

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)特許出願公開番号

特開2023-80989
(P2023-80989A)

(43)公開日

令和5年6月9日(2023.6.9)

(51)Int. Cl.

G 0 6 F 16/31 (2019.01)

F I

G 0 6 F 16/31

テーマコード(参考)

5 B 1 7 5

審査請求 有 請求項の数 21 O L (全 35 頁)

(21)出願番号 特願2021-194605(P2021-194605)

(22)出願日 令和3年11月30日(2021.11.30)

(71)出願人 720008140
先端加速システムズ株式会社
神奈川県横浜市金沢区泥亀一丁目28番B
-1312号

(74)代理人 100103850
弁理士 田中 秀▲てつ▼

(74)代理人 100105854
弁理士 廣瀬 一

(74)代理人 100116012
弁理士 宮坂 徹

(74)代理人 100066980
弁理士 森 哲也

最終頁に続く

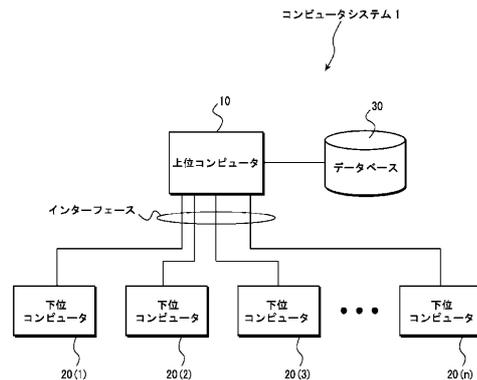
(54)【発明の名称】 近似文字列照合方法及び該方法を実現するためのコンピュータプログラム

(57)【要約】

【課題】 参照データを用いて、与えられるクエリデータの解析を高速及び/又は効率的に行う。

【解決手段】 本発明は、クエリ文字列に基づいて参照文字列における近似文字列を検索する方法である。前記方法は、前記参照文字列に基づいて階層的インデックスを作成することと、前記クエリ文字列の少なくとも一部と一致する前記参照文字列における部分文字列を同定するために、前記階層的インデックスを参照して、前記参照文字列に対するクエリ文字列のマッピングを行うことと、前記マッピングにより同定される少なくとも1以上の前記部分文字列に基づいて、前記近似文字列を導出することを含む。前記階層的インデックスは、前記参照文字列から切り出される各キーに対して、その出現回数に従い追加キーを追加しソートすることを繰り返しながら作成される。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

コンピューティングデバイスに、クエリ文字列に基づいて参照文字列における近似文字列を検索するための方法を実現させるためのコンピュータプログラムであって、

前記方法は、

前記参照文字列に基づいて階層的インデックスを作成することと、

前記クエリ文字列の少なくとも一部と一致する前記参照文字列における部分文字列を同定するために、前記階層的インデックスを参照して、前記参照文字列に対する前記クエリ文字列のマッピングを行うことと、

前記のマッピングにより同定される少なくとも 1 以上の前記部分文字列に基づいて、前記近似文字列を導出することと、を含み、

10

前記階層的インデックスを作成することは、

前記参照文字列から所定長の各第 1 のキーを切り出すことと、

切り出された前記各第 1 のキーについて、所定のハッシュ関数により該第 1 のキーに基づいて算出されるハッシュ値を割り当てた第 1 のキー配列を作成することと、

作成された前記第 1 のキー配列を更新することと、

更新された前記第 1 のキー配列を前記階層的インデックスとして出力することと、を含み、

前記第 1 のキー配列を更新することは、

前記第 1 のキー配列における前記各第 1 のキーについて、前記参照文字列における該第 1 のキーの出現回数を同定することと、

20

同定された前記第 1 のキーの前記出現回数に従って、該第 1 のキーに第 1 の追加キーを追加することにより新たな第 1 のキーを作成し、該新たな第 1 のキーに基づいて前記第 1 のキー配列を更新することと、を含む、

コンピュータプログラム。

【請求項 2】

前記第 1 のキー配列を作成することは、前記ハッシュ値に従って前記第 1 のキー配列における前記各第 1 のキーをソートすることを含む、

請求項 1 に記載のコンピュータプログラム。

【請求項 3】

30

前記第 1 のキー配列を更新することは、

前記同定した出現回数が所定の許容値を超えているか否かを判断することと、

前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記第 1 のキーに対して前記参照文字列における該第 1 のキーに続く少なくとも 1 以上の文字からなる前記第 1 の追加キーを追加することにより前記新たな第 1 のキーを作成することと、

前記新たな第 1 のキーについて、前記参照文字列における該第 1 のキーの出現回数を同定することと、を含む、

請求項 1 又は 2 に記載のコンピュータプログラム。

【請求項 4】

前記第 1 のキー配列を更新することは、前記第 1 の追加キーに従って前記第 1 のキー配列における前記新たな第 1 のキーをソートすることを更に含む、

40

請求項 1 から 3 のいずれか一項に記載のコンピュータプログラム。

【請求項 5】

前記第 1 のキー配列を更新することは、前記同定された前記出現回数が所定の許容値を超えていないと判断されるまで、現在の前記第 1 のキーに新たな前記第 1 の追加キーを順次に追加することにより新たな前記第 1 のキーを作成することを含む、

請求項 3 又は 4 に記載のコンピュータプログラム。

【請求項 6】

前記キー配列を前記階層的インデックスとして出力することは、

前記同定された前記出現回数が所定の許容値を超えていないと判断される場合に、現在

50

の前記キー配列を前記階層的インデックスとして出力することを含む、
請求項 3 から 5 のいずれか一項に記載のコンピュータプログラム。

【請求項 7】

前記マッピングを行うことは、
前記クエリ文字列から所定長の各第 2 のキーを切り出すことと、
前記クエリ文字列から切り出された前記各第 2 のキーについて、前記所定のハッシュ関数により該第 2 のキーに基づいて算出されるハッシュ値を割り当てた第 2 のキー配列を作成することと、
前記各第 2 のキーについて、前記ハッシュ値に従って、所定のサンプリング間隔で、前記階層的インデックスを参照し、該第 2 のキーの出現開始位置及び出現回数を同定することと、を含む、
請求項 1 から 6 のいずれか一項に記載のコンピュータプログラム。

