデータフィールド内の情報を選択的に認識、使用、および抽出することができない。また、既存のアルゴリズムでは、すべての属性を展開せずに特定のフィールドを抽出することはできない。

**【0009】**

【0009】その他の欠点は、既存の圧縮方法が一般的に適用できないことに関連する

10

20

30

40

50

。例えば、これらの方法は1つの標準コンプレッサ、または少数の標準コンプレッサのセットのみに依存するため、1つのタイプのアノテーションデータにしか適用できない。これらの方法は、選択的暗号化を実行することもできず、また、複数のアノテーションデータセットを相互に、およびシーケンシングデータとリンクさせることもできない。いくつかの特殊化された方法が提案されている。しかし、これらの方法の多くは、ディスクベースの配列（アレイ）管理ツール（例えば、TileDBおよびHDF5）に基づいており、これらのツールは、高レベルのフィーチャ、例えば、限定はされないが、メタデータ、リンケージ、および属性固有インデックス付与などを有さない。

#### 【発明の概要】

##### 【0010】

[0010] 様々な例示的实施形態の簡単な要約を以下に提示する。以下の要約において、いくつかの簡略化および省略がなされ得るが、これは様々な例示的实施形態のいくつかの態様を強調および紹介することを意図したものであり、本発明の範囲を限定することは意図されていない。当業者が実施形態を作成および使用することを可能にするのに十分な例示的实施形態の詳細な説明が後になされる。

##### 【0011】

[0011] 1つまたは複数の実施形態によれば、データの圧縮を制御するための方法は、複数の第1のファイル形式のうちの1つの第1のファイル形式のゲノムアノテーションデータにアクセスするステップと、ゲノムアノテーションデータから複数の属性を抽出するステップと、ゲノムアノテーションデータを複数のチャンクに分割するステップと、抽出された属性およびチャンクを処理して関連付けられた情報にするステップと、関連付けられた情報内で識別される属性およびチャンクに対して、異なるコンプレッサを選択するステップと、関連付けられた情報と、関連付けられた情報内に示されるチャンクおよび属性のための異なるコンプレッサを示す情報とを含む第2のファイル形式のファイルを生成するステップとを含み、複数の第1のファイル形式は互いに互換性がなく、第2のファイル形式とは異なり、異なるコンプレッサを示す情報は、関連付けられた情報内に示される属性およびチャンクの選択的展開を可能にするために、第2のファイル形式に処理される。

##### 【0012】

[0012] 関連付けられた情報は、属性のうちの1つまたは複数の属性を示す第1の情報、および第1の情報内に示される1つまたは複数の属性に関連付けられたチャンクに対応する第2の情報を含む少なくとも1つのテーブルを含み得る。第1の情報は二次元セル配列を含み、各セルは、二次元セル配列に対応するチャンク内に含まれる1つまたは複数の対応する属性を特定し得る。第1の情報は、チャンクに関連する次元固有属性の少なくとも1つの一次元テーブルを含み得る。ゲノムアノテーションデータは同じサイズのチャンクに分割されてもよい。ゲノムアノテーションデータは異なるサイズのチャンクに分割されてもよい。

##### 【0013】

[0013] 方法は、関連付けられた情報をリンクする1つまたは複数の第1のインデックスを生成するステップと、関連付けられた情報を、異なるコンプレッサを示す情報にリンクする1つまたは複数の第2のインデックスを生成するステップとを含み、1つまたは複数の第1のインデックスおよび1つまたは複数の第2のインデックスは、第2のファイル形式のファイルに処理される。方法は、第2のファイル形式のファイルをMPEG-Gファイルに統合するステップを含み得る。

##### 【0014】

[0014] 方法は、関連付けられた情報のためのアクセス制御ポリシー情報を生成するステップと、アクセス制御ポリシー情報を第2のファイル形式のファイルに統合するステップとを含み、アクセス制御ポリシー情報は、関連付けられた情報の第1の部分のための第1のアクセスレベルを示す第1の情報、および関連付けられた情報の第2の部分のための第2のアクセスレベルを示す第2の情報を含み、第2のアクセスレベルは第1のアク

10

20

30

40

50

セスレベルとは異なる。抽出された属性は、染色体属性およびゲノム位置属性のうちの少なくとも1つを含み得る。

【0015】

【0015】1つまたは複数の実施形態によれば、データの圧縮を制御するための装置は、命令を保存するように構成されたメモリと、命令を実行して動作を実行するように構成された少なくとも1つのプロセッサとを備え、動作は、複数の第1のファイル形式のうちの1つの第1のファイル形式のゲノムアノテーションデータにアクセスするステップと、ゲノムアノテーションデータから属性を抽出するステップと、ゲノムアノテーションデータをチャンクに分割するステップと、抽出された属性およびチャンクを処理して関連付けられた情報にするステップと、関連付けられた情報内で識別される属性およびチャンクに対して、異なるコンプレッを選択するステップと、関連付けられた情報と、関連付けられた情報内に示されるチャンクおよび属性のための異なるコンプレッサを示す情報とを含む第2のファイル形式のファイルを生成するステップとを含む。複数の第1のファイル形式は互いに互換性がなく、第2のファイル形式とは異なる。異なるコンプレッサを示す情報は第2のファイル形式に処理され、関連付けられた情報内に示される属性およびチャンクの選択的展開を可能にする。

10

【0016】

【0016】関連付けられた情報は、属性のうちの1つまたは複数の属性を示す第1の情報、および第1の情報内に示される1つまたは複数の属性に関連付けられたチャンクに対応する第2の情報を含む少なくとも1つのテーブルを含み得る。第1の情報は二次元セル配列を含み、各セルは、二次元セル配列に対応するチャンク内に含まれる1つまたは複数の対応する属性を特定する。第1の情報は、チャンクに関連する次元固有属性の少なくとも1つの一次元テーブルを含み得る。ゲノムアノテーションデータは同じサイズのチャンクに分割されてもよい。ゲノムアノテーションデータは異なるサイズのチャンクに分割されてもよい。

20

【0017】

【0017】少なくとも1つのプロセッサは、命令を実行して、関連付けられた情報をリンクする1つまたは複数の第1のインデックスを生成し、関連付けられた情報を、異なるコンプレッサを示す情報にリンクする1つまたは複数の第2のインデックスを生成するように構成され、1つまたは複数の第1のインデックスおよび1つまたは複数の第2のインデックスは、第2のファイル形式のファイルに処理される。少なくとも1つのプロセッサは、命令を実行して、第2のファイル形式のファイルをMPEG-Gファイルに統合するように構成される。

30

【0018】

【0018】少なくとも1つのプロセッサは、命令を実行して、関連付けられた情報のためのアクセス制御ポリシー情報を生成し、アクセス制御ポリシー情報を第2のファイル形式のファイルに統合するように構成され、アクセス制御ポリシー情報は、関連付けられた情報の第1の部分のための第1のアクセスレベルを示す第1の情報、および関連付けられた情報の第2の部分のための第2のアクセスレベルを示す第2の情報を含み、第2のアクセスレベルは第1のアクセスレベルとは異なる。抽出された属性は、染色体属性およびゲノム位置属性のうちの少なくとも1つを含み得る。

40

【図面の簡単な説明】

【0019】

【0019】添付図面では、類似する参照番号は、個々の図にわたって、同一のまたは機能的に類似する要素を指す。添付図面は以下の詳細な説明とともに明細書に組み込まれ、その一部を形成する。添付図面は、特許請求の範囲に示される概念の例示的实施形態をさらに説明する役割を果たし、それらの実施形態の様々な原理および利点を説明する。

【0020】

【0020】上記および他のより詳細かつ具体的な特徴が、添付図面を参照しつつ、以下の明細書においてより完全に開示されている。

50

【図 1】 [ 0 0 2 1 ] 図 1 は、ゲノムアノテーションデータの圧縮および展開を制御するためのシステムの実施形態を示す。

【図 2】 [ 0 0 2 2 ] 図 2 は、ゲノムアノテーションデータの圧縮および展開を制御するための方法を示す。

【図 3】 [ 0 0 2 3 ] 図 3 は、ゲノムアノテーションデータのための統一および標準化されたファイル形式の実施形態を示す。

【図 4】 [ 0 0 2 4 ] 図 4 は、ゲノムアノテーションデータのためのアノテーションファイルに情報を処理するための実施形態を示す。

【図 5】 [ 0 0 2 5 ] 図 5 は、統合および標準化されたファイルのファイルヘッダに情報を処理するための実施形態を示す。

10

【図 6】 [ 0 0 2 6 ] 図 6 は、統一および標準化されたファイルのために情報を処理するための実施形態を示す。

【図 7】 [ 0 0 2 7 ] 図 7 は、統合および標準化されたファイルのために圧縮情報を処理するための実施形態を示す。

【図 8】 [ 0 0 2 8 ] 図 8 は、一次元配列のために情報を処理するための実施形態を示す。

【図 9】 [ 0 0 2 9 ] 図 9 は、二次元配列に情報を処理するための実施形態を、次元固有属性のための配列とともに示す。

【図 10】 [ 0 0 3 0 ] 図 10 は、ゲノムアノテーションデータのためにインデックス付与される属性情報を処理するための実施形態を示す。

【図 11】 [ 0 0 3 1 ] 図 11 は、追加の属性関連情報を処理するための実施形態を示す。

20

【図 12】 [ 0 0 3 2 ] 図 12 は、属性情報を処理するための実施形態を示す。

【図 13】 [ 0 0 3 3 ] 図 13 は、データのチャンクを選択的に圧縮するための実施形態を示す。

【図 14】 [ 0 0 3 4 ] 図 14 は、圧縮ファイルインデックス情報を処理するための実施形態を示す。

【図 15】 [ 0 0 3 5 ] 図 15 は、インデックス情報を処理するための実施形態を示す。

【図 16】 [ 0 0 3 6 ] 図 16 は、チャンクを処理するための実施形態を示す。

【図 17】 [ 0 0 3 7 ] 図 17 は、固定サイズのチャンクを生成するための実施形態を示す。

【図 18】 [ 0 0 3 8 ] 図 18 は、可変サイズのチャンクを生成するための実施形態を示す。

30

【図 19】 [ 0 0 3 9 ] 図 19 は、チャンクおよび属性情報を処理するための実施形態を示す。

【図 20】 [ 0 0 4 0 ] 図 20 は、チャンクおよび属性情報を処理するための実施形態を示す。

【図 21】 [ 0 0 4 1 ] 図 21 は、データパイロードを処理するための実施形態を示す。

【図 22】 [ 0 0 4 2 ] 図 22 は、統一および標準化された形式のゲノムアノテーションデータを選択的に展開するための実施形態を示す。

【図 23 A - 23 B】 [ 0 0 4 3 ] 図 23 A および図 23 B は、1 つまたは複数の実施形態に係る、ファイル形式をフォルダ階層に変換するための方法を示す。

40

【発明を実施するための形態】

【 0 0 2 1 】

[ 0 0 4 4 ] 図面は単に概略的なものであり、縮尺通りに描かれていないことを理解されたい。また、各図面を通して、同じまたは類似の部分を示すために同じ参照番号が使用されることを理解されたい。

【 0 0 2 2 】

[ 0 0 4 5 ] 記載および図面は、様々な例示的实施形態の原理を示す。したがって、当業者は、本明細書に明示的に記載または図示されていないものの、本発明の原理を体現し、本発明の範囲内に含まれる様々な構成を考え出すことができることが理解されよう。さらに、本明細書に記載のすべての例は、主に、本発明の原理、および発明者らが技術の発

50

展のために提供する概念を読者が理解するのを助けるという教育的な目的のためのものであることを明確に意図しており、そのような具体的に述べられた例および条件に限定されないものとして解釈されるべきである。また、本明細書で使用される「または」や「もしくは」という用語は、特に明記しない限り（例えば、「さもなければ」または「または代わりに」）、非排他的選言（すなわち、および/または）を指す。また、様々な例示的な実施形態は、必ずしも相互に排他的ではなく、一部の例示的な実施形態を1つまたは複数の他の例示的な実施形態と組み合わせ、新しい例示的な実施形態が形成され得る。「第1」、「第2」、「第3」などの記述語は、議論される要素の順序を限定することを意図したものではなく、ある要素を次の要素と区別するために使用され、一般的に交換可能である。最大値または最小値などの値は事前に決定されてもよく、用途に基づいて異なる値に設定されてもよい。

10

**【0023】**

[0046] 例示的な実施形態は、情報、例えば、これらに限定されないが、表形式データおよび区切られたテキストファイルなどの圧縮を管理するためのシステムおよび方法を説明する。圧縮される情報はゲノム情報（例えば、アノテーション、シーケンシング、または他の形態のゲノム関連データ）を含み得る。他の実施形態では、ゲノム情報に關係のない別のタイプの情報が圧縮のために管理されてもよい。

**【0024】**

[0047] ゲノム用途では、システムおよび方法は、ゲノムデータの圧縮および展開の標準化を管理し得る。これは、相互運用性に関連する問題、およびデータを視覚化できるようにするために頻繁に形式を変換する必要性を低減または排除する形式に準拠するようにデータを処理することを含む。少なくとも一部の用途では、これらのフィーチャにより、特定のタイプのデータのための圧縮アルゴリズムの改善および最適化が達成され得る。

20

**【0025】**

[0048] 1つまたは複数の実施形態によれば、システムおよび方法は、データを処理して、データの異なるフィールド、属性、および/またはセクションの選択的圧縮を可能にし得る。例えば、一部の装形態では、データのフィールド、属性、またはセクションはすべて（例えば、異なる統計的特性を有するものを含む）、すべてが同じ圧縮アルゴリズムを使用して圧縮されない可能性がある。そうではなく、システムおよび方法は、データの1つまたは複数のセクションを選択的に圧縮し、他のセクションは圧縮しないように実装されてもよく、または、異なる圧縮技術を使用して異なるセクションを選択的に圧縮するように実装されてもよい。選択される圧縮技術は、例えば、複数のセクションのうちの対応するそれぞれのセクション内の異なるタイプのデータにとって有益または最適であり得る。

30

**【0026】**

[0049] この選択的圧縮手法は、ひいては、選択的展開を可能にし得る。例えば、データを選択的に圧縮することにより、データファイル全体を展開する必要がない。そうではなく、少なくとも1つの装形態において、関心のあるデータの1つまたは複数のセクションを、関心のない圧縮されたデータの他のセクションを展開することなく、選択的に展開することができる。これにより効率が向上し、研究中の負担やその他の遅延を回避する能力が向上する。

40

**【0027】**

[0050] システムおよび方法によって実行される処理はまた、少なくとも1つの実施形態が、例えば、異なるデータセクションを選択的に適用される多数の標準またはカスタマイズされた圧縮アルゴリズムを使用して、すべてのタイプのデータに一般的に適用可能となることを可能にし得る。システムおよび方法はまた、（例えば、異なる方法で）圧縮されたデータセクションの選択的暗号化を実行し得る。システムおよび方法はまた、複数のアノテーションデータセットを互いに、およびシーケンシングデータとリンクするように実装されてもよい。システムおよび方法はまた、メタデータ、リンケージ、および/または属性固有インデックス付与を含むようにデータを処理し得る。

50

## 【 0 0 2 8 】

[ 0 0 5 1 ] 図 1 は、情報の圧縮を管理するためのシステムの実施形態を示す。システムは、少なくとも部分的に、例えば、実験室のワークステーション、研究施設、データセンター、大学もしくは企業環境、または、意図される用途に関連してデータが処理される別の場所の実装され得る。別の実施形態では、システムは、1つのコンピュータおよび/もしくは他のデバイス、または、本明細書に記載の処理動作を実行するために互いに通信する複数の接続もしくはネットワーク化されたコンピュータおよび/もしくは他のデバイス上に実装され得る。

## 【 0 0 2 9 】

[ 0 0 5 2 ] 図 1 を参照して、システムは、プロセッサ 1 1 0、メモリ 1 2 0、および記憶領域 1 3 0 を含む。プロセッサ 1 1 0 は、情報の圧縮を管理することを含む動作を実行するためにメモリ 1 2 0 内に保存された命令を実行する。メモリ 1 2 0 は、本明細書に記載の動作を実行するようにプロセッサを制御するための命令を保存する、任意の形式の非一時的コンピュータ可読媒体であり得る。情報はデータソース 1 4 0 から受け取られてもよい。データソースは、例えば、処理される情報のタイプによって異なり得る。例えば、一実装形態では、データソース 1 4 0 は、ゲノムアノテーションデータを生成するためのソフトウェアベースのツールであり得る。別の実装形態では、データソースは、別のタイプのゲノム関連データ、または、場合によってはゲノムデータに関連しないデータを提供してもよい。

## 【 0 0 3 0 】

[ 0 0 5 3 ] プロセッサ 1 1 0 はデータソースにローカルに接続されてもよく、またはネットワークを介してデータソースに接続されてもよい。後者の場合、プロセッサとデータソースとの間にネットワーク接続を確立することができれば、システムが処理する情報を世界中のどこからでも取得できる。ネットワーク接続は、ローカル接続、インターネット接続、仮想プライベートネットワーク接続、クラウドコンピューティングもしくはクラウドベースのストレージデバイスへの接続、および/または別のタイプの接続であり得る。

## 【 0 0 3 1 】

[ 0 0 5 4 ] 記憶領域 1 3 0 は、プロセッサ 1 1 0 の処理された結果を保存し得る。後により詳細に説明されるように、処理された結果は、複数のフィールド、属性、セクション、または、選択的に圧縮および/もしくはその他の形で処理され得る他の部分を含む、統一および標準化されたファイル形式に準拠する表形式データまたは区切られたテキストファイルを含む可能性がある。記憶領域は、データベース、アーカイブ、ストレージエリアネットワーク、クラウドベースのストレージシステム、プロセッサのコンピュータ、システム、もしくはデバイスに含まれるもしくは結合されたメモリ、または別のタイプの記憶領域であるか、またはこれらに含まれ得る。

## 【 0 0 3 2 】

[ 0 0 5 5 ] 図 1 のシステムは、統一および標準化されたファイル形式に処理された情報を圧縮するための1つまたは複数の追加フィーチャを含むか、または追加フィーチャに結合され得る。これらのフィーチャは、パーサ 1 5 0、圧縮マネージャ 1 6 0、およびアグリゲータ 1 7 0 を含む得る。

## 【 0 0 3 3 】

[ 0 0 5 6 ] パーサ 1 5 0 は、保存され、後に記憶領域 1 3 0 から取り出されたファイルを解析してもよい。例えば、プロセッサ 1 1 0 は、ユーザ入力に基づいて、または別のデバイスからデータのリクエストを受信し得る。プロセッサは、例えば、要求されたデータに対応する記憶領域 1 3 0 内に保存された1つまたは複数のファイルを決定し、決定された1つまたは複数のファイルを取り出してもよい。プロセッサ 1 1 0 は、例えば、ファイルのコンテンツをリクエストで指定された1つまたは複数の識別子にリンクするインデックスを使用して、どのファイルが要求されたデータに対応するかを決定してもよい。ファイルが見つかり、ファイルは記憶領域 1 3 0 から取り出されてパーサに送られる。

## 【 0 0 3 4 】

10

20

30

40

50

[ 0 0 5 7 ] パーサは様々な形でファイルを解析することができる。例えば、ファイルは、ファイル内のデータ（または他のタイプのコンテンツ）の異なる部分を含む異なるセクション、部分、属性、フィールドなどを区切る区切り文字を含み得る。これらの区切り文字をガイドとして使用して、パーサ 1 5 0 は、ファイルを、リクエストに対応する個別のセクション（またはチャンク）に解析することができる。一実施形態では、パーサは、圧縮のためにファイルがどのように解析されるかを示す所定の圧縮スキーマに従ってこれらの動作を実行してもよい。

【 0 0 3 5 】

[ 0 0 5 8 ] 圧縮マネージャ 1 6 0 は、ファイルが記憶領域 1 3 0 から取り出されると、統一された標準形式でファイルを受け取る。圧縮マネージャは複数のコンプレッサ 1 6 1 1、1 6 1 2、・・・、1 6 1 N を含み、ここで、N 2 である。各コンプレッサは、解析されたファイルの異なるセクションを圧縮するために、異なる圧縮アルゴリズムを実行してもよい。各コンプレッサによって適用される圧縮アルゴリズムは、ファイルの解析されたデータまたはコンテンツに基づいて事前に決定されていてもよい。例えば、効率の観点から、または他の理由で、特定のタイプの情報を圧縮するのに、ある圧縮アルゴリズムがより適している可能性がある。他のタイプの情報を圧縮するには他の圧縮アルゴリズムの方が適している可能性がある。圧縮スキーマは、ファイル内の特定のデータまたはコンテンツを圧縮するためにどの圧縮アルゴリズムを使用するかを制御し得る。（また、パーサが、圧縮スキーマから導出された命令に基づいて、ファイルの異なる解析されたセクションを特定のコンプレッサにルーティングしてもよい。）したがって、このようにして、圧縮マネージャ 1 6 0 は、ファイルの異なる解析された部分を選択的に圧縮することができる。

【 0 0 3 6 】

[ 0 0 5 9 ] アグリゲータ 1 7 0 は、圧縮ファイルの統一および標準化されたファイル形式に従って、または異なる形式に従って、複数の異なる圧縮されたセクションを圧縮ファイルに集約または結合する。圧縮後、選択的に圧縮ファイル（例えば、圧縮された表形式データファイルまたは圧縮された区切られたテキストファイル）は、記憶領域 1 3 0 に記憶されるか、異なる記憶領域に記憶されるか、所定の場所もしくはデバイスに送られるか、または、別のスキーム、例えば、システムのための所定の制御命令によって決定されるスキームに従ってルーティングされ得る。

【 0 0 3 7 】

統一された / 標準化されたデータ処理

[ 0 0 6 0 ] 1 つまたは複数の実施形態によれば、プロセッサ 1 1 0 は、ファイルの異なる部分を選択的に圧縮すること（したがって、後に選択的に展開すること）を可能にする、統一および / または標準化されたファイル形式に準拠するようにデータソースからの情報を処理する。処理は、1 つまたは複数の所定の機能をサポートすることを可能にするファイル形式でゲノムアノテーションデータを保存することを可能にするために実行され得る。これらの機能の例としては、高速クエリ、ランダムアクセス、複数の解像度（例えば、ズーム）、選択的暗号化、認証、アクセス制御、およびトレーサビリティなどが挙げられる。

【 0 0 3 8 】

[ 0 0 6 1 ] 一実施形態では、ゲノム用途データは、大幅な圧縮ゲインを可能にする形式に処理される。これは、例えば、データの複数の異なる属性を複数のセクション（例えば、複数のテーブル）に分離することによって実現され得る。その場合、複数のコンプレッサは、ゲノムアノテーションデータ内の分離された属性のそれぞれを圧縮するために使用され得る。一実施形態では、すべてのまたは一部のコンプレッサは、コンプレッサに割り当てられたファイルの特定の属性を圧縮するように設計された、特殊化またはカスタマイズされたコンプレッサであり得る。プロセッサ 1 1 0 によって実行される処理はまた、アノテーションに関連するシーケンシングデータのためのメタデータおよびリンケージ、ならびにシーケンシングデータに関連していても関連していなくてもよい他のアノテーション

10

20

30

40

50

ョンデータへのリンケージを生成し得る。メタデータおよび/またはリンケージは、ゲノム用途データと同じ統一および標準化された形式で保存されてもよく、これにより、シーケンシングデータのための既存のファイル形式（例えば、MPEG-G）とのシームレスな統合が可能になる。

【0039】

[0062] 図2は、統一および標準化されたファイル形式に従って情報の圧縮を管理するための方法を示す。方法は、例えば、図1のシステムまたは別のシステムによって、後述される動作を実行するように1つまたは複数のプロセッサを制御するためのメモリに保存された命令に従って、実行され得る。説明を目的として、方法が図1のシステムによって部分的にまたは全体として実行され、また、データがゲノムアノテーションデータであると仮定する。

10

【0040】

[0063] 図2を参照して、方法は、210において、データソース140からデータ（例えば、ゲノムアノテーションデータ）を受け取ることを含む。ゲノムアノテーションデータは、例えば、データを生成するために使用されるアノテーションツールによって生成される既存の形式で受け取られ得る。既存の形式の一例はMPEG-G形式であるが、別の実施形態では、アノテーションデータは異なる既存の形式で受け取られてもよい。データは、ローカル接続またはネットワーク接続を介して、データソースからプロセッサ110によって受け取られ得る。

【0041】

20

[0064] 220において、統一および標準化されたファイル形式への処理のためにゲノムアノテーションデータの複数の部分が取り出される。ゲノムアノテーションデータは既存の既知の形式に従って編成されているため、既知の形式の各部分で取り出されるコンテンツは、プロセッサ110に知られている特定のタイプのデータを有する。したがって、プロセッサ110を制御する命令は、1つまたは複数の実施形態に係る統一および標準化された形式で新しいファイルを生成するために、既知の形式の各部分において特定のタイプのデータを取り出すことができる。

【0042】

[0065] 230において、プロセッサ110は、例えば図3に示されるように、統一された標準ファイル形式300に準拠するようにゲノムアノテーションデータを処理するために、メモリ120内の命令を実行する。例えば、プロセッサは、ファイルヘッダセクション310、ファイル保護セクション320、ファイルメタデータセクション330、ファイルトレーサビリティセクション340、テーブル情報およびインデックスセクション350、1つまたは複数の圧縮パラメータを示すセクション360、ならびに1つまたは複数のテーブル370を生成するようにデータを処理することができる。図4のチャートには、これらのセクションおよび対応する説明とともに、後により詳細に説明される他の情報が示されている。

30

【0043】

[0066] ファイルヘッダセクション310は様々なタイプの識別情報を含む。例えば、ファイルヘッダセクションは、ファイル名フィールド311、ファイルタイプフィールド312、およびファイルバージョンフィールド313を含み得る（例えば、図5を参照されたい）。ファイル名フィールドは、例えば、アノテーションツールによって、またはユーザからの情報に基づいて生成された名前をファイルに付ける文字列を含み得る。ファイルタイプフィールドは、ファイルのタイプ（例えば、バリエーション、遺伝子発現など）を示す文字列を含み得る。ファイルバージョンフィールドは、ファイルの目的および/またはバージョンを識別する文字列（例えば、このファイルはアップデートを記録するためのものである）を含み得る。

40

【0044】

[0067] ファイル保護セクション320は、ファイルのアクセス制御ポリシーに関連する情報を含み、これは、例えば、ユーザ、管理者、データ所有者、またはデータへの

50

アクセスおよび/もしくはデータの配布を管理する他のエンティティによって決定され得る。アクセス制御ポリシーは、ファイル全体に対して同じレベルのアクセスを示してもよく、またはファイルの異なる部分に対して異なるレベルのアクセスを示してもよい。一実施形態では、ファイル保護セクションは、ペイロードサイズフィールド321およびペイロードフィールド322を含み得る(例えば、図6を参照されたい)。ペイロードサイズフィールドは、例えば、このフィールドまたはそれに関連する情報がスキップされる、または考慮に入れられることを可能にする保護情報ペイロードのサイズを示す1つまたは複数の整数値を含み得る。ペイロードフィールドは、ファイルの保護およびアクセス制御ポリシーに関する情報を含む。

【0045】

【0068】ファイルメタデータセクション330およびファイルトレサビリティセクション340は、後により詳細に説明されるような情報を含み得る。一実施形態では、ファイル保護情報、メタデータ、バージョンング、およびトレサビリティ情報は、例えば一般的な圧縮アルゴリズム、例えば7zipを使用して圧縮されてもよい。一部の実装形態では、(例えば、MPEG-G part 3で実行されるように)URI(uniform resource identifier)表記とともにこの情報を圧縮するためにJSON/XML/XACMLベースのスキーマが使用されてもよい。

【0046】

【0069】テーブル情報およびインデックスセクション350は、ファイル内のテーブルの数(n)を示すnTablesフィールド351を含み得る。一実施形態では、統一および標準化された形式は、アノテーションデータを複数のテーブル(例えば、n2)内に保存し得る。この場合、複数のテーブルの異なるテーブルは様々なタイプのアノテーション付きデータを保存し得る。あるケースでは、同じゲノムアノテーションデータを異なる解像度で保存するためにファイル内の異なるテーブルが使用されてもよい。これはテーブルの内容の一例として提供されているに過ぎない。他の実施形態では、テーブルは他のタイプのアノテーションデータを保存することができる。

【0047】

【0070】nTablesフィールド351に加えて、テーブル情報およびインデックスセクションは、TableID[i]セクション352、TableInfo[i]セクション353、およびByteOffset[i]セクション354を保存することができる。TableID[i]セクション352は、ファイル内のn個のテーブルのうちの対応するテーブルiの一意の識別子を示す1つまたは複数の整数値を保存し得る。TableInfo[i]セクション353は、n個のテーブルのうちのテーブルiのうちの対応するものの解像度を示す情報を保存し得る。ByteOffset[i]セクション354は、ファイル内のn個のテーブルのうちのテーブルiのうちのそれぞれのバイトオフセットを示す1つ以上の整数値を保存する。この形式を使用することで、例えばJSON/XMLのような形式でTableInfoフィールドを使用して、ファイル全体を読み取ることなく、ファイル内のテーブルのうちの1つまたは複数に関する基本情報(例えば、解像度レベル)が抽出され得る。同様に、圧縮ファイル内のテーブルのバイトオフセットは特定のテーブルに直接ジャンプするために利用可能であり得る。

【0048】

【0071】240において、受け取られた(ソース)ファイル内のゲノムデータが複数のチャンクまたはチャンクグループに分割される。複数のチャンクは同じ固定サイズを有してもよいし、または、例えば、チャンクが属性依存であるか、および/もしくはデータに関連する他のパラメータに基づくかに依存する可変サイズを有し得る。チャンクは、例えば、上記のようにパーサ150によって生成されてもよい。

【0049】

【0072】250において、分割されたチャンクまたはチャンクグループを圧縮するために異なるコンプレッサ(または圧縮アルゴリズム)が指定される。コンプレッサは、統一および標準化されたファイル外であってもよく、かつ/または、例えば、埋め込まれ

10

20

30

40

50

たコードによって、もしくは圧縮パラメータテーブル内に保存された選択可能なものによって内部的に開始されてもよい。

【 0 0 5 0 】

[ 0 0 7 3 ] したがって、圧縮パラメータセクション 3 6 0 は、ファイルのテーブル内に含まれる異なる情報（例えば、属性）を圧縮するために使用される異なるコンプレッサ（または圧縮アルゴリズム）の数を示すフィールド 3 6 1 を含み得る。この情報は、例えば、1 つまたは複数の整数値として表されてもよい。セクション 3 6 0 はまた、一意の識別子によってインデックス付与された異なるコンプレッサを列挙する情報を含むフィールド 3 6 2 を含み得る。これらはテーブル内で参照できるため、複数のテーブルにおいて、または複数の属性のために使用されるコンプレッサの繰り返しの記述が回避される。

10

【 0 0 5 1 】

[ 0 0 7 4 ] 一実施形態では、図 7 に示されるように、フィールド 3 6 2（または、圧縮パラメータセクション 3 6 0 内のフィールド 3 6 2 とは別の、もしくはフィールド 3 6 2 に階層的に関連するフィールド）は、CompressorID フィールド 3 6 3、nDependencies フィールド 3 6 4、CompressorNameList フィールド 3 6 5、および CompressorParametersList フィールド 3 6 6 を含み得る。

【 0 0 5 2 】

[ 0 0 7 5 ] CompressorID フィールド 3 6 3 は、複数のコンプレッサのうちの対応するコンプレッサの一意の識別子を示す情報（例えば、文字列）を含んでもよい。したがって、CompressorID フィールドは、ファイルのそれぞれの属性または他の部分の圧縮に使用されるコンプレッサを識別するために使用され得る。一実施形態では、一意の識別子は、対応するコンプレッサを指し示すためにテーブル内で使用され得る。