10

【請求項 8】

前記第 2 のキーの前記出現開始位置及び前記出現回数を同定することは、
前記第 2 のキーの前記出現回数が前記所定の許容値を超えているか否かを判断することと、
前記第 2 のキーの前記出現回数が前記所定の許容値を超えていると判断される場合に、前記第 2 のキーに対して前記クエリ文字列における該第 2 のキーに続く少なくとも 1 以上の文字からなる第 2 の追加キーを追加することにより新たな第 2 のキーを作成することと、
前記第 2 のキーの前記出現回数が前記所定の許容値を超えていないと判断される場合に、同定された現在の前記第 2 のキーを一致文字列として出力するとともに該第 2 のキーの前記出現開始位置を出力することと、を含む、
請求項 7 に記載のコンピュータプログラム。

20

【請求項 9】

前記第 2 のキーの前記出現開始位置及び前記出現回数を同定することは、前記第 2 のキーの前記同定された前記出現回数が前記所定の許容値を超えていないと判断されるまで、現在の前記第 2 のキーに新たな前記第 2 の追加キーを順次に追加して、前記新たな第 2 のキーを作成することを更に含む、
請求項 8 に記載のコンピュータプログラム。

30

【請求項 10】

前記第 2 のキーの前記出現回数が前記所定の許容値を超えていると判断される場合に、該第 2 のキーの前記所定のサンプリング間隔を大きくする、
請求項 8 又は 9 に記載のコンピュータプログラム。

【請求項 11】

前記近似文字列を導出することは、
前記マッピングにより同定された前記一致文字列に基づく、被照合文字列と照合文字列とからなる文字列ペアを受信することと、
前記文字列ペアに基づいて少なくとも 1 つの近似文字列を導出するために、所定のアライメント処理を実行することと、
導出された前記少なくとも 1 つの近似文字列を出力することと、を含む、
請求項 8 から 10 のいずれか一項に記載のコンピュータプログラム。

40

【請求項 12】

前記所定のアライメント処理を実行することは、
前記被照合文字列と前記照合文字列とに基づいて所定のアライメント表を作成することと、
前記アライメント表の対角線上の要素を中心にした幅 m を有する計算領域を設定することと、
設定された前記計算領域における各要素について、変異度を算出することと、
算出された前記変異度に基づいて、最大変異度を決定することと、

50

決定された前記最大変異度に基づいて、前記少なくとも1つの近似文字列を導出することを含む、

請求項11に記載のコンピュータプログラム。

【請求項13】

前記所定のラインメント処理を実行することは、

前記最大変異度と所定の下限値とを比較して、前記最大変異度が前記所定の下限値を超えているかを判断することと、

前記最大変異度が前記所定の下限値を超えていないと判断される場合に、新たな計算領域を設定するために、前記計算領域の前記幅 m を拡幅することと、

前記最大変異度が前記所定の下限値を超えていると判断される場合に、前記最大変異度を有する要素に基づいて、前記少なくとも1つの近似文字列を導出することと、を含み、

前記最大変異度が前記下限値を超えると判断されるまで、前記計算領域を拡幅することにより新たな計算領域を設定して前記変異度を算出することを繰り返す、

請求項12に記載のコンピュータプログラム。

【請求項14】

前記所定の下限値は、所定の要素列に m 個の連続したギャップがあり、それ以外の部分は一致したと仮定した場合の変異度の値である、

請求項13に記載のコンピュータプログラム。

【請求項15】

前記一致文字列に対して前記参照文字列における対応する所定の文字列を追加することにより前記被照合文字列を作成することと、

前記一致文字列に対して前記クエリ文字列における対応する所定の文字列を追加することにより前記照合文字列を作成することと、を更に含む、

請求項11から14のいずれか一項に記載のコンピュータプログラム。

【請求項16】

コンピューティングデバイスに、クエリ文字列に基づいて参照文字列を探索するための階層的インデックスを作成する方法を実現させるためのコンピュータプログラムであって、

前記方法は、

前記参照文字列から所定長の各第1のキーを切り出すことと、

切り出された前記各第1のキーについて、所定のハッシュ関数により該第1のキーに基づいて算出されるハッシュ値を割り当てた第1のキー配列を作成することと、

作成された前記第1のキー配列を更新することと、

更新された前記第1のキー配列を前記階層的インデックスとして出力することと、を含み、

前記第1のキー配列を更新することは、

前記第1のキー配列における前記各第1のキーについて、前記参照文字列における該第1のキーの出現開始位置及び出現回数を同定することと、

同定された前記第1のキーの前記出現開始位置及び前記出現回数に従って、該第1のキーに第1の追加キーを追加することにより新たな第1のキーを作成し、該新たな第1のキーに基づいて前記第1のキー配列を更新することと、を含み、

コンピュータプログラム。

【請求項17】

コンピューティングデバイスに、参照文字列に対してクエリ文字列のマッピングを行う方法を実現させるためのコンピュータプログラムであって、

前記方法は、

前記参照文字列に基づく階層的インデックスを読み出すことと、

前記クエリ文字列から所定のキー長を有する各キーを切り出して、キー配列を作成することと、

前記クエリ文字列から切り出された前記各キーについて、前記所定のハッシュ関数によ

10

20

30

40

50

り該キーに基づいて算出されるハッシュ値を割り当てたキー配列を作成することと、
前記各キーについて、前記ハッシュ値に従って、所定のサンプリング間隔で、前記階層的インデックスを参照し、該キーの出現開始位置及び出現回数を同定することと、
前記同定した出現回数が所定のしきい値を超えているか否かを判断することと、
前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記キーに対して前記クエリ文字列における該キーに続く少なくとも1以上の文字からなる追加キーを追加することにより新たなキーを作成することと、
前記同定された前記出現回数が所定のしきい値を超えていないと判断される場合に、同定された現在の前記キーの出現開始位置及び該キーを出力することと、を含み、
前記キーの前記出現開始位置及び前記出現回数を同定することは、前記同定された前記出現回数が所定のしきい値を超えていないと判断されるまで、現在の前記キーに新たな前記追加キーを順次に追加して、前記新たなキーを作成することを含む、
コンピュータプログラム

10

【請求項18】

前記方法は、前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記キーの前記所定のサンプリング間隔を大きくするように構成される、
請求項17に記載のコンピュータプログラム。

【請求項19】

コンピューティングデバイスに、参照文字列における部分文字列とクエリ文字列との間の変異を所定のアラインメント処理により同定する方法を実現させるためのコンピュータプログラムであって、