20

【 0 0 5 3 】

[ 0 0 7 6 ] nDependencies フィールド 3 6 4 は、属性の圧縮が他の属性に基づいて実行されることを示す情報（例えば、1 つまたは複数の整数値）を含み得る。例えば、システムプロセッサは、1 つまたは複数の他の属性を副次的情報として使用して 1 つの属性を圧縮するための情報を処理してもよい（循環依存関係がないことを条件として）。その場合、この情報は、統一および標準化されたファイル形式の nDependencies フィールド 3 6 4 内に含まれ得る。一実施形態では、変数 nDependencies の指標は、展開に使用される依存関係属性の数を示し得る（例えば、図 1 2 の属性情報構造には、対応する attributeID の例が示されている）。コンテキストベースの算術符号化などのコンプレッサは、副次的情報の組み込みを簡単にサポートし得る。他の属性からの副次的情報を組み込む別のメカニズムは、他の属性の値に基づいて現在の属性の値を並べ替えるか、または分割することである。これにより、類似する値をまとめて、より良い圧縮を提供することができる。パラメータは、リスト内の各コンプレッサによって使用される依存関係を記述する。

30

【 0 0 5 4 】

[ 0 0 7 7 ] CompressorNameList フィールド 3 6 5 は、コンプレッサ名のリストを示す情報（例えば、文字列）を含み得る。コンプレッサ名は、ファイルの外部の（例えば、圧縮モジュールまたはシステムプロセッサ 1 1 0 によって実装される）1 つまたは複数のコンプレッサ、（例えば、埋め込まれたコードまたは実行コードを介して）ファイルに埋め込まれた 1 つまたは複数のコンプレッサ、または両方を指す可能性がある。コンプレッサ名およびパラメータのうちの 1 つまたは複数は、標準コンプレッサに対応するか、または展開メカニズムを記述し得る。一実施形態では、コンプレッサ実行コードおよびパラメータのうちの 1 つまたは複数は、例えば CompressorName を「EMBEDDED」に設定することによって、CompressorParameters 内に示され得る。一部の実施形態では、同じ属性またはファイル部分を圧縮するために複数のコンプレッサが使用され得る。この場合、複数のコンプレッサが順番に適用さ

40

50

れ、統一および標準化されたファイル形式に含まれるリストに表示され得る。

【0055】

[0078] Compressor Parameters List フィールド 366 は、指定されたコンプレッサの全部または一部によって圧縮された情報の展開を実行するために使用されるべき1つまたは複数のパラメータを示す情報（例えば、リストまたは他の形式）を含んでもよい。

【0056】

[0079] 260において、方法は、本明細書に記載されるような対応するフィールド内に保存されるメタデータ、トレーサビリティ情報、リンケージ情報、インデックス、および保護情報を含む追加情報を生成するように、データの所有者もしくは管理者によって生成された、受け取られた（ソース）ファイルおよび/またはポリシー情報を処理することを含む。

10

【0057】

[0080] 270において、後に取り出すために、統一および標準化されたゲノムアノテーションデータファイルが保存される。複数の異なるデータチャンクが選択的に圧縮されているため、ファイル内の他の情報を展開する必要なく、圧縮ファイルの指定された部分のみを展開して、ファイル内の特定の情報にアクセスしてもよい。これは、任意の所与の時点で特に関心が持たれる可能性のあるゲノムデータの部分のみをターゲットにすることで、処理速度およびユーザの利便性を向上させる。

【0058】

[0081] 統一および標準化されたファイルにおいて指定される情報のすべてまたは一部は、テーブル、構造、インデックス内に示された関連付けられた情報、またはゲノムアノテーションデータの選択的圧縮および展開、ならびにアクセスを可能にする他のタイプの関連付けられた情報に対応し得る。関連付けられた情報内のテーブルの例については、後により詳細に説明する。

20

【0059】

テーブル

[0082] 各統一および標準化されたファイル内のテーブル370は、保存されている属性および/または他の情報に関して、互いに関連していても独立していてもよい。各テーブルは複数のセルを含み、各セルは1つまたは複数の属性を含み得る。各属性は特定のデータ型を有し、コンプレッサのうちの1つを使用して圧縮され得る。テーブル内のセルが複数の属性を含む場合、圧縮を向上させるために、および後に展開するときに属性への選択的アクセスを可能にするために、テーブルの各セル内の複数の属性に対応するデータは別々に圧縮され得る。ゲノム機能的アノテーションファイルでは、属性は、例えば、染色体、開始位置、終了位置、フィーチャID、およびフィーチャ名を含み得る。

30

【0060】

[0083] 圧縮は、異なる実施形態では異なる方法で実行され、例えば、属性ごとに、セルごとに、テーブルごとに、などで実行され得る。例えば、これらの目的のために、同じテーブル内のセルグループが他のセルグループとは別に圧縮されてもよい。一緒にグループ化されたセルは、同じデータに対応する属性を有し、他のセルグループは、同じまたは異なるデータに対応する異なる属性を有し得る。一実装形態では、各セル内の属性は異なる圧縮アルゴリズムを使用して圧縮されてもよく、または、1つまたは複数のセル内の属性は第1の圧縮アルゴリズムによって圧縮され、1つまたは複数の他のセル内の属性は第2の圧縮アルゴリズムを使用して圧縮され、1つまたは複数の他のセルの属性は第3の圧縮アルゴリズムを使用して圧縮されてもよい（以下同様）。

40

【0061】

[0084] ある追加の実装形態では、コンプレッサ1611、1612、・・・、161Nの全部または一部は、複数のセルのうちのそれぞれのセル内の異なるデータ/属性タイプに特化した異なる圧縮アルゴリズムを実行することができる。以下、統一および標準化されたファイル形式の例が複数のテーブルを含むものとして議論されるが、一実施形

50

態では、ファイル形式は1つのテーブルのみを含んでもよい。

【0062】

[0085]一部の例では、各セル内のコンテンツは属性として説明されるが、コンテンツは、他の実施形態では他のタイプのデータ（例えば、データのタイプ、属性、特性など）または複数のタイプのデータの組み合わせを含み得る。例えば、1つのファイル内のテーブル370は、すべてが同じゲノム配列を表すコンテンツまたはデータを保存し得るが、配列内の各テーブルは、異なる解像度でそのコンテンツまたはデータを表し得る。

【0063】

[0086]システムプロセッサ110は、ファイルテーブル370内の情報を様々な形で処理および保存することができる。例えば、処理は、一次元テーブルおよび/または多次元テーブルを生成するために実行され得る。所与の標準化されたファイル内のすべてのテーブルが同じタイプのものであってもよく、またはタイプの組み合わせを含んでもよく、例えば、一次元テーブルおよび多次元テーブルの両方を含み得る。ゲノムアノテーション用途では、一次元テーブルは、例えば、ゲノムアノテーションデータまたは定量的ブラウザトラックを含み、多次元テーブルは、例えば、1つまたは複数のサンプルに関するバリエーションデータおよび遺伝子発現データを含み得る。一実施形態では、一次元テーブルは複数の属性を含み得る。

10

【0064】

[0087]図8は、セルの一次元配列810内にデータを保存するテーブルの例を示す。配列内の各セル820は、属性1、属性2、および属性3とラベル付けされた複数の属性を含む。テーブル内の各セルは同じ数の属性を有しているが、一部の実施形態では、テーブルのセル内の属性の数は異なってもよい。テーブル構造の例、およびテーブル構造に含めるためにシステムプロセッサがコンテンツを処理し得る方法の例を後により詳細に説明する。

20

【0065】

[0088]図9は、セルの二次元配列910内にデータを保存する多次元テーブルの例を示す。テーブル内の各セル920は、データに対応する1つまたは複数の属性を含む。セル内の属性の数は同じであっても異なってもよい。テーブル構造の例、およびテーブル構造に含めるためにシステムプロセッサがコンテンツを処理し得る方法の例を後により詳細に説明する。

30

【0066】

[0089]一実施形態では、各ファイル内の多次元テーブル370のうちの1つまたは複数は、二次元配列910に加えて、次元固有属性を保存し得る。次元固有属性を保存するための例示的な構成が同様に図9に示されている。ここでは、配列910内に保存された属性とともに、テーブル370は、次元固有属性を2つの追加テーブル内に保存する。次元固有属性の第1の配列940は行属性を保存するセルを含む。次元固有属性の第2の配列950は列属性を保存するセルを含む。したがって、各配列940および950は、テーブル370内の追加の一次元テーブルの一部であると考えられることができる。これらのフィーチャを含むテーブル370の一例は、二次元配列テーブル910において複数のサンプルのバリエーションデータを表すことを企図し、サンプルの遺伝子型およびサンプルレベルの尤度が2D配列内の各セルの属性である。さらに、配列テーブル940は、バリエーション位置に対応する次元固有属性を保存してもよく、配列テーブル950は、サンプル名に対応する次元固有属性を保存してもよい。このテーブル構造の例、およびテーブル構造に含めるためにシステムプロセッサがコンテンツを処理し得る方法の例を後により詳細に説明する。

40

【0067】

[0090]図10は、システムプロセッサ110がどのようにしてゲノムアノテーションデータを特定のテーブル形式に処理し得るかの実施形態を示す。図11は、テーブルを生成するために処理される情報のためのテーブル形式の実施形態を示す。

【0068】

50

[ 0 0 9 1 ] 図 1 0 および図 1 1 に示されるように、表形式は、統一および標準化されたファイル形式のセクションに類似する多数のセクションを含む。例えば、表形式は、テーブルヘッダセクション 1 0 1 0、テーブル保護セクション 1 0 2 0、およびテーブルメタデータセクション 1 0 3 0 を含み得る。テーブルヘッダセクション 1 0 1 0 は、テーブルの名前を示す `Table ID` フィールド 1 1 0 5、ならびにペイロード(データ)タイプおよび/または他の情報を示す `Table Info` フィールド 1 1 1 0 を含み得る。テーブル保護セクション 1 0 2 0 は、テーブルのアクセス制御ポリシーを示す情報を含み得る。テーブルメタデータセクション 1 0 3 0 は、リンケージ、トレーサビリティ、および/または他の情報を含み得る。

【 0 0 6 9 】

10

[ 0 0 9 2 ] これらのフィーチャに加えて、テーブルは、要約統計量セクション 1 0 4 0、属性セクション 1 0 5 0、インデックスセクション 1 0 6 0、およびデータセクション 1 0 7 0 を含み得る。要約統計量セクション 1 0 4 0 は、平均、カウント、分布、および/または、1 つまたは複数のキー値に対応し得るテーブルに保存されたデータ/属性に関連する他の情報を示す情報を含み得る。統計は、例えば、関連性のある統計量への高速アクセスを可能にするために使用され得る。この情報は図 1 1 の参照番号 1 1 4 0 に対応する。

【 0 0 7 0 】

[ 0 0 9 3 ] 属性セクション 1 0 5 0 は、本明細書に記載されるような一次元および/または多次元配列テーブルのための様々な属性を保存し、固有次元属性を有しても有さなくてもよい。例えば、次元固有属性がテーブルに含まれている場合、セクション 1 0 5 0 は、各次元  $i$  ( $i = 1, \dots, N, N - 2$ ) のサイズ、名前、およびメタデータ、ならびに、テーブルに示されているように各次元に対応する次元固有属性を示すフィールド 1 1 5 0 を含み得る。二次元データの場合、`Symmetry Flag` フィールドはまた、次元配列が対称であるか否か、例えば、対称である `Hi-C` データであるか否かを示す情報を保存し得る。

20

【 0 0 7 1 】

[ 0 0 9 4 ] テーブルは、メイン配列(または二次元配列)内の属性の数を示す `n Attributes Main` フィールド 1 1 6 0 を含み、その後には各属性の属性情報が続く。属性情報は、本明細書でより詳細に説明されるコンテンツおよび構造を含み得る。また、フィールド 1 1 6 0 は、2 D 配列テーブルの属性/データのバイトオフセット情報を含み得る。異なるアルゴリズムと、属性に対応するバイトオフセットとに基づく属性の選択的圧縮は、他の属性、テーブルセクション、または統一および標準化されたファイルに含まれる他のデータを展開する必要なく、展開中に該属性(ゲノムアノテーションデータの特定の部分)に選択的にアクセスすることを可能にし得る。インデックスセクション 1 0 6 0 は、本明細書でより詳細に記載されるようなインデックスを保存し、データセクション 1 0 7 0 は、例えば、本明細書に記載されるようなチャンクまたはチャンクグループを含むデータを保存し得る。

30

【 0 0 7 2 】

[ 0 0 9 5 ] 図 1 2 は、統一および標準化されたファイルのテーブルまたは他のセクション内に含まれ得る属性情報、およびそれに関連付けられた構造を生成するために、システムプロセッサがどのようにしてゲノムアノテーションデータを処理するかの実施形態を示す。

40

【 0 0 7 3 】

[ 0 0 9 6 ] 図 1 2 を参照して、属性情報は、`Attribute Info Size` フィールド 1 2 0 5、`Attribute ID` フィールド 1 2 1 0、`Attribute Name` フィールド 1 2 1 5、`Attribute Metadata` 1 2 2 0 フィールド、および `Attribute Type` フィールド 1 2 2 5 を含む構造に処理され得る。`Attribute Info Size` フィールドは、この情報のサイズを示す情報を含み、これは、プロセッサがこの構造をスキップすることを可能にする。`Attribute`

50

IDフィールドは、属性のための一意の識別子を含み得る。Attribute Nameフィールドは、属性の名前を示す情報を含み、Attribute Metadataフィールドは、属性を1つまたは複数の他の属性とリンクまたはグループ化するための情報を含み得る。

【0074】

[0097] Attribute Typeフィールド1225は、属性が2つのタイプのうちの1つであることを識別する情報を含み得る。第1のタイプは基本型であり、属性が文字、文字列（null終端）、float、double、ブーリアン、異なるビット幅を有する符号付きおよび符号なし整数に対応するかを示す。第2のタイプは派生型であり、属性が可変長配列および固定長配列のどちらに対応するかを示す。

10

【0075】

[0098] これらのフィーチャに加えて、属性情報構造は、Default Valueフィールド1230、Summary Statisticsフィールド1235、およびCompressor IDフィールド1240を含み得る。Default Valueフィールドは、属性のほとんどの値が所定のデフォルト値と等しいか否かを示す。そうである場合、所定のタイプの符号化（例えば、スペース符号化）を使用して、少なくともファイル内の対応する属性が符号化され得る。要約統計量フィールドは、平均、カウント、分布、および/または、属性の高速分析を可能にし得る他の統計データを示す情報を含み得る。

【0076】

20

[0099] Compression IDフィールド1240は、属性に対して使用されるべき圧縮アルゴリズムのタイプを示す情報を含む。この情報は、例えば、属性を圧縮するために割り当てられた、コンプレッサ1611、1612、・・・、161Nのうちの1つを識別し得る。この情報は、システムプロセッサによって圧縮形式で取り出されたときに属性を復元するためにどの展開アルゴリズムが使用されるべきかを制御する。展開プロセス中にコンプレッサが副次的情報/コンテキストを使用する場合、対応する依存関係属性も属性情報内で指定され得る。多次元配列の場合、副次的情報は、多次元配列属性から、または次元固有属性から取得され得る。例えば、VCFファイルでは、（二次元配列テーブルの属性である）遺伝子型データの圧縮のための副次的情報としてバリエーション固有フィールド（例えば、次元固有属性）が使用されてもよい。一実施形態では、システムプロセッサは、展開に必要な追加データを変数Compressor Common Data内に含めるように属性情報を処理し得る。この追加データは、例えば、すべてのチャンクに共通であり得る。この追加データは、ゲノムアノテーションデータファイルの全部または一部から計算されたコードブック、辞書、または統計モデルを保存するのに有用である可能性がある。チャンクおよびチャンクに関連する処理については、後により詳細に説明する。

30

【0077】

コンプレッサ

[0100] 圧縮マネージャ160によって受け取られた情報（例えば、属性、データ、および他の情報）は、例えば、図1に示されるシステムの異なるコンプレッサ1611、1612、・・・、161Nを使用して選択的に圧縮され得る。これらのコンプレッサの非限定的なリストは以下を含む。

40

・ランレングス圧縮：この圧縮アルゴリズムは、例えば、同じ値を有する長いランのために使用され、ランの値および長さに基づいて置換動作を実行することを含み得る。

・差分符号化：この圧縮アルゴリズムは、例えば、数値のシーケンスを増加させるために使用され、連続する値の間の差に基づく置換動作を含み得る。

・辞書ベース/列挙：この圧縮アルゴリズムは、例えば、少数のオプションのセットに対応する値を有する属性のために使用され得る。この場合、アルゴリズムは、セット内のインデックスに基づいて置換動作を実行し、結果として得られた辞書をCompressor Common Dataフィールド内に保存し得る。