20

前記方法は、

マッピングにより同定された一致文字列に基づく、被照合文字列と照合文字列とからなる文字列ペアを受信することと、

前記文字列ペアに基づいて少なくとも1つの近似文字列を導出するために、所定のアラインメント処理を実行することと、

導出された前記少なくとも1つの近似文字列を出力することと、を含み、

前記所定のアラインメント処理を実行することは、

前記被照合文字列と前記照合文字列とに基づいて所定のアラインメント表を作成することと、

30

前記アラインメント表の対角線上の要素を中心にした幅 m を有する計算領域を設定することと、

設定された前記計算領域における各要素について、変異度を算出することと、

算出された前記変異度に基づいて、最大変異度を決定することと、

決定された前記最大変異度に基づいて、前記少なくとも1つの近似文字列を導出することを含む、
コンピュータプログラム。

【請求項20】

前記所定のアラインメント処理を実行することは、

前記最大変異度と所定の下限値とを比較して、前記最大変異度が前記所定の下限値を超えているかを判断することと、

40

前記最大変異度が前記所定の下限値を超えていないと判断される場合に、新たな計算領域を設定するために、前記計算領域の前記幅 m を拡幅することと、

前記最大変異度が前記所定の下限値を超えていると判断される場合に、前記最大変異度を有する要素に基づいて、前記少なくとも1つの近似文字列を導出することと、を含み、

前記最大変異度が前記下限値を超えると判断されるまで、前記計算領域を拡幅することにより新たな計算領域を設定して前記変異度を算出することを繰り返す、
請求項19に記載のコンピュータプログラム。

【請求項21】

前記所定のアラインメント処理を実行することは、所定の要素列に m 個の連続したギャ

50

ップがあり、それ以外の部分は一致したと仮定した場合の変異度の値を前記所定の下限值として設定することを更に含む、
請求項 20 に記載のコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、近似文字列照合技術に関し、特に、ヒトゲノムの解析に利用可能な近似文字列照合装置及び近似文字列照合方法並びに該方法を実現するためのコンピュータプログラムに関する。

10

【背景技術】

【0002】

ヒトゲノムは、人が持つ遺伝情報のセットであり、これを担っている物質が、約 30 億対の塩基が連なった DNA (デオキシリボ核酸) である。塩基は、アデニン (A)、グアニン (G)、シトシン (C)、及びチミン (T) がある。すなわち、人の遺伝情報は、これらの塩基の並び (配列) によって決定される。

【0003】

ヒトゲノムの読み取りにはシーケンサと称される装置が用いられる。シーケンサは、サンプルとなるヒトゲノムを読み取って、これを所定の上限值 (数百塩基対程度) に細断して増幅し、データ片からなる膨大なデータ配列として出力する。現行のシーケンサは、一人分のヒトゲノムを 1 時間ほどで読み出すことができる。

20

【0004】

シーケンサにより読み出されるばらばらのデータ片は、人の標準的なゲノム配列として定められたヒトゲノム参照配列と比較されることによって、元の長さのヒトゲノム配列に再構築され解析される。例えば、各データ片が、ヒトゲノム参照配列との比較において、どの位置にあるかが調べられ (マッピング)、また、どのような変異があるかといった解析がなされる。通常、シーケンサから読み出されるデータには誤差が含まれるため、1 サンプルあたりヒトゲノム配列一人分の例えば 30 倍の冗長データを用いて統計処理を行うことにより誤差を小さくしている。したがって、ヒトゲノム配列の解析には膨大な計算量が必要とされるため、典型的には、スーパーコンピュータやクラスタコンピュータ、FPGA (Field-Programmable Gate Array) ベースのコンピュータといった高性能なコンピュータが用いられる。

30

【0005】

データ片のマッピングには、例えば BWA (Burrow-Wheeler Aligner) といったプログラムツールが用いられる。BWA は、3 つのアルゴリズム、BWA - b a c k t r a c k、BWA - S W、及び BWA - M E M から構成される。このうち、BWA - M E M は、In del (挿入欠失) に対応した高速アルゴリズムとして広く利用されている。BWA - M E M は、読み出したデータ片に基づくクエリ文字列のうち、ヒトゲノム参照配列に繰り返し現れる部分に対して、接尾辞配列を用いてインデックスを作成し、マッピングを行うアルゴリズムである (非特許文献 1)。

40

【0006】

また、データ片は、マッピングにより得られるヒトゲノム参照配列の部分配列と照合され、どのような変異があるかが解析される。ヒトゲノム参照配列の部分配列とデータ片との照合には、例えば、アラインメントアルゴリズムが用いられる。アラインメントアルゴリズムでは、動的計画法に従って、アラインメント表と称される配列の各要素について変異度 (類似度) を算出しながら、近似文字列を同定する (非特許文献 2)。

【先行技術文献】

【特許文献】

【0007】

50

【非特許文献1】 Heng Li, “Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM”, May 26, 2013, arXiv:1303.3997 (q-bio.GN)

【非特許文献2】 内山将夫 他, 「近似文字列照合による全文検索のための接尾辞配列の高速走査法」, 2002年9月15日, 情報処理学会, Vol. 43 No. SIG 9(TOD 15)

【発明の概要】

【発明が解決しようとする課題】

【0008】

上述したBWA-MEMは、ギャップを許容しかつ高速マッピングが可能なアルゴリズムとして広く利用されているが、作成されるインデックスのサイズが非常に大きくなるため、大量のメモリリソースを必要とするという問題があった。また、BWA-MEMでは、クエリ文字列の文字数（すなわち、塩基数）に応じた回数だけメモリへのアクセスが必要となるため、メモリへのランダムアクセスが頻繁に発生し、これにより、プロセッサの処理速度がメモリへのアクセス時間により律速されてしまうという問題があった。

10

【0009】

より具体的には、BWA-MEMでは、クエリ文字列における部分文字列（これを「キー」と称することがある。）がヒトゲノム参照配列内に出現する位置を同定するために、メインメモリへのランダムアクセスがクエリ文字列の総文字数分の回数だけ実行される。これは、ヒトゲノムの解析に用いるデータの量が非常に大きいため、必要なデータ配列をより高速アクセス可能なキャッシュメモリに一度に収容しきれないからである。ここで、メインメモリへの1回のランダムアクセスの待ち時間（アクセスが発生してから実際にデータが得られるまでの時間）は、現在のコンピュータでは、1バンクあたり約1μ秒（0.000001秒）かかっている。したがって、30億塩基対のヒトゲノムの解析の場合、冗長性を30とすると、総文字数に対する総アクセス時間Tは、