50

・スパーズ：この圧縮アルゴリズムは、例えば、例えば図12の属性情報において示されるように、デフォルト値とわずかに（例えば、所定の値未満）異なる属性を圧縮するために使用され得る。このタイプの圧縮は、属性を、非デフォルト値のための座標位置および値として表すことを含み得る。圧縮を改善するために、座標位置は各チャンク内でさらに差分符号化されてもよく、例えば、二次元スパーズ配列において、行インデックスを差分符号化し、各行内で列インデックスを差分符号化してもよい。

・可変長配列：この圧縮アルゴリズムは、例えば、可変長配列を値のストリームと長さのストリームとに分離するために使用されてもよい。

・トークン化：この圧縮アルゴリズムは、例えば、構造化文字列属性について使用され、属性を異なるタイプのトークンに分割し、前の値（例えば、前のトークン、差分、新しい値などと合致）に基づいて各トークンを符号化することを含み得る。

・汎用圧縮/エントロピー符号化法：これらのタイプの圧縮法は、例えば、gzip、bzip2、7-zip、適応算術符号化、および他の可逆データ圧縮（例えば、限定はされないが、BSCの可逆ブロックソーティング圧縮）を含む。

【0078】

[0101]上記の圧縮アルゴリズムのリストは限定的なリストであることを意図したものではない。一実施形態では、上記コンプレッサのうちの1つまたは複数に加えて、またはその代わりに、1つまたは複数の特殊化またはカスタマイズされた圧縮アルゴリズムが含まれていてもよい。そのような特殊化された圧縮アルゴリズムの一例はGTC（遺伝子型コンプレッサ）であり、バリエーションデータを圧縮するために使用され得る。そのような圧縮は、圧縮されたテーブルの行/列に含まれる情報への高速ランダムアクセスをサポートする可能性がある。場合によっては、例えば、より高速な選択的アクセスを可能にするために、ファイルの一部が圧縮されなくてもよい。

【0079】

[0102]コンプレッサ1611、1612、・・・、161Nのうちの1つまたは複数は複数のストリームを生成し得る。一例は、座標のストリームおよび値のストリームを生成するスパーズコンプレッサである。一実施形態では、これらのストリームは、例えば、適切に指定されたパラメータを使用して実装される異なるエントロピーコードを使用して圧縮されてもよい。例えば、以下を検討する。

```
CompressorNameList = [ 'sparse', 'gzip', '7-zip' ]
CompressorParameterList = [
  { "outStreams": [ "coordinate", "value" ] },
  { "inStreams": [ "coordinate" ] },
  { "inStreams": [ "value" ] }
]
```

【0080】

[0103]この場合、gzip圧縮が座標のストリームに適用され、7-zipが値のストリームに適用されてもよい（例えば、パラメータを表すためにJSONが使用され得る）。これにより、所定の、または最適な結果を生成し得る、データストリームごとのコンプレッサの適用が可能になる。ストリームが指定されていない場合、同じ圧縮法がすべての受け取られたストリームに適用されてもよい。

【0081】

[0104]上記したように、1つまたは複数の外部コンプレッサに加えて、またはその代わりに、（統一および標準化されたファイル形式に処理された）圧縮ファイルが1つまたは複数の埋め込まれたコンプレッサを含んでもよい。すなわち、システムプロセッサは、統一および標準化されたファイル形式内で実行可能な埋め込まれたコンプレッサを、好ましくは適切なセキュリティ保護とともに含むよう、ゲノムアノテーションデータを処理し得る。埋め込まれたコンプレッサの場合、対応する展開実行コードが、悪意のあるソフトウェアから保護するための発信元および真正性の証明としての1つまたは複数デジタル署名とともに、圧縮パラメータに含まれ得る。異なるプラットフォーム間での相互運用

10

20

40

50

性的ために、展開実行コードのために標準化された仮想マシンバイトコードが使用されてもよい。

【0082】

チャンクおよびインデックス付与

[0105] 図13は、統一および標準化されたファイル形式に構成された情報を処理するための方法の実施形態を示す。この方法では、圧縮ファイルは、展開および情報復元中に圧縮された情報に効率的にランダムアクセスできるように処理される。

【0083】

[0106] 図13を参照して、方法は、1310において、属性、データ、(1Dおよび/または2D) テーブル、および/または各テーブル370 (またはファイル全体) 内の他の情報をチャンクに分割することを含む。チャンクはすべて同じ固定サイズを有してもよく、または可変サイズを有してもよい。すべての属性に同じチャンクが使用されてもよいし、そうでなくてもよい。一実施形態では、効率的なランダムアクセスを可能にするために、各テーブル370内の属性は矩形のチャンクに分割される。

10

【0084】

[0107] 1320において、分割後、チャンクは、複数のコンプレッサのうちの異なるものによって選択的に圧縮される(例えば、システムプロセッサ110によって実行される命令によって指し示されるように)。

【0085】

[0108] 次に、1330において、プロセッサは、統一および標準化された形式の圧縮ファイルに含められる圧縮されたチャンクまたはチャンクグループごとに、少なくとも1つのインデックスを生成する。インデックスは、圧縮ファイル内の各特定のチャンクの位置を効率的に決定するための情報を含み得る。データをインデックスに処理することにより、関心のある1つまたは複数のチャンクとは無関係な、または関連性の無い他のチャンクを展開することなく、対応するチャンクを選択的に展開するだけで、データ内の任意の位置に高速アクセスすることが可能となり得る。特定の属性の値に基づく高速ランダムアクセスをサポートするために、システムプロセッサはまた、属性固有インデックスを生成するようにデータを処理してもよい。一実施形態では、各圧縮ファイルはまた、複数のチャンクにわたって展開を実行するためのコードブックまたは統計モデルの共有を可能にする情報を含み得る。

20

【0086】

[0109] 図14は、プロセッサによって生成されるインデックス構造の例を示す。この例は、以下で説明するように、フラグAttributeDependentChunksフィールドが偽である場合の一次元テーブルのためのインデックス構造に対応する。インデックスは第1のセクション1410および第2のセクション1420を含む。第1のセクション1410は複数のチャンクに対応する情報を含み、各チャンクは、開始インデックス情報1411、終了インデックス情報1412、およびバイトオフセット情報1413を含む。第2のセクション1420は、チャンク、または関連付けられた圧縮ファイルの他の部分に含まれる属性のための複数の追加インデックスを含む。追加インデックスはそれぞれ、インデックス付与された属性情報1421、インデックスタイプ情報1422、およびインデックスデータ情報1423を含み得る。インデックスファイル構造は、例えば、(例えば、図10に示されるような) 各テーブルのインデックスセクションに、または圧縮された統合および標準化されたファイルの選択されたテーブルに含まれ得る。

30

40

【0087】

[0110] 図15は、1つまたは複数の実施形態に係るインデックス構造の例を示し、図16は、1つまたは複数の実施形態に係る、図14のインデックス構造に対応し、そのより良い説明を提供し得るチャンク構造の例を示す。

【0088】

[0111] まず図16を参照して、チャンク構造はnChunksフィールド160

50

5 および `VariableSizeChunks` フィールド 1610 を含み得る。 `nChunks` フィールド 1605 は、特定の属性に対応するチャンクの数を示す情報（例えば、整数）を含む。 `VariableSizeChunks` フィールド 1610 は、 `nChunks` が同じ固定サイズを有するか、または異なるサイズを有するかを示すフラグを含む。固定サイズのチャンクは処理がより単純であるが（特に、多次元配列テーブルに含まれている場合）、可変サイズのチャンクは、例えば、データのスパース性が大きく変動するため、固定チャンクサイズの選択が最適でない場合に有用であり得る。場合によっては、可変サイズのチャンクがゲノムアノテーションデータの特定の属性（例えば、染色体、ゲノム位置など）により適している可能性があり、したがって、これらの属性に関してより高速なランダムアクセスを可能にし得る。

10

## 【0089】

[0112] チャンク構造が、可変サイズのチャンクが使用されていることを示す場合、各次元に沿った各チャンクの開始インデックス 1615 および終了インデックス 1620 がチャンク構造内に含まれ得る。この情報は図 14 の圧縮ファイルインデックス構造内にも示されている。固定サイズのチャンクの場合、各次元に沿ってチャンクサイズが示され得る。いずれの場合も（固定サイズまたは可変サイズ）、ランダムアクセスを可能にするために、チャンクごとにファイル内にバイトオフセット情報 1625 が含まれ得る。

## 【0090】

[0113] 図 17 は、特定の属性に関連するゲノムアノテーションデータファイルの一部分の例を示す。この例では、データは二次元配列に配置され、特定の `ByteOffset[j]` とともに、同じ固定サイズのチャンク 1710 に分割されている（例えば、 `ChunkSize(1)` は 5 であり、 `ChunkSize(2)` は 11 である）。図示されているこのデータ部分は計 15 個のチャンク (`nChunks = 15`) を含み、 `VariableSizeChunks` フィールドは固定サイズを指定するために偽値を示す。

20

## 【0091】

[0114] 図 18 は、特定の属性に関連するゲノムアノテーションデータファイルの一部分の例を示す。この例では、データは二次元配列に配置され、可変サイズを有するチャンク 1810 ~ 1840 に分割されている。チャンクサイズは可変であるため、開始インデックスおよび終了インデックスが、関連付けられた `ByteOffset` とともにブロックごとに提供される。例えば、第 1 の矩形チャンク（チャンク 1）の第 1 の方向については `Startindex[1][1] = 1` および `Endindex[1][1] = 17` が示され、第 1 の矩形チャンクの第 2 の方向については `Startindex[1][2] = 1` および `Endindex[1][2] = 17` が示される。チャンク 1 についても対応する `ByteOffset` 情報が提供される。同様に、第 2 の矩形チャンク（チャンク 2）の第 1 の方向については `Startindex[2][1] = 1` および `Endindex[2][1] = 11` が示され、第 2 の矩形チャンクの第 2 の方向については `Startindex[2][2] = 18` および `Endindex[2][2] = 31` が示される。図示されているこのデータ部分は計 4 個のチャンク (`nChunks = 4`) を含み、 `VariableSizeChunks` フィールドは可変サイズを指定するために真値を示す。残りのチャンク 3 および 4 も同様に指定され得る。該属性についての 4 つのチャンクすべてに対応する情報が、チャンク構造、インデックス構造、対応するテーブル、および/または、ゲノムアノテーションデータの統一および標準化されたファイルに関連する他の情報に組み込まれ得る。

30

40

## 【0092】

[0115] 再び図 15 を参照して、インデックス構造は、 `AttributeDependentChunks` フィールド 1510 を含む情報を含む。このフィールドは、チャンクサイズが 1 つまたは複数の属性に依存しているか、または、ファイルに関連付けられているすべての属性について同じチャンク化（例えば、同じチャンクサイズ）が使用されているかを示すフラグを含み得る。フラグが第 1 の値を有する場合、1515 において、すべての属性について同じチャンク化が使用される。すべての属性について同じチャンク

50

ク化を使用することで、インデックス構造のサイズが大幅に削減され、例えば、ほとんどの場合においてチャンク内のすべての属性が照会される場合に有用であり得る。

【0093】

[0116] 図19は、AttributeDependentChunksフラグが第1の値(例えば、偽)を有する場合における、一次元のケースのためのデータペイロード構造に含まれる情報の例を示す。この例では、データ1910は複数のチャンク1920(例えば、nChunks=5)に分割され、各チャンクは1つまたは複数の属性1930に対応する。このケースでは、チャンク1は3つの属性に対応する。また、各属性はペイロードサイズおよびペイロードを示す情報を含む。

【0094】

[0117] 再び図15を参照して、フラグが第2の値を有する場合、1520において、図16のチャートに示されるようにして属性依存チャンク化が使用される。属性依存チャンク化は、例えば、圧縮およびランダムアクセスに関する互いに異なる属性の最適なチャンクサイズが大幅に異なる場合(例えば、所定の値を上回る場合)に有用であり得る。例えば、一部の属性がスパース(疎)で他の属性がデンス(密)である場合、同じチャンクサイズを使用すると、圧縮が最適化されない可能性がある。属性依存チャンク化は、単一の属性のすべてのチャンクが頻繁に照会される場合、例えば、所定の頻度値を上回る頻度で照会される場合にも有用であり得る。

【0095】

[0118] 図20は、AttributeDependentChunksフラグが第2の値(例えば、真)を有する場合における、一次元のケースのためのデータペイロード構造の例を示す。この例では、データ2010は1つまたは複数の属性2020に対応する。各属性は、データを分割して得られる1つまたは複数のチャンク2030を含むか、またはそれに対応する。また、各チャンクはペイロードサイズおよびペイロードを示す情報を含む。図15および図20の比較から明らかなように、一実施形態によれば、チャンクおよび属性の編成は動作モードに依存し、例えば、属性依存チャンク化が実行されるか否かに依存する。

【0096】

[0119] 一実施形態では、選択的展開プロセス中のランダムアクセスは、行番号および/または列番号を使用せずに実行されてもよい。この実施形態は、例えば、特定の属性に関してランダムアクセスが実行される用途のために、例えば、ゲノム位置に関してランダムアクセスが実行される用途のために実行され得る。

【0097】

[0120] これらの場合、図15のインデックス構造は1つまたは複数の属性固有インデックス1525を含み得る。例えば、属性依存チャンク化が実行される場合、nAdditionalIndexesが使用され得る(n=1)。行/列ごとではなく、属性ごとのより高速な照会を実行するために、インデックス構造は1つまたは複数の属性インデックスを含み得る。例えば、ゲノムアノテーションデータの場合、染色体および位置の属性のためのインデックスがシステムプロセッサ110によって生成され、図15のインデックス構造に含められ得る。システムプロセッサによって照会が実行されると、クエリの染色体および位置属性にマッチするチャンクまたはチャンク番号が選択的展開のために返され得る。他の実施形態では異なる属性および/または異なる数の属性が使用されてもよい。

【0098】

[0121] インデックス構造はまた、nAdditionalインデックスに関連する他の情報を含み得る。例えば、属性固有インデックスごとに、インデックス構造は、AttributeIDsIndexedフィールド1530、IndexTypeフィールド1535、IndexSizeフィールド1540、およびIndexDataフィールド1545を含み得る。AttributeIDsIndexedフィールド1530は、インデックスが付与された1つまたは複数の属性(例えば、染色体、ゲノム位置な

10

20

30

40

50

ど)のリストを含み得る。

【0099】

【0122】IndexTypeフィールド1535は、AttributeIDsIndexedフィールド内にリストされる各属性について実行されたインデックス付与のタイプを示す情報を含み得る。一実施形態では、1つまたは複数の属性のインデックス付与のタイプは、標準セットとは異なる可能性がある。例えば、染色体属性およびゲノム位置属性のそれぞれについて、ならびに範囲クエリについてRツリーまたはCSIインデックスが使用されてもよい。データベースタイプのクエリについてBツリーインデックスが使用されてもよい。ゲノム範囲インデックス付与は、各チャンクの左端および右端の座標を保存でき、これは、照会された範囲と重複するチャンクを高速で識別することを可能にする。同様に、Bツリーインデックスは、属性値から、値およびチャンク内の値の位置を含むチャンクへのマップを保存できる。

10

【0100】

【0123】IndexSizeフィールド1540は、リストされた属性ごとにシステムプロセッサによって生成されたインデックスのサイズを示す情報を含み得る。

【0101】

【0124】IndexDataフィールド1545は、実際のインデックス付与データを所定の(例えば、バイナリ形式)で含むことができる。形式のタイプは、例えば、属性について生成されたインデックスのタイプに依存し得る。

【0102】

【0125】図21は、(例えば、図19および/または図20のデータ内に示されるような)ペイロード情報を生成するためにシステムプロセッサ110によって実行され得る例示的なロジックを示す。このロジックは、AttributeDependentChunksフィールド内の値に基づいて実装される。より具体的には、図19と図20との比較から明らかなように、ペイロードの構造はAttributeDependentChunksフィールドの値によって異なる。ペイロードのデータのチャンクが属性に依存しない場合、システムプロセッサは、上記のような図19に示される構造を生成するためのロジックを実装する。ペイロードのデータのチャンクが属性に依存する場合、システムプロセッサは、上記のような図20に示される構造を生成するための論理を実装する。