20

$$T = 3,000,000,000 \times 0.000001 \times 30 \\ = 90,000 \text{ (秒)}$$

となる。つまり、ヒトゲノムの解析において、メインメモリへのランダムアクセスだけで、約9万秒（約25時間）かかることになる。このため、たとえ、高性能なプロセッサを備えたコンピュータを用いたとしても、メモリアクセス時間が制約となって、マッピング時間を短縮するには限界があった。

【0010】

また、上述したように、シーケンサにより一人分のヒトゲノムを読み出すためには、現状、1時間ほど要している。一方で、高性能コンピュータを用いてBWA-MEMを実行した場合、シーケンサの読み出し時間以内にマッピングが完了し、両者のバランスは概ね保たれている。

30

【0011】

一方で、次世代型シーケンサは、低コスト化が進み、また、読み出し時間を更に短縮し得ると言われており、これに伴って、解析時間もまた短縮することが望まれる。解析時間を短縮するための一つのアプローチとして、コンピュータの更なる高性能化が考えられるが、コンピュータの高性能化のためには非常にコストが高いため実用化へのハードルが高い。

40

【0012】

更に、マッピングされたデータ片の解析に用いられる従前のアラインメントアルゴリズムでは、アラインメント表の全ての要素について変異度を算出するため、計算量が多くなり、時間がかかるという問題がある。

【0013】

そこで、本発明は、参照データを用いて、与えられるクエリデータの解析を高速及び/又は効率的に行うことができる新たな技術を提供することを目的とする。

【0014】

より具体的には、本発明の一つの目的は、与えられるクエリ文字列と参照文字列との間の近似文字列照合を高速及び/又は効率的に行うことができる近似文字列照合装置及びこ

50

れを用いた近似文字列照合方法を提供することである。

【0015】

また、本発明の一つの目的は、シーケンサによって読み出されたデータ片に基づくヒトゲノム参照配列を用いた解析を高速及び/又は効率的に行うことができる近似文字列照合装置及びこれを用いた近似文字列照合方法を提供することである。

【0016】

また、本発明の一つの目的は、前記近似文字列照合に適合した参照文字列に基づく階層的インデックスを作成する技術を提供することである。

【課題を解決するための手段】

【0017】

上記課題を解決するための本発明は、以下に示す発明特定事項乃至は技術的特徴を含んで構成される。

【0018】

ある観点に従う本発明は、コンピューティングデバイスに、クエリ文字列に基づいて参照文字列における近似文字列を検索するための方法を実現させるためのコンピュータプログラムである。

前記方法は、前記参照文字列に基づいて階層的インデックスを作成することと、前記クエリ文字列の少なくとも一部と一致する前記参照文字列における部分文字列を同定するために、前記階層的インデックスを参照して、前記参照文字列に対する前記クエリ文字列のマッピングを行うことと、前記マッピングにより同定される少なくとも1以上の前記部分文字列に基づいて、前記近似文字列を導出することと、を含む。

ここで、前記階層的インデックスを作成することは、前記参照文字列から所定長の各第1のキーを切り出すことと、切り出された前記各第1のキーについて、所定のハッシュ関数により該第1のキーに基づいて算出されるハッシュ値を割り当てた第1のキー配列を作成することと、作成された前記第1のキー配列を更新することと、更新された前記第1のキー配列を前記階層的インデックスとして出力することと、を含む。

また、前記第1のキー配列を更新することは、前記第1のキー配列における前記各第1のキーについて、前記参照文字列における該第1のキーの出現回数を同定することと、同定された前記第1のキーの前記出現回数に従って、該第1のキーに第1の追加キーを追加することにより新たな第1のキーを作成し、該新たな第1のキーに基づいて前記第1のキー配列を更新することと、を含む。

【0019】

前記第1のキー配列を作成することは、前記ハッシュ値に従って前記第1のキー配列における前記各第1のキーをソートすることを含み得る。

【0020】

また、前記第1のキー配列を更新することは、前記同定した出現回数が所定の許容値を超えているか否かを判断することと、前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記第1のキーに対して前記参照文字列における該第1のキーに続く少なくとも1以上の文字からなる前記第1の追加キーを追加することにより前記新たな第1のキーを作成することと、前記新たな第1のキーについて、前記参照文字列における該第1のキーの出現回数を同定することと、を含み得る。

【0021】

また、前記第1のキー配列を更新することは、前記第1の追加キーに従って前記第1のキー配列における前記新たな第1のキーをソートすることを更に含み得る。

【0022】

また、前記第1のキー配列を更新することは、前記同定された前記出現回数が所定の許容値を超えていないと判断されるまで、現在の前記第1のキーに新たな前記第1の追加キーを順次に追加することにより新たな前記第1のキーを作成することを含み得る。

【0023】

また、前記キー配列を前記階層的インデックスとして出力することは、前記同定された

10

20

30

40

50

前記出現回数が所定の許容値を超えていないと判断される場合に、現在の前記キー配列を前記階層的インデックスとして出力することを含み得る。

【0024】

前記マッピングを行うことは、前記クエリ文字列から所定長の各第2のキーを切り出すことと、前記クエリ文字列から切り出された前記各第2のキーについて、前記所定のハッシュ関数により該第2のキーに基づいて算出されるハッシュ値を割り当てた第2のキー配列を作成することと、前記各第2のキーについて、前記ハッシュ値に従って、所定のサンプリング間隔で、前記階層的インデックスを参照し、該第2のキーの出現開始位置及び出現回数を同定することと、を含み得る。

【0025】

前記第2のキーの前記出現開始位置及び前記出現回数を同定することは、前記第2のキーの前記出現回数が前記所定の許容値を超えているか否かを判断することと、前記第2のキーの前記出現回数が前記所定の許容値を超えていると判断される場合に、前記第2のキーに対して前記クエリ文字列における該第2のキーに続く少なくとも1以上の文字からなる第2の追加キーを追加することにより新たな第2のキーを作成することと、前記第2のキーの前記出現回数が前記所定の許容値を超えていないと判断される場合に、同定された現在の前記第2のキーを一致文字列として出力するとともに該第2のキーの前記出現開始位置を出力することと、を含み得る。

【0026】

また、前記第2のキーの前記出現開始位置及び前記出現回数を同定することは、前記第2のキーの前記同定された前記出現回数が前記所定の許容値を超えていないと判断されるまで、現在の前記第2のキーに新たな前記第2の追加キーを順次に追加して、前記新たな第2のキーを作成することを含み得る。

【0027】

また、前記第2のキーの前記出現回数が前記所定の許容値を超えていると判断される場合に、該第2のキーの前記所定のサンプリング間隔が大きくなるように変更され得る。

【0028】

前記近似文字列を導出することは、前記マッピングにより同定された前記一致文字列に基づく、被照合文字列と照合文字列とからなる文字列ペアを受信することと、前記文字列ペアに基づいて少なくとも1つの近似文字列を導出するために、所定のアラインメント処理を実行することと、導出された前記少なくとも1つの近似文字列を出力することと、を含み得る。

【0029】

前記所定のアラインメント処理を実行することは、前記被照合文字列と前記照合文字列とに基づいて所定のアラインメント表を作成することと、前記アラインメント表の対角線上の要素を中心にした幅 m を有する計算領域を設定することと、設定された前記計算領域における各要素について、変異度を算出することと、算出された前記変異度に基づいて、最大変異度を決定することと、決定された前記最大変異度に基づいて、前記少なくとも1つの近似文字列を導出することを含み得る。

【0030】

また、前記所定のアラインメント処理を実行することは、前記最大変異度と所定の下限値とを比較して、前記最大変異度が前記所定の下限値を超えているかを判断することと、前記最大変異度が前記所定の下限値を超えていないと判断される場合に、新たな計算領域を設定するために、前記計算領域の前記幅 m を拡幅することと、前記最大変異度が前記所定の下限値を超えていると判断される場合に、前記最大変異度を有する要素に基づいて、前記少なくとも1つの近似文字列を導出することと、を含み得る。そして、前記方法は、前記最大変異度が前記下限値を超えると判断されるまで、前記計算領域を拡幅することにより新たな計算領域を設定して前記変異度を算出することを繰り返すように構成され得る。

【0031】

10

20

30

40

50

前記所定の下限值は、所定の要素列にm個の連続したギャップがあり、それ以外の部分は完全又は実質的に一致したと仮定した場合の変異度の値であり得る。

【0032】

また、前記方法は、前記一致文字列に対して前記参照文字列における対応する所定の文字列を追加することにより前記被照合文字列を作成することと、前記一致文字列に対して前記クエリ文字列における対応する所定の文字列を追加することにより前記照合文字列を作成することと、を更に含む得る。

【0033】

また、ある観点に従う本発明は、コンピューティングデバイスに、クエリ文字列に基づいて参照文字列を探索するための階層的インデックスを作成する方法を実現させるためのコンピュータプログラムである。

前記方法は、前記参照文字列から所定長の各第1のキーを切り出すことと、切り出された前記各第1のキーについて、所定のハッシュ関数により該第1のキーに基づいて算出されるハッシュ値を割り当てた第1のキー配列を作成することと、作成された前記第1のキー配列を更新することと、更新された前記第1のキー配列を前記階層的インデックスとして出力することと、を含む。

ここで、前記第1のキー配列を更新することは、前記第1のキー配列における前記各第1のキーについて、前記参照文字列における該第1のキーの出現開始位置及び出現回数を同定することと、同定された前記第1のキーの前記出現開始位置及び前記出現回数に従って、該第1のキーに第1の追加キーを追加することにより新たな第1のキーを作成し、該新たな第1のキーに基づいて前記第1のキー配列を更新することと、を含む。

【0034】

また、ある観点に従う本発明は、コンピューティングデバイスに、参照文字列に対してクエリ文字列のマッピングを行う方法を実現させるためのコンピュータプログラムである。

前記方法は、前記参照文字列に基づく階層的インデックスを読み出すことと、前記クエリ文字列から所定のキー長を有する各キーを切り出して、キー配列を作成することと、前記クエリ文字列から切り出された前記各キーについて、前記所定のハッシュ関数により該キーに基づいて算出されるハッシュ値を割り当てたキー配列を作成することと、前記各キーについて、前記ハッシュ値に従って、所定のサンプリング間隔で、前記階層的インデックスを参照し、該キーの出現開始位置及び出現回数を同定することと、前記同定した出現回数が所定のしきい値を超えているか否かを判断することと、前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記キーに対して前記クエリ文字列における該キーに続く少なくとも1以上の文字からなる追加キーを追加することにより新たなキーを作成することと、前記同定された前記出現回数が所定のしきい値を超えていないと判断される場合に、同定された現在の前記キーの出現開始位置及び該キーを出力することと、を含む。

そして、前記キーの前記出現開始位置及び前記出現回数を同定することは、前記同定された前記出現回数が所定のしきい値を超えていないと判断されるまで、現在の前記キーに新たな前記追加キーを順次に追加して、前記新たなキーを作成することを含む。

【0035】

ここで、前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記キーの前記所定のサンプリング間隔が大きくなるように変更され得る。

【0036】

また、ある観点に従う本発明は、コンピューティングデバイスに、参照文字列における部分文字列とクエリ文字列との間の変異を所定のアラインメント処理により同定する方法を実現させるためのコンピュータプログラムである。

前記方法は、マッピングにより同定された一致文字列に基づく、被照合文字列と照合文字列とからなる文字列ペアを受信することと、前記文字列ペアに基づいて少なくとも1つの近似文字列を導出するために、所定のアラインメント処理を実行することと、

導出された前記少なくとも1つの近似文字列を出力することと、を含む。

そして、前記所定のアラインメント処理を実行することは、前記被照合文字列と前記照合文字列とに基づいて所定のアラインメント表を作成することと、前記アラインメント表の対角線上の要素を中心にした幅 m を有する計算領域を設定することと、設定された前記計算領域における各要素について、変異度を算出することと、算出された前記変異度に基づいて、最大変異度を決定することと、決定された前記最大変異度に基づいて、前記少なくとも1つの近似文字列を導出することを含む。

【0037】

また、前記所定のアラインメント処理を実行することは、前記最大変異度と所定の下限値とを比較して、前記最大変異度が前記所定の下限値を超えているかを判断することと、前記最大変異度が前記所定の下限値を超えていないと判断される場合に、新たな計算領域を設定するために、前記計算領域の前記幅 m を拡幅することと、前記最大変異度が前記所定の下限値を超えていると判断される場合に、前記最大変異度を有する要素に基づいて、前記少なくとも1つの近似文字列を導出することと、を含み得る。

10

そして、前記最大変異度が前記下限値を超えると判断されるまで、前記計算領域を拡幅することにより新たな計算領域を設定して前記変異度を算出することが繰り返され得る。

【0038】

また、前記所定のアラインメント処理を実行することは、所定の要素列に m 個の連続したギャップがあり、それ以外の部分は完全に又は実質的に一致したと仮定した場合の変異度の値を前記所定の下限値として設定することを更に含み得る。