20

【0103】

【0126】図22は、インデックス構造(例えば、限定はされないが、図15のインデックス構造)に基づいてアノテーション付きゲノムデータを選択的に識別および展開するためにシステムプロセッサ110によって実行され得る照会方法の一実施形態を示す。照会は、次のように実装されるルックアップ動作に基づいて実行され得る。まず、2210において、システムプロセッサ110は、関心のある1つまたは複数の属性を示す情報を含むユーザクエリを受け取る。2220において、システムプロセッサはクエリを解析し、クエリ内の1つまたは複数の属性を識別する。クエリ内において、1つまたは複数の属性は様々な形で指定され、例えば、「abcd」などの文字列、またはそれぞれの属性を識別する一連の番号のうち特定の属性番号などの形で指定され得る。2230において、システムプロセッサは、属性固有インデックスを検索して、クエリ内で識別される属性に対応する1つもしくは複数のチャンクまたは1つもしくは複数のチャンク番号を決定する。2240において、システムプロセッサは、ゲノムアノテーションデータを検索して、クエリ属性とマッチする1つもしくは複数のチャンクまたは1つもしくは複数のチャンク番号を取り出す(または他の方法でそのようなチャンクまたはチャンク番号へのアクセスを取得する)。チャンクインデックスを使用してチャンクを復元することで、条件に合致する値がフィルタリングされ得る(チャンクには合致しない値も含まれる可能性があるため)。

30

40

【0104】

【0127】2250において、復元されたチャンクが選択的に展開され得る(例えば、指定された属性に関係のない他のチャンクを展開することなく展開される)。複数のデ

50

ータチャンクが先に互いに独立して選択的に圧縮されていたため、選択的展開が可能になる。展開は、図12の属性情報構造、および/またはインデックス、テーブル、もしくは本明細書に記載される他の構造において示されるアルゴリズムに基づいて実行され得る。チャンクがグローバル圧縮された場合、グローバル圧縮データは、例えば、属性情報構造内に示されるCompressorCommonDataメカニズムを使用して、チャンク間で共有され得る。一実施形態では、対称二次元配列テーブルの場合（例えば、図11のテーブルのSymmetryFlagが真である場合）、チャンクは一部、例えば下の三角部分および対角線のみをカバーしてもよい。この場合、展開プロセスは、対応する下三角値で埋めることにより、上対角値に対処し得る。その他の場合、チャンクは重複することなく、インデックスの全範囲をカバーしてもよい。

10

## 【0105】

## 相互運用性

[0128] ゲノムアノテーションデータの統一および標準化されたファイル形式は様々な追加のフィーチャを有し得る。1つの追加のフィーチャは、この形式で生成されたファイルに関連するリンケージおよび相互運用性に関連する。一実施形態では、システムプロセッサは、ソースファイル（例えば、図1のデータソース140によって提供されるデータ）の任意の符号化またはフォーマットから独立するように、ファイル情報を本明細書に記載されるように統一および標準化された形式に処理することができる。

## 【0106】

[0129] 別の実施形態では、システムプロセッサは、ソースファイルの符号化および/またはフォーマットとリンクされるように、かつ相互運用可能になるように、ファイル情報を統一および標準化された形式に処理することができる。データソースファイルは、MP EG - G形式/符号化または別のタイプの形式または符号化であってもよい。

20

## 【0107】

[0130] リンクされた解釈可能な状態に処理されると、ファイルは、関連付けられた情報をデータセットに保存することにより、ソースファイル（MP EG - Gファイルなど）の一部として処理され得る。MP EG - Gファイルは、調査全体のデータを保存し、各データセットグループが異なる個人に対応する。各MP EG - Gデータセットグループは、さらに、複数の異なるシーケンシングランに対応する複数のデータセットに分割される。

30

## 【0108】

[0131] 単一の個人に対応するデータを保存するために、複数の異なるアノテーションファイルが複数の別個のデータセットとして組み込まれ、各データセットは単一のアノテーションファイルまたはシーケンシングデータを含む。別個のデータセットの例を以下に示す。

## データセットグループ（単一の個人）

- データセット1（シーケンシングデータ）
- データセット2（シーケンシングデータ）
- データセット3（バリエーションコールデータ）
- データセット4（遺伝子発現データ）

40

...

## 【0109】

[0132] 大規模な調査からアノテーションデータを収集する場合、データセットは次のように編成されてもよい。

## データセットグループ（大規模調査）

- データセット1（バリエーションコールデータ）
- アノテーションファイル（サンプル1）
- アノテーションファイル（サンプル2）

...

## データセット2（遺伝子発現データ）

50

アノテーションファイル ( サンプル 1 )  
 アノテーションファイル ( サンプル 2 )  
 . . .  
 . . .

## 【 0 1 1 0 】

[ 0 1 3 3 ] 一実施形態では、例えば、以下に示すようにして、圧縮および分析パフォーマンスを向上させるために複数の異なるアノテーションファイルがマージされてもよい。

データセットグループ ( 大規模調査 )

データセット 1 ( バリエーションコールデータ )  
 アノテーションファイル ( すべてのサンプル )  
 データセット 2 ( 遺伝子発現データ )  
 アノテーションファイル ( すべてのサンプル )  
 . . .

10

## 【 0 1 1 1 】

[ 0 1 3 4 ] この実装形態を実行するために、システムプロセッサは、データ型 ( シーケンシング / バリエーション / 遺伝子発現 / . . . )、データセット内のアノテーションファイルの数、およびこれらのファイルの各々のバイトオフセットをサポートするために、既存のデータセットヘッダー構造を追加フィールドで拡張し得る。コンプレッサがアノテーションファイル間またはデータセット間で共有される場合、コンプレッサのパラメータは、それぞれ、データセットレベルまたはデータセットグループレベルで保存され得る。一実施形態では、アノテーションファイルのうちの 1 つまたは複数は、コンプレッサ名「 P O I N T E R 」を有するコンプレッサ構造を含み、位置を保存する圧縮パラメータ、例えば、{ “ D a t a s e t G r o u p I d ” : 1 , “ D a t a s e t I d ” : 2 , “ C o m p r e s s o r I d ” : 5 } は、コンプレッサが、データセットグループ 1、データセット 2 の中の 5 番目のコンプレッサで指定されているとおりであることを示す。

20

## 【 0 1 1 2 】

リンケージ

[ 0 1 3 5 ] 上記のフィーチャに加えて、システムプロセッサは、異なるタイプのアノテーションデータと対応するシーケンシングデータとの間のリンケージを含む形式で統一および標準化されたファイルを生成するように、データを処理してもよい。一実施形態では、このリンケージは、ファイルに保存された、またはファイルに関連付けて保存されたメタデータに基づいて提供され得る。

30

## 【 0 1 1 3 】

[ 0 1 3 6 ] これは、システムプロセッサが、データセットグループ、またはシーケンシングデータもしくは関連するアノテーションデータを保存するデータセットを、例えば M P E G - G p a r t 3 で説明されている U R I ( u n i f o r m r e s o u r c e i d e n t i f i e r ) 表記を使用するか、または J S O N を使用することによって図 4 の F i l e M e t a d a t a フィールド、または T a b l e M e t a d a t a フィールド内で指定することによって達成され得る。例えば、シーケンシングデータセットへのリンケージを提供するために、次の J S O N が F i l e M e t a d a t a 内で使用されてもよい。

40

```
“ L i n k a g e s ” : [ {
  “ D a t a T y p e ”      : “ S e q u e n c i n g ” ,
  “ D a t a s e t G r o u p ” : 5 ,
  “ D a t a s e t ”      : 2
} ]
```

上記の例は単一のリンケージのみを示すが、別の実施形態では複数のリンケージが提供されてもよい。

## 【 0 1 1 4 】

[ 0 1 3 7 ] 追加で、または代わりに、システムプロセッサはテーブルレベルのリンケ

50

ージを生成してもよい。一実施形態では、システムプロセッサは、インデックスによってテーブルレベルのリンケージを生成することができる。この場合、例えば、あるテーブルのn番目の行(列)が別のテーブルのm番目の行(列)に対応し得る。このタイプのリンケージは、複数のアノテーションファイル/テーブルが同じ行/列を共有する場合、繰り返しを防ぐことができる(例えば、まだマージされておらず、同じバリエーションからなる複数のVCF)。同様に、このタイプのリンクは、サンプルに関連する情報が単一のテーブルに保存されており、VCFテーブルおよび遺伝子発現テーブルの両方がこれにリンクしている場合に有用であり得る。

#### 【0115】

[0138] 別の実施形態では、システムプロセッサは、インデックスによってテーブルレベルのリンケージを生成することができる。この場合、値を別のテーブル内の属性とマッチさせることによって特定の属性がリンクされ得る。例えば、遺伝子発現データは、遺伝子に関する詳細な情報を伴わずに遺伝子名を含み得る。遺伝子に関する詳細な情報は別のファイルで入手可能である。このようなリンケージの使用例としては、自己免疫疾患に対応し、ヒト6番染色体内の座標範囲を指定するMHC(主要組織適合遺伝子複合体)のすべての遺伝子の遺伝子発現データを要求するクエリが考えられる。このクエリに対処するために、ゲノム座標インデックスに基づいて遺伝子情報ファイルから座標範囲のための複数の遺伝子名が取得され、これらの名前を遺伝子発現ファイル内で照会することで必要なデータが取得され得る。以下の例はこれらのフィーチャに関連しており、具体的には、行(次元1)を別のテーブル(例えば、同じアノテーションファイル内のテーブルNo. 3)の行にリンクし得る。

```
“Linkages” : [ {
  “Type” : “byIndex” ,
  “DimensionInCurrentTable” : 1 ,
  “Table” : 3 ,
  “DimensionInLinkedTable” : 1
} ]
```

#### 【0116】

[0139] 属性値によって列(次元2)を別のテーブルの行にリンクする例(現在のテーブルの次元2の属性2が、ファイル2のデータセット4内のテーブル3の次元1の属性5にリンクされる)は、システムプロセッサによって次のようにして実行され得る。

```
“Linkages” : [ {
  “Type” : “byValue” ,
  “DimensionInCurrentTable” : 2 ,
  “AttributeInCurrentTable” : 4 ,
  “Dataset” : 4 ,
  “AnnotationFile” : 2 ,
  “Table” : 3 ,
  “DimensionInLinkedTable” : 1 ,
  “AttributeInLinkedTable” : 5 ,
} ]
```

#### 【0117】

[0140] メタデータ構造は任意の情報ストレージをサポートするため、システムプロセッサは、例えば標準化された形式を使用して、さらに3つ以上のテーブルのリンクにリンケージをさらに拡張し得る(例えば、テーブル3が、テーブル1で使用されている遺伝子IDをテーブル2内の遺伝子名に変換してもよい)。上記の例ではリンケージに特定のJSONベースの形式が使用されているが、別の形式が使用されてもよく、例えば、限定はされないが、XMLが使用されてもよい。

#### 【0118】

[0141] 一実施形態では、属性ベースのリンケージが実施されてもよい。メタデー

10

20

30

40

50

データベースのリンケージは高レベルのリンケージに有用である可能性があるが、場合によっては、行/列ごとのリンケージが有益である可能性がある。例えば、複数のサンプルを有するVCFファイルでは、属性SequencingDatasetGroupおよびSequencingDatasetを列属性に追加することによって、特定のサンプルに対応するシーケンシングデータがリンクされ得る。そのようなリンケージ属性では、デコンプレッサがリンケージ属性を通常の属性と区別できるよう、メタデータ内で「LinkageAttributeFlag」が真に設定されていてもよい。

【0119】

[0142] 場合によっては、システムプロセッサは、ゲノム領域に応じて、アノテーションデータセット間のマッピング動作を実行し得る。これは、例えば、各データセットに別々にインデックスを付けることで実現されてもよい。例えば、VCFファイル内のある領域に対応するシーケンシングデータを見つけるために、シーケンシングデータのマスターインデックステーブルにアクセスすることで、1つまたは複数の適切なアクセスユニットが決定されてもよい。異なるデータセットについて別々のインデックス付与を用いることにより、各データセットについて、所定の（例えば、最適な）チャンクサイズおよびその他のパラメータを選択できる。さらに、場合によっては、AUクラスが異なるため、バリエーションをAUに直接リンクできない可能性がある。同様に、複数のサンプルを有するVCFファイルでは、バリエーションは複数のデータセットにわたってアクセスユニットにマッピングされ、この情報の保存は大きなストレージを占める可能性がある。一実施形態では、シーケンシングデータのAUIdまたはバイトオフセットは、VCFファイル内に行属性として保存され、これは、現在のバリエーションに対応するアクセスユニットの高速な検索を可能にする。また、遺伝子に対してリスト型属性を使用することによって、遺伝子がバリエーションのリストにマッピングされてもよい。

【0120】

[0143] 一実施形態では、アクセス制御ポリシーは、XACMLなどの所定の形式を使用してい、ファイルレベル（例えば、図4のアノテーションファイル内）、テーブルレベル（例えば、図10および図11の情報内）、またはその両方で指定され得る。特定のユーザはすべてのデータにアクセスできる可能性があるが、他のユーザは粗い解像度のデータにしかアクセスできない可能性がある（異なる解像度が異なるテーブルに保存されていることを思い出されたい）。このタイプのポリシーはファイルレベルで指定され得る。テーブル内の属性に固有のポリシーが、例えばテーブルレベルで指定されてもよい。これは、属性のサブセットのみへのアクセス、または属性の値に基づく特定のチャンクのみへのアクセスを含み得る。別のタイプのポリシーでは、メタデータおよび情報へのアクセスは許可されるが、実際のデータへのアクセスは許可されない可能性がある。

【0121】

展開

[0144] 本明細書に記載される統一および標準化された形式のファイルの選択された部分、またはファイル全体を展開するためのプロセスは、1つまたは複数のタイプのクエリを実行することをまず含み得る。クエリおよび/または展開アルゴリズムは相互に排他的ではない可能性があり、一部の実施形態では組み合わせられてもよい。例えば、メタデータおよび特定の属性の両方が展開されてもよいし、または選択されたチャンクからの選択された属性が展開されてもよい。展開はシステムプロセッサによって、またはシステムプロセッサに結合された別の処理エンティティによって実行され得る。場合によっては、アクセス制御ポリシーによってこれらのクエリの一部が制限されてもよい。これらをサポートするためにアプリケーション・プログラミング・インターフェース（API）が使用されてもよい。そのようなAPIはMP EG - G part 3と類似していてもよく、または異なるタイプのAPIであってもよい。

【0122】

[0145] メタデータ/情報クエリ：これらのタイプのクエリは、要求されたテーブル（例えば、解像度レベル）、コンプレッサ、属性、および/またはチャンクに対応する

10

20

30

40

50

メタデータおよび情報のみを照会し得る。クエリを実行する際、まず、図4のアノテーションファイルの初めに記載されているようなトップレベルの情報に直接アクセスし得る。テーブル固有メタデータ/属性の詳細は、図4のテーブルのByteOffsetを使用して選択的にアクセスされ得る。

#### 【0123】

【0146】完全データ展開：このタイプの展開は、すべてのテーブルおよび属性を含むデータ全体の展開を含み得る。これは、上記したように、まず、ファイルのトップレベルのメタデータおよびテーブル情報を読み取ることによって実行され得る。次に、圧縮パラメータが展開マネージャにロードされる。そして、テーブルごとに、テーブル情報、次元、および属性が読み取られた後、インデックスが読み取られて、各次元に沿ったチャンクの位置が決定される。次に、各チャンクおよび各属性についてデータペイロードが（シリアルまたはパラレルで）選択的に展開される。別の属性を依存関係/コンテキストとして使用して属性が圧縮されている場合、最初にその別の属性を展開することによって展開が実行されてもよい。属性がCompressorCommonDataを使用する場合、チャンクを展開する前にこの情報がロードされてもよい。二次元対称配列の場合（例えば、上記のSymmetryFlagを参照されたい）、対角および下三角行列のみが展開され、対称性を利用して上三角部分が埋められてもよい。

10

#### 【0124】

【0147】1つのテーブルのみの展開：このタイプの展開は完全データ展開と類似しているが、要求されたテーブルにジャンプするために（例えば、図4の）アノテーションファイル内のByteOffsetフィールドが使用され、そのテーブルのみが展開される点で異なる。

20

#### 【0125】

【0148】テーブルの選択された属性の照会：この方法は1つのテーブルのみの展開と類似しているが、要求された属性に対応する情報のみが読み取られる点で異なる。図12の属性情報構造内のAttributeInfoSize変数に基づいて、他の属性はスキップされる。一実施形態では、図22のPayloadSize[i][j]情報に基づいて他の属性をスキップすることによって、要求された属性のみが展開される。属性依存チャンクが使用される場合、所与の属性のすべてのチャンクが、例えば図20のデータペイロード構造に示されるように、一緒に保存されてもよい。