20

【0039】

なお、本発明は、前記コンピュータプログラムを記憶した記録媒体としても成立する。また、本発明は、前記方法を実行するように構成されたハードウェア及び/又はファームウェアからなる装置としても成立する。近似文字列照合装置は、本発明の一形態である。

【0040】

なお、本明細書等において、手段又は部(unit)とは、単に物理的手段を意味するものではなく、その手段又は部が有する機能をソフトウェアによって実現する場合も含む。また、1つの手段又は部が有する機能が2つ以上の物理的手段により実現されても、2つ以上の手段又は部の機能が1つの物理的手段により実現されても良い。また、「システム」とは、複数の装置(又は特定の機能を実現する機能モジュール)が論理的に集合した物のことをいい、各装置や機能モジュールが単一の筐体内にあるか否かは特に問わない。

30

【発明の効果】

【0041】

本発明によれば、与えられるクエリデータに対する参照データを用いた解析を高速及び/又は効率的に行うことができる。とりわけ、本発明によれば、与えられるクエリ文字列と参照文字列との間の近似文字列照合を高速及び/又は効率的に行うことができる。

【0042】

また、本発明によれば、シークエンサによって読み出されたデータ片に基づくヒトゲノム参照配列を用いた解析を高速及び/又は効率的に行うことができる。

【0043】

40

更に、本発明によれば、近似文字列照合に適合した参照文字列に基づく階層的インデックスを提供することができる。

【0044】

本発明の他の技術的特徴、目的、及び作用効果乃至は利点は、添付した図面を参照して説明される以下の実施形態により明らかにされる。本明細書に記載された効果はあくまで例示であって限定されるものではなく、また他の効果があっても良い。

【図面の簡単な説明】

【0045】

【図1】図1は、本発明の一実施形態に係るコンピュータシステムの概略的構成の一例を示すブロックダイアグラムである。

50

【図2】図2は、本発明の一実施形態に係るコンピュータシステムによる近似文字列照合処理の概略の一例を説明するフローチャートである。

【図3】図3は、本発明の一実施形態に係るコンピュータシステムによるインデックス作成処理の一例を説明するフローチャートである。

【図4】図4は、本発明の一実施形態に係るコンピュータシステムにおいて用いられる参照文字列の一例を示す図である。

【図4A】図4Aは、図4に示される参照文字列に基づく階層的インデックスの作成過程におけるデータ配列構造の一例を説明するための図である。

【図4B】図4Bは、図4に示される参照文字列に基づく階層的インデックスの作成過程におけるデータ配列構造の一例を説明するための図である。

【図4C】図4Cは、図4に示される参照文字列に基づく階層的インデックスの作成過程におけるデータ配列構造の一例を説明するための図である。

【図5】図5は、図4に示される参照文字列に基づく階層的インデックスの作成過程におけるキー開始位置及び出現回数を示すテーブル構造の一例を示す図である。

【図6】図6は、本発明の一実施形態に係るコンピュータシステムにおいて用いられる参照文字列の一例を示す図である。

【図7】図7は、本発明の一実施形態に係るコンピュータシステムによるクエリ文字列に基づく参照文字列の探索処理の一例を説明するフローチャートである。

【図8A】図8Aは、Needleman-Wunschアルゴリズムを説明するためのアラインメント表の一例を示す図である。

【図8B】図8Bは、Needleman-Wunschアルゴリズムを説明するためのアラインメント表の一例を示す図である。

【図8C】図8Cは、Needleman-Wunschアルゴリズムを説明するためのアラインメント表の一例を示す図である。

【図8D】図8Dは、Needleman-Wunschアルゴリズムを説明するためのアラインメント表の一例を示す図である。

【図9】図9は、本発明の一実施形態に係るコンピュータシステムによる動的計画法を用いたアラインメント処理の一例を説明するフローチャートである。

【図10A】図10Aは、本発明の一実施形態に係る動的計画法を用いたアラインメントを説明するためのアラインメント表の一例を示す図である。

【図10B】図10Bは、本発明の一実施形態に係る動的計画法を用いたアラインメントを説明するためのアラインメント表の一例を示す図である。

【図10C】図10Cは、本発明の一実施形態に係る動的計画法を用いたアラインメントを説明するためのアラインメント表の一例を示す図である。

【図11A】図11Aは、本発明の一実施形態に係る動的計画法を用いたアラインメントにより得られる近似文字列の一例を示す図である。

【図11B】図11Bは、本発明の一実施形態に係る動的計画法を用いたアラインメントにより得られる近似文字列の一例を示す図である。

【図12】本発明の一実施形態に係るシステムにおけるコンピューティングデバイスのハードウェア構成の一例を示す図である。

【発明を実施するための形態】

【0046】

以下、図面を参照して本発明の実施の形態を説明する。ただし、以下に説明する実施形態は、あくまでも例示であり、以下に明示しない種々の変形や技術の適用を排除する意図はない。本発明は、その趣旨を逸脱しない範囲で種々変形（例えば各実施形態を組み合わせる等）して実施することができる。また、以下の図面の記載において、同一又は類似の部分には同一又は類似の符号を付して表している。図面は模式的なものであり、必ずしも実際の寸法や比率等とは一致しない。図面相互間においても互いの寸法の関係や比率が異なる部分が含まれていることがある。

【0047】

10

20

30

40

50

図1は、本発明の一実施形態に係るコンピュータシステムの概略的構成の一例を示すブロックダイアグラムである。同図に示すように、コンピュータシステム1は、例えば、上位コンピュータ10と、複数の下位コンピュータ20(1)~20(n)と、データベース30とを含み構成される。上位コンピュータ10と下位コンピュータ20(1)~20(n)とは、所定のインターフェースを介して通信可能に接続される。本開示では、コンピュータシステム1は、上位コンピュータ10が複数の下位コンピュータ20(1)~20(n)を統括的に制御する中央集権型コンピュータシステムとして構成されるが、これに限られず、例えば分散型コンピュータシステムとして構成されても良い。分散型コンピュータシステムにおいては、各コンピュータが並列分散処理により協調的に動作し得るが、特定の処理に関しては、代表する一のコンピュータのみが該処理を実行する場合があっても良い。上位コンピュータ10及び下位コンピュータ20のハードウェア構成は、図12に例示されるが、当業者にとって自明であるため、その詳細な説明は省略する。本開示では、下位コンピュータ20は、上位コンピュータ10の制御の下、並列的にタスクを処理する。例えば、複数の下位コンピュータ20のそれぞれは、与えられた個々のクエリ文字列に基づいて参照文字列との比較において解析処理を行う。以下では、下位コンピュータ20(1)~20(n)について、それらを特に区別する必要がない限り、単に、「下位コンピュータ20」と表記することがある。