30

#### 【0126】

【0149】配列内の選択されたインデックス範囲のみの照会：この方法は1つのテーブルのみの展開と類似しているが、インデックスがロードされ、チャンク化のタイプ（固定サイズ/可変サイズ）に依存して、要求された範囲と重複するチャンクが決定される点で異なる。また、図16のチャンク構造テーブル内のByteOffset情報を使用して、上記の決定されたチャンクのペイロードにジャンプしてもよい。場合によっては、属性依存チャンクが使用されておらず、所与のチャンクのすべての属性と一緒に保存されている場合、方法はより効率的に実施され得る。また、要求されたチャンクが展開され、重複するインデックスのみが返される。ある属性のコンプレッサがチャンク内で効率的なランダムアクセスを可能にする場合、このランダムアクセスは、展開速度をさらに高めるための基礎として利用され得る。これが起こる可能性のあるシナリオの例として、スパース配列またはGTCなどの遺伝子型のための特化コンプレッサが挙げられる。

40

#### 【0127】

【0150】特定の属性の値/範囲に基づく照会：この手法は、配列内の選択されたインデックス範囲のみの照会と類似しているが、（例えば、図15のインデックス構造内に示されている）追加の属性固有インデックスが対象のチャンクについて利用可能である場合、このインデックスを使用して該1つまたは複数のチャンクが決定されてもよい。そのようなインデックスが利用可能でない場合、すべてのチャンクについて対象の属性が展開された後、関連するチャンクが決定されてもよい。追加のインデックスを使用しない場合であっても、すべてのチャンクについて一部の属性のみが展開されるため、展開速度が向

50

上する可能性がある。一実施形態では、残りの属性は、関連するチャンクについてのみ展開されてもよい。

【0128】

フォルダ構造および編集

[0151] 上記のような、システムプロセッサによるゲノムアノテーションデータの統一および標準化されたファイル形式への処理は、効率、利便性、システム要件の低減、高速照会、およびデータアクセサビリティに関して多くの利点を提供する。特定の用途へのより良い適合のために、実施形態に追加の変形および調整が加えられてもよい。

【0129】

[0152] 例えば、データが単一のマシンに保存され、頻繁に編集される場合、データはファイルマネージャを使用してディレクトリ/フォルダ階層で保存されてもよい。この階層は、ファイル全体を上書きするのではなく、単一のチャンクおよび属性に対応するファイルのみを変更することによって、データの一部を操作することを容易にし得る。編集が完了し、データが送信されるとき、ファイルは単一のファイル形式に変換し直されてもよい。これは、データペイロードサイズに基づいてインデックスを再計算し、フォルダ階層をまとめて1つのファイルに戻すことで実行できる。

10

【0130】

[0153] ファイル形式（例えば、図23A）をフォルダ階層（例えば、図23B）に変換し、そして単一のテーブルに戻す方法の例が示されている。この例では、変換は、ファイル形式のテーブルセクション2310内の各テーブルをフォルダ階層のフォルダとして扱い、ファイル形式のデータセクション2320内のデータの各チャンクをフォルダ階層のフォルダとして扱うことによって実行される（例えば、Attribute Dependent Chunks フィールドに偽値が保存されていると仮定して）。フォルダ階層において、ブロック2331～2343はファイルに対応し、ブロック2321～2326はフォルダに対応する。フォルダ階層において、チャンクは既に個別のフォルダに保存されているため、インデックスは属性固有インデックスを保存するだけでよい。単純な例を図9に示す。

20

【0131】

[0154] フォルダ階層として考えると、提案される方式は、用語を次のように変更することによって、上記の関連特許出願で提案されている方式と大まかに関連付けることができる。「属性」：「1つまたは複数の行/列を有する領域」および「チャンク」：「ブロック」。

30

【0132】

[0155] 上記実施形態のうちの1つまたは複数がいかにして実装され得るかを説明するために、関連する機能を提供しつつ、様々なアノテーションデータを保存するための以下の2つの例、（1）バリエーションコールデータ、および（2）ゲノム機能的アノテーションデータについて議論する。

40

50

【表 1】

バリエントコールデータ

```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NC3136
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO
FORMAT NA00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2
GT:GQ:DP:HQ 0|0:48:1:51,51 1 0:48:8:51,51 1/1:43:5:.,.
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017
GT:GQ:DP:HQ 0|0:49:3:58,50 0 1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB
GT:GQ:DP:HQ 1|2:21:6:23,27 2 1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T
GT:GQ:DP:HQ 0|0:54:7:56,60 0 0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTCT G,GTACT 50 PASS NS=3;DP=9;AA=G
GT:GQ:DP 0/1:35:4 0/2:17:2 1/1:40:3
```

10

20

表 1 : VCFファイルの例 (IGSRより)

【 0 1 3 3 】

30

[ 0 1 5 6 ] 上記の表 1 は、5つのバリエントおよび3つのサンプルを含むVCFファイルのセクションを示す。VCFファイルはシステムプロセッサ110によって、データを保持するとともに、関連付けられた追加の機能を提供しつつ、統一および標準化されたファイル形式に処理され得る。そのようなファイルは以下のフィーチャを有し得る。

【 0 1 3 4 】

[ 0 1 5 7 ] メタデータ：(##で始まる)コメント行はFileMetadataの一部として保持されてもよい。これが、シーケンシングデータとともにMPEG-Gファイルの一部として保存される場合、メタデータも、このバリエントコールデータに対応するシーケンシングデータを含む、対応するデータセットグループを含み得る。

【 0 1 3 5 】

40

[ 0 1 5 8 ] トレーサビリティ：トレーサビリティ情報が、シーケンシングデータとともにMPEG-Gファイルの一部として保存される場合、トレーサビリティ情報は、RAWシーケンシングデータと、使用されるツールおよびそれらのバージョンとのURIから開始する、1つまたは複数のバリエントコールの生成のための1つまたは複数のコマンドを含み得る。トレーサビリティ情報は、再現可能な方法でファイルを検証するために使用されてもよい。

【 0 1 3 6 】

[ 0 1 5 9 ] テーブル：バリエントデータが単一の解像度で保存される場合、バリエントデータは、統合されたファイル形式で単一のテーブルに保存されてもよい(nDimensions = 2)。

50

## 【 0 1 3 7 】

【 0 1 6 0 】次元属性：第 1 の次元（バリエーション）については、CHROM、POS、ID、REF、ALT、QUAL、FILTER、およびINFOフィールドなどの複数の次元属性が存在し得る。コメントに記述されているように、INFOフィールドはNS、DP、AFなどの複数の属性に分割され得る。これらの属性のタイプもコメントフィールド内で言及されていてもよい。これらをグループ化するために属性メタデータが使用されてもよい（例えば、NS、DP、AFはグループINFOに属し得る）。デフォルト値は属性に依存し、例えば、FILTER属性の場合、デフォルト値は「PASS」に設定され得る。

## 【 0 1 3 8 】

【 0 1 6 1 】第 2 の次元（サンプル）については、サンプル名（例えば、NA00001）がオリジナルのVCFファイル内に存在する唯一の属性である可能性がある。シーケンシングデータへのリンケージをサポートするために、追加の属性、例えば、このサンプルに対応するシーケンシングデータを含むデータセットグループおよびデータセットが追加されてもよい。特定の数量、例えば、特定のバリエーションに対応するカウントや平均数量などへの高速アクセスをサポートするために、追加の次元属性が追加されてもよい。コメント内のINFO属性の記述がAttributeMetadataの一部として保存されてもよい。

## 【 0 1 3 9 】

【 0 1 6 2 】2Dテーブル属性：これらの属性は、それぞれが二次元配列であるGT、GQ、DPなどのFORMATフィールド内に記述されている。これらの属性のタイプもコメント内に記述されていてもよい。ほとんどのバリエーションが表現されていない場合、GT属性のデフォルト値は例えば0/0に設定されてもよい。コメント内のこれらの属性の記述がAttributeMetadataの一部として保存されてもよい。

## 【 0 1 4 0 】

【 0 1 6 3 】コンプレッサ：属性のためのコンプレッサは、属性のタイプおよび特性に基づいて選択されてもよい。例えば、CHROMは列挙ベースのスキームの後にgzipを使用して圧縮され、POSは差分符号化の後にgzipを使用して圧縮され得る。サンプル名（NA00001など）は、例えば、トークン化ベースの文字列コンプレッサを使用して効率的に圧縮され得る。一部のINFOフィールドは少数のバリエーションのためのみ存在する可能性があるため、スパース表現で符号化されてもよい。同様に、遺伝子型（GT）は、スパース表現、または遺伝子型専用のコンプレッサ（例えば、GTC）を用いて符号化されてもよい。

## 【 0 1 4 1 】

【 0 1 6 4 】特定の可変長属性の長さは、1つまたは複数の他の属性に依存し得る。例えば、AF（アレル頻度）属性の長さは、変異アレルの数に等しい可能性がある。そのような場合、コンプレッサのndependenciesは1に設定され、この依存関係を利用して圧縮が強化され得る。同様に、両者間の依存関係を利用して、他のFORMATフィールドの圧縮のための副次的情報としてGTフィールドの値が使用されてもよい。

## 【 0 1 4 2 】

【 0 1 6 5 】チャンク化およびインデックス付与：メイン2D配列のチャンク化はアクセスパターンに依存して実行されてもよい。ほとんどのアクセスが特定の領域内のバリエーションに対するものである場合、各チャンクはすべてのサンプルと少数のバリエーションとを含み得る（例えば、水平チャンク）。ほとんどのアクセスが特定のサンプルのすべてのバリエーションに対するものである場合、チャンクはすべてのバリエーションと少数のサンプルとを含み得る（例えば、垂直チャンク）。両方のタイプのクエリが頻繁に生じる場合、場合によっては、少数のバリエーションとサンプルとを含む矩形チャンクを使用する方がよい可能性がある。チャンクのサイズを大きくすることにより、ランダムアクセスのパフォーマンスが圧縮率とトレードオフされる可能性がある。

## 【 0 1 4 3 】

[ 0 1 6 6 ] ゲノム領域に基づくランダムアクセスの場合、表 2 に示されるように追加のインデックスが使用されてもよい (例えば、C S I インデックス付与に基づく)。

【表 2】

|                                       |                  |
|---------------------------------------|------------------|
| A t t r i b u t e I D s I n d e x e d | C H R O M, P O S |
| I n d e x T y p e                     | C S I            |
| I n d e x S i z e                     | インデックスのサイズ       |
| I n d e x D a t a                     | C S I インデックス構造   |

10

【 0 1 4 4 】

[ 0 1 6 7 ] C S I で行われるように実際のファイル位置を指定する代わりに、対象のゲノム領域と重複するチャンク I D のリストが返されてもよい (または示されてもよい)。そして、デフォルトインデックス構造から、ファイル内のこれらのチャンクの位置が決定され得る。I n d e l バリエーションまたは構造バリエーションが一般的である場合、バリエーションの S T A R T 位置および E N D 位置の両方に基づいて C S I インデックス付与が行われ得る。高速ランダムアクセス照会を可能にするために、より多くの属性にインデックスが付けられてもよい。例えば、F I L T E R 属性にインデックスを付けることで、F I L T E R = P A S S であるか否かに基づくバリエーションのより高速なフィルタリングが可能にされてもよい。

20

【 0 1 4 5 】

[ 0 1 6 8 ] 保護 : アクセス制御ポリシーは、例えばユースケースに応じて、様々な形態をとることができる。例えば、特定のユーザはすべてのデータにアクセスできるが、他のユーザは特定のゲノム領域内のバリエーションにのみアクセスできる場合がある (例えば、C H R O M および P O S によって指定される)。同様に、アクセスは特定のサンプルのみに制限されてもよい。あるケースでは、これは、チャンクがそれに従って選択されることを要する可能性がある。アクセス制御は属性レベルで課されてもよく、例えば、I N F O フィールドへのアクセスは許可されるが、個別のサンプルデータへのアクセスは許可されない。

【表 3】

ゲノム機能的アノテーションデータ (B E D)

30

```

browser position chr7:127471196-127495720
browser hide all
track name="ItemRGBDemo" description="Item RGB demonstration" visibility=2
itemRgb="Cn"
chr7 127471196 127472363 Pos1 0 - 127471196 127472363 255,0,0
chr7 127472363 127473530 Pos2 0 - 127472363 127473530 255,0,0
chr7 127473530 127474697 Pos3 0 - 127473530 127474697 255,0,0
chr7 127474697 127475864 Pos4 0 - 127474697 127475864 255,0,0
chr7 127475864 127477031 Neg1 0 - 127475864 127477031 0,0,255
chr7 127477031 127478198 Neg2 0 - 127477031 127478198 0,0,255
chr7 127478198 127479365 Neg3 0 - 127478198 127479365 0,0,255
chr7 127479365 127480532 Pos5 0 - 127479365 127480532 255,0,0
chr7 127480532 127481699 Neg4 0 - 127480532 127481699 0,0,255

```

40

表 3 : 簡単な B E D ファイルの例 (UCSC ゲノムブラウザ F A Q より)

【 0 1 4 6 】

[ 0 1 6 9 ] 上記の表 3 は、アノテーションデータとともに B E D ファイルのセクショ

50

ンを示す。システムプロセッサは、後述されるように、データを保持するとともに、追加の機能を提供しつつ、統一および標準化されたファイルに準拠するようにこのファイル内の情報を処理し得る。

【0147】

[0170]メタデータ：コメント行（例えば、最初の3つの行）はFileMetadataの一部として保持されてもよい。これが、シーケンシングデータとともにMP EG-Gファイルの一部として保存される場合、メタデータも、このアノテーションデータに対応するシーケンシングデータを含む、対応するデータセットグループを含み得る。

【0148】

[0171]テーブル：データを異なるスケールや解像度で表示できるようにするために、システムプロセッサによって、異なる解像度のための事前に計算された値とともに、複数のテーブルが生成されてもよい。例えば、TableInfoフィールドは、解像度を示すパラメータおよび他の情報を事前定義された形式で保存してもよい。これにより、ユーザはファイル全体を読み取ることなく、利用可能な解像度のリストを照会できる可能性がある。また、テーブルごとのByteOffset変数を使用することにより、所望の解像度に直接アクセスできる可能性がある。複数のテーブルのうちの1つまたは複数は、例えば、単一の次元を有し得る。

10

【0149】

[0172]属性：一実施形態では、各列は次のような属性として機能することができる：chrom（文字列）、chromStart（整数）、chromEnd（整数）、name（文字列）、score（整数）、strand（文字）、thickStart（整数）、thickEnd（整数）、itemRGB（長さ3の8ビット整数配列）。

20

【0150】

[0173]コンプレッサ：属性のためのコンプレッサは、属性のタイプおよび特性に基づいて選択されてもよい。例えば、chromは列挙ベースのスキームの後にgzipまたはランレングス圧縮を使用して圧縮されてもよく、chromStartおよびchromEndは、差分符号化の後にgzipを使用して圧縮されてもよい。thickStartおよびthickEndの値がchromStartおよびchromEndに近い場合、例えば、これらの値を副次的情報として使用することによって圧縮が改善される可能性がある。

30

【0151】

[0174]この例では、chromStartの値は前の行のchromEndの値と合致する。これを利用する1つの方法は、chromStartおよびchromEndを「長さ2の整数配列」タイプの単一の属性と見なすことである。これは、例えば、視覚化ツールがこの代替表現を理解する場合に実行され得る。

【0152】

[0175]チャンク化およびインデックス付与：ゲノム領域に基づくランダムアクセスの場合、表4に示されるように追加のインデックスが使用されてもよい（CSIインデックス付与に基づき）。CSIで行われるように実際のファイル位置を指定する代わりに対象のゲノム領域と重複するチャンクIDのリストが示されてもよい。そして、デフォルトインデックス構造から、ファイル内のこれらのチャンクの位置が決定され得る。

40

【表 4】

|                     |                                |
|---------------------|--------------------------------|
| AttributeIDsIndexed | chrom, chromStart,<br>chromEnd |
| IndexType           | CSI                            |
| IndexSize           | インデックスのサイズ                     |
| IndexData           | CSIインデックス構造                    |