10

【0048】

データベース30は、上位コンピュータ10の制御の下、各種のデータ、例えば参照文字列を格納する。一例として、データベース30は、ヒトゲノム参照配列を格納する。ヒトゲノム参照配列は、人の標準的なゲノム配列として定められた塩基対の配列を示すデータである。また、データベース30は、参照文字列に基づいて作成されたインデックスを格納する。インデックスは、参照文字列を探索するために用いられるある種のデータ配列構造である。なお、図中、データベース30は、上位コンピュータ10にのみアクセス可能に接続される構成となっているが、これに限られず、下位コンピュータ20もまたアクセス可能に接続される構成であっても良い。

20

【0049】

図2は、本発明の一実施形態に係るコンピュータシステムによる近似文字列照合処理の概略の一例を説明するフローチャートである。かかる処理は、例えば、上位コンピュータ10及び複数の下位コンピュータ20が、プロセッサの制御の下、近似文字列照合プログラムを実行することにより他のハードウェアコンポーネントと協働し、実現される。

30

【0050】

すなわち、同図に示すように、まず、上位コンピュータ10は、データベース30に格納された参照文字列を読み出し、読み出した参照文字列に基づいて、インデックスを作成する(S201)。具体的には、上位コンピュータ10は、参照文字列における部分文字列を所定の条件に従って拡張し及び/又はソートして配列化することにより、所定のデータ配列構造を有するインデックスを作成する。データ配列構造は、階層的なツリー構造からなる。本開示では、このようなインデックスを階層的インデックスと称するものとする。上位コンピュータ10は、作成した参照文字列に基づくインデックスをデータベース30に格納する。参照文字列に基づく階層的インデックスの作成処理の詳細は後述される。なお、参照文字列に基づく階層的インデックスが既に作成されデータベース30に格納されている場合には、処理S201は省略され得る。

40

【0051】

続いて、上位コンピュータ10は、図示しない外部装置からクエリ文字列を受信し、受信したクエリ文字列に基づいて、階層的インデックスを参照して、マッピングを行う(S202)。一例では、外部装置は、ヒトゲノムを読み出すシーケンサであり、クエリ文字列は、シーケンサから出力される例えば150~200塩基程度の塩基対のデータ片である。本開示では、下位コンピュータ20もまた、上位コンピュータ10の制御の下、割り当てられたクエリ文字列に基づいて、階層的インデックスを参照して、マッピングを行う。

50

【 0 0 5 2 】

マッピングは、クエリ文字列の少なくとも一部と一致する参照文字列における部分文字列（一致文字列）を同定する処理である。つまり、マッピングにより、参照文字列におけるクエリ文字列の少なくとも一部と一致する部分文字列の出現開始位置及び長さ（文字数）が同定される。本開示に係るマッピングでは、クエリ文字列における各キーについて、階層的インデックスを検索ないしは探索することにより、参照文字列におけるクエリ文字列の各キー及びその出現開始位置が同定される。なお、処理の高速化のため、作成された階層的インデックスは、メインメモリ又はキャッシュメモリ等の高速メモリに展開され得る。

【 0 0 5 3 】

また、本開示に係るマッピングでは、階層的インデックスの探索に際して、クエリ文字列におけるキーの出現開始位置のサンプリング間隔が従来技術に比較してある程度大きくなるように調整される。例えば B L A S T (Basic Local Alignment Search Tool) と称される D N A の塩基配列やタンパク質のアミノ酸配列のシーケンスアラインメントを行うためのアルゴリズムでは、キーの出現開始位置の間隔は 3 ~ 5 文字に設定されるが、本開示に係る探索では、例えば 1 0 文字又はそれ以上に設定され得る。探索に際して、キーの出現開始位置のサンプリング間隔が小さいと、探索の回数が増え、効率が低下してしまうのに対して、キーの出現開始位置のサンプリング間隔を常に大きくしてしまうと、探索の高速化と引き換えに、見落としの確率が上昇してしまう。そこで、本開示では、キーの長さを長くしながら、それが一定長を超える場合には、出現開始位置のサンプリング間隔を大きくすることにより、キーの一致を見落とすことを防ぐとともに、探索効率の向上を図っている。

【 0 0 5 4 】

なお、上記の例では、上位コンピュータ 1 0 及び下位コンピュータ 2 0 が、それぞれ、割り当てられたクエリ文字列に従ってマッピングを行っているが、これに限られず、例えば、上位コンピュータ 1 0 のみが、割り当てられたクエリ文字列に基づいて、階層的インデックスを参照して、探索を行っても良い。

【 0 0 5 5 】

次に、上位コンピュータ 1 0 及び下位コンピュータ 2 0 は、それぞれ、後述する近似文字列照合のための文字列ペアを作成する (S 2 0 3) 。文字列ペアは、マッピングにより同定される一致文字列を含む、参照文字列における被照合文字列とクエリ文字列における照合文字列とからなる。

【 0 0 5 6 】

具体的には、上位コンピュータ 1 0 及び下位コンピュータ 2 0 は、マッピングにより同定された一致文字列の先頭及び末尾のそれぞれに、参照文字列における対応する所定長の文字列を追加することにより、被照合文字列を作成する。つまり、被照合文字列は、参照文字列において同定された一致文字列を含む該一致文字列近傍の文字列である。例えば、参照文字列が「 C C G A T C T G T A T A C C C T A C G A 」であって、一致文字列が「 T A C C 」である場合に、例えば前後 2 文字ずつ追加した文字列「 T A T A C C C T 」が被照合文字列となる。ヒトゲノムの解析の場合、参照文字列について、一致文字列の先頭及び末尾に追加する塩基の長さは、それぞれ、例えば 5 0 塩基程度であり得る。

【 0 0 5 7 】

また、上位コンピュータ 1 0 は、一致文字列の末尾にクエリ文字列における対応する所定長の文字列を追加することにより、照合文字列を作成する。ヒトゲノムの解析の場合、クエリ文字列について、一致文字列の末尾に追加する文字列の長さは、例えば 5 0 塩基程度である。