10

## 【0153】

【0176】CSIで行われるように実際のファイル位置を指定する代わりに対象のゲノム領域と重複するチャンクIDのリストが示されてもよい。そして、デフォルトインデックス構造から、ファイル内のこれらのチャンクの位置が決定され得る。

## 【0154】

【0177】保護：また、アクセス制御ポリシーはユースケースに応じて様々な形態をとることができる。例えば、特定のユーザはすべてのデータにアクセスできる可能性があるが、他のユーザは粗い解像度のデータにしかアクセスできない可能性がある（例えば、異なる解像度が異なるテーブルに保存され得ることを思い出されたい）。同様に、アクセスは特定のゲノム領域のみに制限されてもよい。この場合、チャンクがそれに従って選択されてもよい。

20

## 【0155】

シングルセルRNAseq発現データ

【0178】シングルセルRNAseq発現データは、整数/float発現値のスパース二次元行列を含み得る。各行は遺伝子に対応し、各列は細胞を表すバーコードに対応する。この場合、発現値は二次元のスパース属性配列として保存され、一方、遺伝子に関連する情報は次元固有の行属性になり、バーコードに関連する情報は次元固有の列属性になる。スパース配列は座標ストリームと値ストリームに分割されてもよい。これらは別々に圧縮され、行座標および各行の列座標は圧縮前に差分符号化される。高速ランダムアクセスを可能にするために、チャンクあたり固定数の遺伝子（行）にデータがチャンク化されてもよい。最後に、遺伝子idに基づく選択的照会を実行するために追加のBツリーインデックスが使用されてもよい。このインデックスは、遺伝子idを、遺伝子idおよびチャンク内の遺伝子idの位置を含むチャンクにマッピングする。

30

## 【0156】

【0179】この手法をE18マウスの1万個の脳細胞からなるデータセットに適用した。約31,000個の遺伝子、6,800,000個のバーコード、および4,000万個の整数エントリがスパース配列に含まれていた。最終圧縮レイヤーとしてBSCを用いた提案される手法は、750MB（未圧縮）から67MBにサイズを縮小する。これは、列を属性に分割したり、座標の差分符号化を行うことなく、オリジナルのファイルにgzipまたはBSCを直接適用した場合の圧縮後サイズの2分の1未満である。

40

## 【0157】

【0180】本明細書に記載の方法、プロセス、および/または動作は、コンピュータ、プロセッサ、コントローラ、または他の信号処理デバイスによって実行されるコードもしくは命令によって実行され得る。コードまたは命令は、1つまたは複数の実施形態に従って、非一時的なコンピュータ可読媒体に保存されてもよい。方法（または、コンピュータ、プロセッサ、コントローラ、もしくは他の信号処理デバイスの動作）の基礎をなすアルゴリズムが詳細に説明されているので、方法の実施形態の動作を実装するためのコードまたは命令は、コンピュータ、プロセッサ、コントローラ、または他の信号処理デバイスを、本明細書に記載の方法を実行するための専用プロセッサに変換することができる。

## 【0158】

50

[ 0 1 8 1 ] 本明細書に開示される実施形態のプロセッサ、コンプレッサ、デコンプレッサ、マネージャ、セクタ、パーサ、ならびに他の情報生成、処理、および計算フィーチャは、例えば、ハードウェア、ソフトウェア、または両方を含み得る論理で実装され得る。少なくとも部分的にハードウェアで実装される場合、プロセッサ、コンプレッサ、デコンプレッサ、マネージャ、セクタ、パーサ、およびその他の情報生成、処理、および計算フィーチャは、例えば、様々な集積回路、例えば、限定はされないが、特定用途向け集積回路、フィールドプログラマブルゲートアレイ、論理ゲートの組み合わせ、システムオンチップ、マイクロプロセッサ、または別のタイプの処理もしくは制御回路のうちのいずれかであり得る。

【 0 1 5 9 】

[ 0 1 8 2 ] 少なくとも部分的にソフトウェアで実装される場合、プロセッサ、コンプレッサ、デコンプレッサ、マネージャ、セクタ、パーサ、およびその他の情報生成、処理、および計算フィーチャは、例えば、コンピュータ、プロセッサ、マイクロプロセッサ、コントローラ、または他の信号処理デバイスなどによって実行されるコードまたは命令を保存するためのメモリまたは他のストレージデバイスを含み得る。方法（または、コンピュータ、プロセッサ、マイクロプロセッサ、コントローラ、もしくは他の信号処理デバイスの動作）の基礎をなすアルゴリズムが詳細に説明されているので、方法の実施形態の動作を実装するためのコードまたは命令は、コンピュータ、プロセッサ、コントローラ、または他の信号処理デバイスを、本明細書に記載の方法を実行するための専用プロセッサに変換することができる。

【 0 1 6 0 】

[ 0 1 8 3 ] 上記の説明から、本発明の様々な例示的实施形態がハードウェアまたはファームウェアで実装され得ることは明らかであろう。また、様々な例示的实施形態は、機械可読記憶媒体上に記憶された命令として実装されてもよい。これらの命令は、本明細書において詳細に説明される動作を実行するために少なくとも1つのプロセッサによって読み取られ、実行され得る。機械可読記憶媒体は、パーソナルもしくはラップトップコンピュータ、サーバ、または他のコンピューティングデバイスなどの機械によって読み取り可能な形態で情報を保存するための任意のメカニズムを含み得る。したがって、機械可読記憶媒体は、リードオンリーメモリ（ROM）、ランダムアクセスメモリ（RAM）、磁気ディスク記憶媒体、光学記憶媒体、フラッシュメモリデバイス、および同様な記憶媒体を含むことができる。

【 0 1 6 1 】

[ 0 1 8 4 ] 当業者は、本明細書に示されるブロック図が、本発明の原理を具現化する例示的回路の概念図であることを理解するであろう。同様に、あらゆるフローチャート、フロー図、状態遷移図および疑似コードなどは、機械可読媒体内に実体的に表され、よって（コンピュータまたはプロセッサが明示的に示されているか否かに関わらず）コンピュータまたはプロセッサによって実行され得る様々なプロセスを表す。

【 0 1 6 2 】

[ 0 1 8 5 ] 様々な例示的实施形態が、それらの例示的实施形態の特定の例示的態様に特に言及して詳細に説明されているが、本発明は他の実施形態も可能であり、その細部は様々な明白な点において変更可能であることを理解されたい。1つまたは複数の実施形態を1つまたは複数の他の実施形態と組み合わせることで新たな実施形態が形成されてもよい。当業者には容易に明らかであるように、本発明の趣旨および範囲から逸脱することなく変形および変更を加えることが可能である。したがって、上記の開示、記載、および図面は説明のみを目的としており、請求項によってのみ定義される本発明をいかにようにも制限しない。

【 0 1 6 3 】

[ 0 1 8 6 ] 1つまたは複数の実施形態によれば、ゲノムアノテーションデータの選択的圧縮および展開を制御するためのシステムおよび方法は、複数の互換性のないファイル形式のうちの1つのファイル形式の情報を処理して、選択的圧縮および展開を制御する統

10

20

30

40

50

合されたファイルにすることを含む。高速クエリ、ランダムアクセス、複数の解像度（ズーム）、選択的暗号化、認証、アクセス制御、およびトレーサビリティなどの機能をサポートするように、アノテーションデータが処理され、情報が抽出さおよびさらに処理される。また、処理は、データの複数の異なる属性を分離し、これらに特化したコンプレッサを使用できるようにすることによって、大幅な圧縮ゲインを可能にするように実行されてもよい。アノテーションに関連するシーケンシングデータへのメタデータおよびリンクエッジ、ならびに同じ調査からの他のアノテーションデータへのリンクエッジをサポートするために追加の処理が実行されてもよい。これは、シーケンシングデータのための既存の M P E G - G ファイル形式とのシームレスな統合を可能にする。

【 0 1 6 4 】

10

[ 0 1 8 7 ] 1 つまたは複数の実施形態では、きめ細かいセキュリティ設定を可能にする、階層内の複数のレベルにおける保護（アクセス制御）情報を生成することために追加の処理が実行されてもよい。同様に、メタデータおよび属性は、複数の異なるタイプのアノテーションデータやシーケンシングデータセットを効果的にリンクすることを可能にする。処理によって生成されたファイルは、スタンドアロンファイルとして、または M P E G - G ファイルの一部として使用され得る。また、そのようなファイル、特にゲノムアノテーションデータのためのそのようなファイルの生成は、対象の属性に適した圧縮技術を組み込むことにより、様々なデータ型において最高水準の圧縮パフォーマンスを達成するのに十分な柔軟性を提供する。

20

30

40

50

【図面】

【図 1】

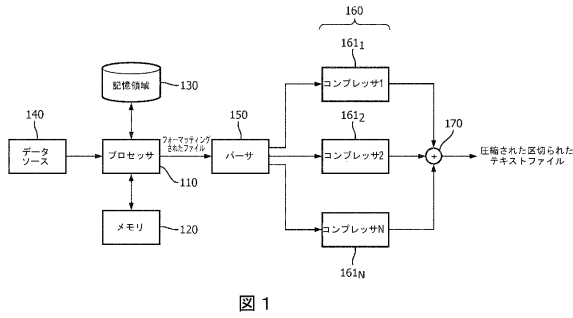


図 1

【図 2】

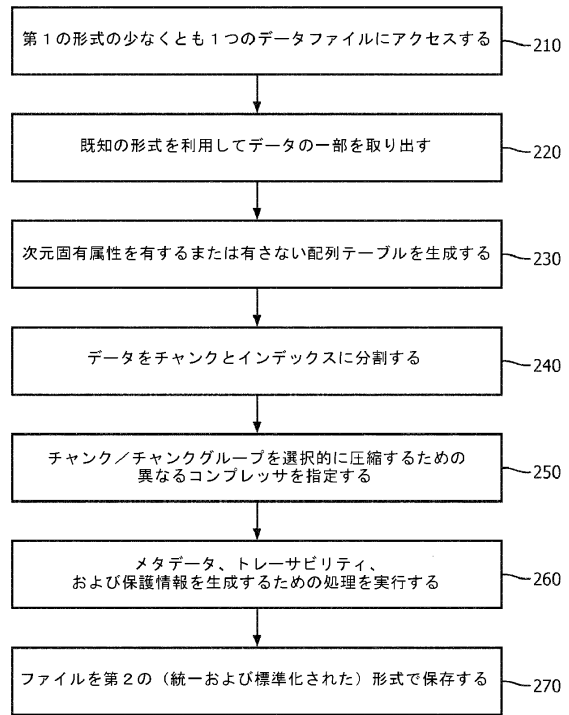


図 2

【図 3】

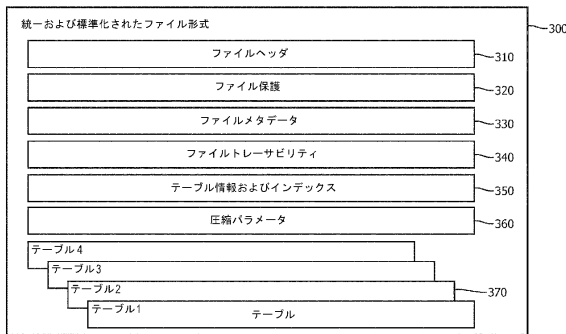


図 3

【図 4】

アノテーションファイル

| フィールド | 簡単な説明                    | タイプ                              |
|-------|--------------------------|----------------------------------|
| 310   | File header              | ファイルヘッダ(\$T2)                    |
| 320   | File protection info     | アクセス制御ポリシー                       |
| 330   | File metadata            | メタデータ/リンケージ                      |
| 340   | File traceability info   | データの生成に使用されるコマンド                 |
| 351   | nTables                  | ファイルに保存されているテーブルの数 (例えば複数の解像度)   |
| 352   | For i in 1...nTables:    |                                  |
| 353   | Table ID [i]             | 一意のテーブル識別子                       |
| 354   | TableInfo [i]            | テーブル情報 (例えば解像度)                  |
| 361   | ByteOffset [i]           | ファイル内のテーブル i のバイトオフセット           |
| 362   | nCompressors             | 後に保存される様々な属性において使用される異なるコンプレッサの数 |
|       | For i in nCompressors:   |                                  |
|       | Compressor [i]           | コンプレッサ情報(\$T4)                   |
|       | For i in 1...nTables [i] |                                  |
|       | Table [i]                | テーブル(\$T5)                       |

図 4

10

20

30

40

50

【図 5】

ファイルヘッダ 310

| フィールド            | 簡単な説明                | タイプ |
|------------------|----------------------|-----|
| 311 File name    |                      | 文字列 |
| 312 File type    | 例えば「バリエント」、「遺伝子発現」など | 文字列 |
| 313 File version | 更新を記録するため            | 文字列 |

図 5

【図 6】

一般的情報構造 320

| フィールド            | 簡単な説明               | タイプ |
|------------------|---------------------|-----|
| 321 Payload size | これをスキップすることを可能にするため | 整数  |
| 322 Payload      | 所定のコンプレッサ（例えば7z）で圧縮 | バイト |

図 6

【図 7】

コンプレッサ情報構造 360

| フィールド                        | 簡単な説明                        | タイプ        |
|------------------------------|------------------------------|------------|
| 363 CompressorID             | 一意のコンプレッサ識別子                 | 文字列        |
| 364 nDependencies            | 他の属性の値に基づいて属性を圧縮することを可能にするため | 整数         |
| 365 CompressorNameList       | コンプレッサ名のリスト（または「埋め込み」）       | リスト（文字列）   |
| 366 CompressorParametersList | 展開に必要なパラメータ                  | リスト（タプル情報） |

図 7

【図 8】

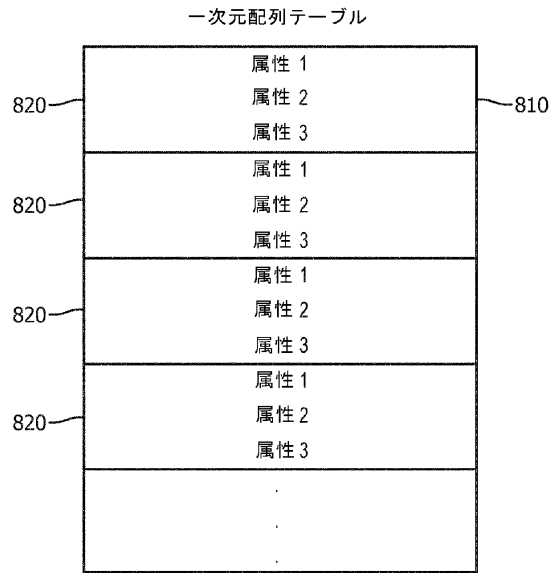


図 8

10

20

【図 9】

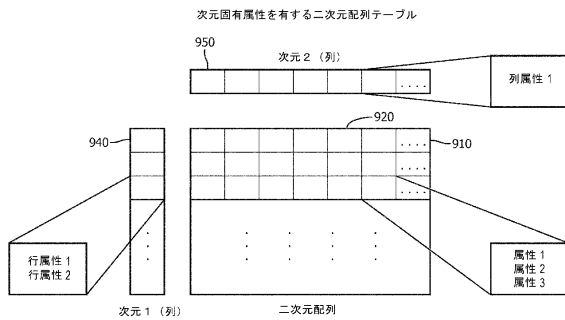


図 9

【図 10】

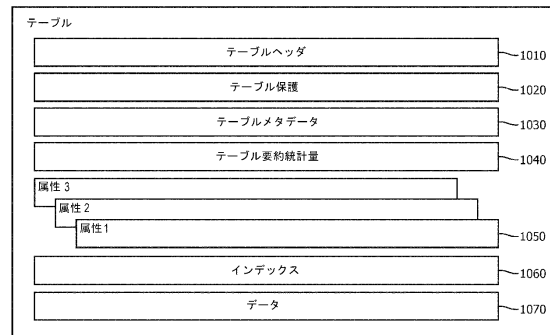


図 10

30

40

50

【 図 1 1 】

| フィールド                               | 簡単な説明                   | タイプ          |
|-------------------------------------|-------------------------|--------------|
| TableID                             | 1105 § T 1 と同じ          | ゲノム情報 (§T3)  |
| TableInfo                           | 1110 § T 1 と同じ          | ゲノム情報 (§T3)  |
| Table protection                    | 1020 アクセス制御ポリシー         | ゲノム情報 (§T3)  |
| Table metadata                      | 1030 メタデータ/リンケージ        | ゲノム情報 (§T3)  |
| Summary statistics                  | 1040 例えばカウント、平均値 1140   | リスト(キー値)     |
| nDimensions                         | 1150 次元の数               | 整数           |
| For i in 1 ... nDimensions:         |                         |              |
| Size [i]                            | 次元 i のサイズ               | 整数           |
| Dimension name [i]                  |                         |              |
| Dimension metadata [i]              | メタデータ/リンケージ             | ゲノム情報 (§T3)  |
| If nDimensions == 2:                |                         |              |
| Symmetry flag                       | 二次元配列が対称な場合は真           | ブール          |
| nAttributesMain                     | 1160                    |              |
| For i in 1 ... nAttributesMain:     |                         |              |
| Attribute info main [i]             |                         | 属性情報 (§T6)   |
| Byte offset main                    | Index main のバイトオフセット    | 整数           |
| If n dimensions > 1:                |                         |              |
| For i in 1 ... nDimensions:         |                         |              |
| nAttributes dim [i]                 | 次元固有属性の数                | 整数           |
| For j in 1 ... nAttributes dim [i]: |                         |              |
| Attribute info dim [i][j]           |                         | 属性情報 (§T6)   |
| Byte offset dim [i][j]              | Index dim [i] のバイトオフセット | 整数           |
| Index main                          |                         | インデックス (§T7) |
| Data payloads main                  |                         | データ (§T9)    |

図 1 1

【 図 1 2 】

属性情報構造

| フィールド                         | 簡単な説明   | タイプ         |
|-------------------------------|---|-------------|
| AttributeInfoSize             | 1205 構造をスキップすることを可能にするため  | 整数          |
| Attribute ID                  | 1210 一意の属性識別子   | 整数          |
| AttributeName                 | 1215  | 文字列         |
| AttributeMetadata             | 1220 属性のメタデータ/リンケージ/グループ化   | ゲノム情報 (§T3) |
| Attribute type                | 1225 基本型 (例えば整数、文字、float、文字列) または派生型 (例えば固定長、可変長配列)   | 文字列         |
| DefaultValue                  | 1230 ほとんどの値がデフォルトと合致する場合のスパース符号化のため   | 属性タイプ       |
| Summary statistics            | 1235 例えばカウント、平均値  | リスト(キー値)    |
| CompressorID                  | 1240 この属性に使用されるコンプレッサ   | 整数          |
| For i in 1 ... nDependencies: | § T 4 で定義されている nDependencies  |             |
| If nDimensions > 1:           |   |             |
| Dimension                     | これがメイン n 次元テーブルの属性である場合、これは、どの次元が依存関係属性を含むかを示す (依存関係属性がメイン Dimensional テーブルにもある場合は nDimensions + 1 に設定) | 整数          |
| AttributeID                   | 依存関係を含む AttributeID   | 整数          |
| CompressorCommonData size     |   | 整数          |
| CompressorCommonData          | すべてのチャンクに共通するコンプレッサのコードブック/統計モデルを保存するため   | バイト         |

図 1 2

【 図 1 3 】

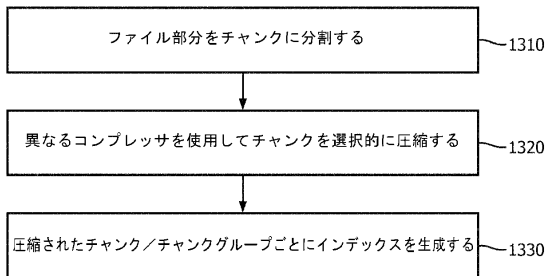


図 1 3

【 図 1 4 】

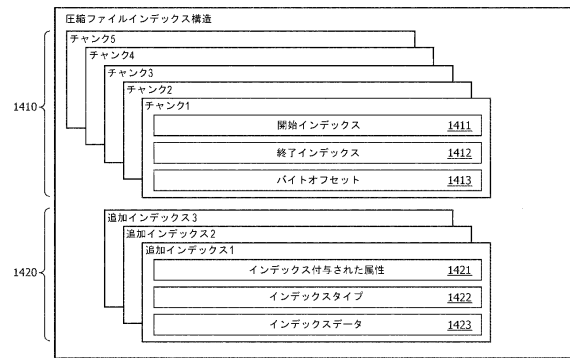


図 1 4

10

20

30

40

50

【 図 1 5 】

インデックス構造

| フィールド   | 簡単な説明  | タイプ        |
|---|--|------------|
| AttributeDependentChunks<br><br>1510          | チャンクサイズが属性に依存しているか、またはすべての属性に同じチャンク化が使用されているかを示すフラグ                  | ブール        |
| If not AttributeDependentChunks:              |  |            |
| ChunksStructure 1515                          |  |            |
| Else:   |  | チャンク(\$T8) |
| For i in 1 ... nAttributes:                   |  |            |
| Chunks structure [i] 1520                     |  |            |
| // Additional attribute specific indexes 1525 |  |            |
| n Additional indexes                          | 特定の属性 (例えば染色体、位置) に基づくより高速な照会のための追加インデックスの数。これらは、所望の照会結果を含むチャンク番号を返す | 整数         |
| For i in 1 ... nAdditionalIndexes:            |  |            |
| AttributeIDsIndexed [i] 1530                  | インデックス付与された属性のリスト  | リスト(整数)    |
| IndexType [i] 1535                            | インデックスタイプ (例えば、染色体およびゲノム位置のためのCSIインデックス、またはデータベース型クエリのためのBツリー)       | 文字列        |
| IndexSize [i] 1540                            | インデックスをスキップすることを可能にするため  | 整数         |
| IndexData [i] 1545                            | 実際のインデックスデータ、詳細はIndexType [i] に依存する                                  | バイト        |

図 1 5

【 図 1 6 】

チャンク構造

| フィールド                       | 簡単な説明                               | タイプ |
|-----------------------------|-------------------------------------|-----|
| nChunks 1605                | チャンクの数                              | 整数  |
| VariableSizeChunks 1610     | チャンクサイズが可変か固定かを示すフラグ (各次元の境界を除く)    | ブール |
| If VariableSizeChunks:      |                                     |     |
| For j in 1 ... nChunks:     |                                     |     |
| For k in 1 ... nDimensions: |                                     |     |
| StartIndex [j][k] 1615      | 次元 k 沿いのチャンクの開始位置                   | 整数  |
| EndIndex [j][k] 1620        | 次元 k 沿いのチャンクの終了位置                   | 整数  |
| ByteOffset [j] 1625         | ファイル内のチャンク j のバイトオフセット              | 整数  |
| Else:                       |                                     |     |
| For k in 1 ... nDimensions: |                                     |     |
| ChunkSize [k]               | 固定サイズのチャンクの場合、各次元でチャンクのサイズを格納するのに十分 | 整数  |
| For j in 1 ... nChunks:     |                                     |     |
| ByteOffset [j]              | ファイル内のチャンク i のバイトオフセット              | 整数  |

図 1 6

【 図 1 7 】

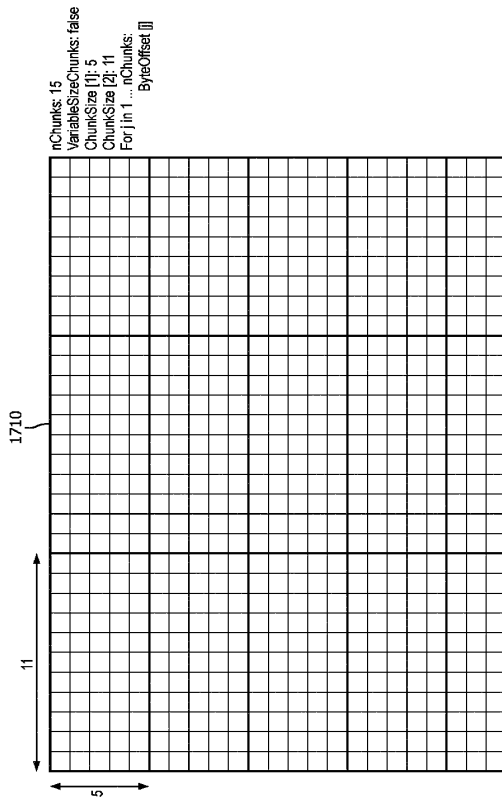


FIG. 17

【 図 1 8 】

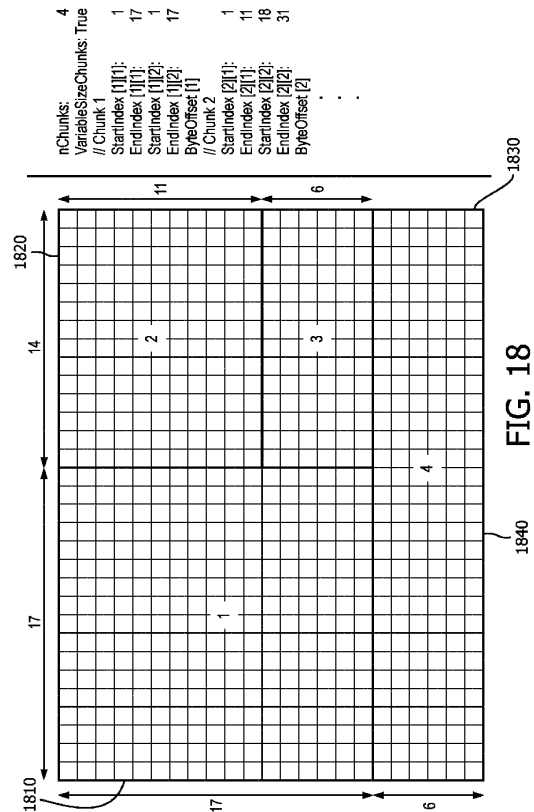


FIG. 18

【図 19】

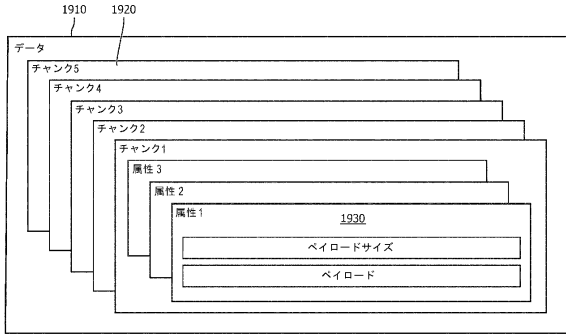


図 19

【図 20】

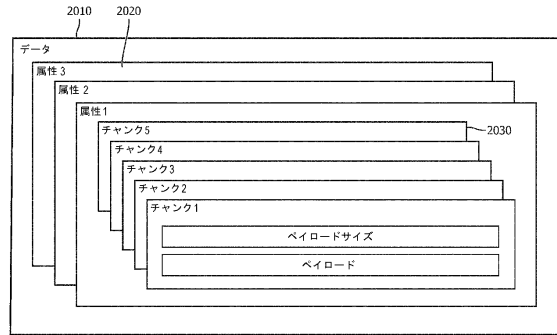


図 20

10

【図 21】

データペイロード

| フィールド                            | 簡単な説明                  | タイプ |
|----------------------------------|------------------------|-----|
| If not AttributeDependentChunks: |                        |     |
| For i in 1... nChunks:           |                        |     |
| For j in 1... nAttributes:       |                        |     |
| Payload size [i][j]              | 特定の属性をスキップすることを可能にするため | 整数  |
| Payload [i][j]                   | 圧縮されたペイロード             | バイト |
| Else:                            |                        |     |
| For i in 1... nAttributes:       |                        |     |
| For j in 1... nChunks:           |                        |     |
| Payload size [i][j]              |                        | 整数  |
| Payload [i][j]                   | 圧縮されたペイロード             | バイト |

図 21

【図 22】

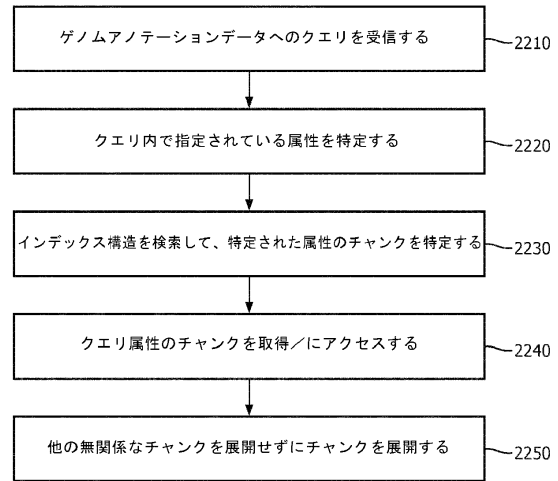


図 22

20

30

40

50

【図 2 3 A】

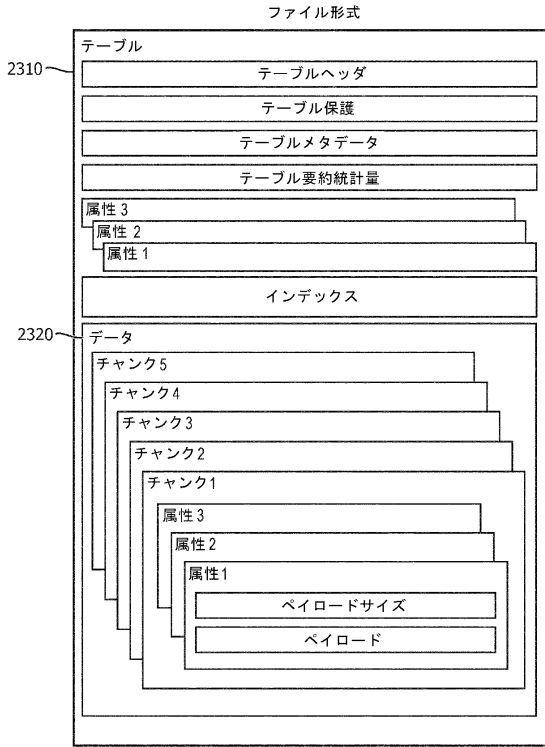


図 2 3 A

【図 2 3 B】

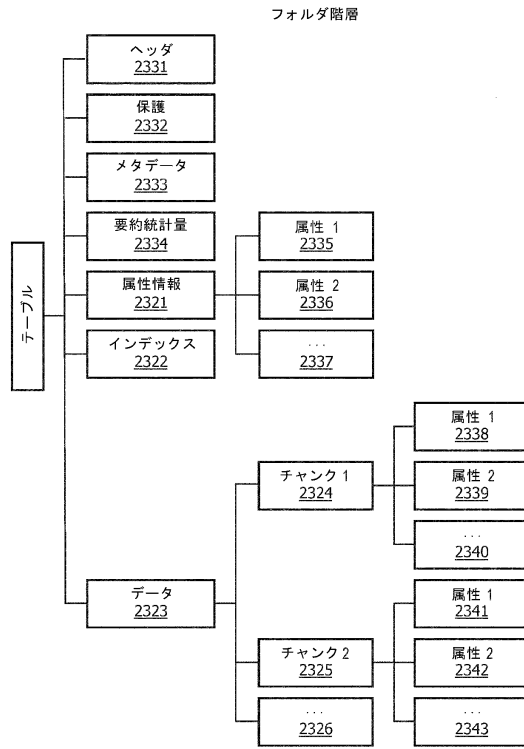


図 2 3 B

10

20

30

40

50

---

フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

フィリップス インターナショナル ビー． ヴィ． インテレクチュアル プロパティー アンド ス  
タンダーズ

(72)発明者 チャン イー ヒム

オランダ国 5 6 5 6 アーエー アインドーフェン ハイ テック キャンパス 5 フィリップス イ  
ンターナショナル ビー． ヴィ． インテレクチュアル プロパティー アンド スタンダーズ

審査官 渡邊 加寿磨

(56)参考文献 米国特許出願公開第 2 0 1 3 / 0 0 3 1 0 9 2 ( U S , A 1 )

特表 2 0 1 9 - 5 3 7 7 8 1 ( J P , A )

米国特許出願公開第 2 0 1 8 / 0 0 8 9 3 6 9 ( U S , A 1 )

国際公開第 2 0 1 8 / 0 7 1 0 5 5 ( W O , 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

G 0 6 F 1 2 / 0 0

G 0 6 F 1 6 / 0 0 - 1 6 / 9 5 8

H 0 3 M 7 / 3 0