【 0 0 5 8 】

なお、本開示では、参照文字列について、一致文字列の先頭及び末尾に所定長の文字列が追加されるものとしたが、これに限られず、例えば、先頭又は末尾の一方にのみ所定長の文字列が追加されても良い。また、クエリ文字列について、一致文字列の先頭及び末尾

10

20

30

40

50

に、それぞれ、所定長の文字列を追加するようにしても良いし、或いは、文字列を追加せずに、一致文字列そのものを照合文字列として扱っても良い。

【 0 0 5 9 】

次に、上位コンピュータ 10 及び下位コンピュータ 20 は、参照文字列に基づく被照合文字列とクエリ文字列に基づく照合文字列とからなる文字列ペアに基づいて少なくとも 1 つの近似文字列を導出する (S 2 0 4)。近似文字列の導出には、動的計画法を用いた所定のアラインメントが適用される。アラインメントは、2 つの配列 (文字列) をその要素どうしで置換、挿入及び欠損を許容しつつ比較して、定義されたスコア / ペナルティに従って変異度 (類似度) を算出する手法である。アラインメントを実現するアルゴリズムとしては、Smith-Waterman アルゴリズムや Needleman-Wunsch アルゴリズムが知られている。また、動的計画法とは、ある段階で得られた最適解に基づいて次の段階の最適解を算出する手法である。つまり、変異度の算出では、動的計画法に従って、配列状のアラインメント表の各要素に対してスコアが算出され、最大スコアを持つ要素が決定され、これにより、少なくとも 1 つ以上の近似文字列が導出される。

10

【 0 0 6 0 】

以上のように、本実施形態の近似文字列照合では、参照文字列に基づく階層的インデックスが作成された後、与えられたクエリ文字列に従って、該階層的インデックスを探索することにより一致文字列 (及びその長さ) が同定され、同定された一致文字列に基づく被照合文字列と照合文字列とからなる文字列ペアに対して近似文字列照合がなされることにより、近似文字列が導出される。

20

【 0 0 6 1 】

図 3 は、本発明の一実施形態に係るコンピュータシステムによるインデックス作成処理の一例を説明するフローチャートである。すなわち、図 3 は、図 2 に示した階層的インデックスの作成処理 (S 2 0 1) の詳細を示している。また、図 4 は、階層的インデックスを作成するための参照文字列の一例を示す図であり、図 4 A ~ 4 C は、参照文字列に基づく階層的インデックスの作成過程におけるデータ配列構造の一例を示す図である。更に、図 5 は、階層的インデックスの作成過程におけるキーの出現開始位置及び出現回数を示すテーブル構造の一例を示す図であり、図 6 は、階層的インデックスを説明するための図である。

【 0 0 6 2 】

まず、図 3 に示すように、上位コンピュータ 10 は、階層的インデックスを作成するための参照文字列を受信する (S 3 0 1)。例えば、ヒトゲノム参照配列であれば、上位コンピュータ 10 は、データベース 30 にアクセスし、格納されているヒトゲノム参照配列を読み出す。以下では、理解容易のため、4 種類の文字「A」、「C」、「G」、及び「T」から構成される参照文字列「CCGATCTGTATACCCCTACGA」を例にして説明する (図 4 参照)。

30

【 0 0 6 3 】

上位コンピュータ 10 は、受信した参照文字列について、所定のキー長に従って各部分文字列 (すなわち、キー) を切り出して、キー配列を作成する (S 3 0 2)。例えば、キー長が「2」である場合、キー配列は、図 4 A (a) のようになる。同図中、左端の番号は、配列番号である。また、参照文字列の末端である 19 番目の「A」で始まるキーは、説明の簡略化のため、ここでは省略している。

40

【 0 0 6 4 】

次に、上位コンピュータ 10 は、作成したキー配列の各キーについて、所定のハッシュ関数を用いてハッシュ値を算出し、これを該キーに割り当てる (S 3 0 3)。これにより、キー配列は、図 4 A (b) のようになる。本開示におけるハッシュ関数は、4 種類の文字「A」、「C」、「G」、及び「T」にそれぞれ割り当てた「0」~「3」の数値により、4 進数で表現した値を出力する関数として定義されるが、これに限られない。例えば、1 番目のキー「CG」については、ハッシュ値は、 $h(CG) = 1 \times 4 + 2 = 6$ となり、また、11 番目のキー「AC」については、ハッシュ値は、 $h(AC) = 0 \times 4 + 1 =$

50

1 となる。

【 0 0 6 5 】

次に、上位コンピュータ 10 は、キー配列の各キーを、割り当てたハッシュ値に従って、例えば昇順にソートする (S 3 0 4)。これにより、キー配列は、図 4 A (c) のようになる。本例では、ソート後の各キー配列は、ソート前の配列番号を含み得る。例えば、図 4 A (c) 中、ソート後のキー配列における 0 番目のキー「AC」は、ソート前の(元の)配列番号「11」を保持し、また、ソート後の 1 番目のキー「AC」は、ソート前の配列番号「16」を保持している。

【 0 0 6 6 】

なお、上記の例では、切り出された各キーについて、算出したハッシュ値を割り当てて、ソートするものとしているが、これに限られない。例えば、参照文字列から切り出されるキーに拘わらず、参照文字列に現れる全ての文字の組み合わせに基づいてハッシュ値を算出して割り当てたキー配列を用意し、切り出されるキーに対応する出現開始位置を割り当てても良い。すなわち、4 種類の文字「A」、「C」、「G」、及び「T」から構成される参照文字列において、例えば、キーの長さが 2 文字であれば、 4^2 個の要素を有するキー配列がまず作成され、更に、ハッシュ値がそれぞれ割り当てられる。続いて、切り出されたキーは、参照文字列における出現開始位置とともに、作成されたキー配列における対応する要素(同じキーの要素)に割り当てられることにより、図 4 A (c) に示すようなキー配列が得られる。

【 0 0 6 7 】

次に、上位コンピュータ 10 は、現在のキー配列における各キーの出現開始位置及び出現回数を同定する (S 3 0 5)。図 5 (a) は、図 4 A (c) に示されるキー配列における各キーの出現開始位置及び出現回数を示している。例えば、キー「AC」は、キー配列において、出現開始位置「0」(配列番号「0」)を基点にして 2 回出現することが示されている。また、キー「CC」は、出現開始位置「4」を基点にして 3 回出現することが示されている。なお、ここでは、各文字どうしの全ての組み合わせからなるキーのパターン(すなわち、16 パターン)に対するその出現開始位置及び出現回数が見示されており、例示したキー配列に含まれていない例えばキー「AA」については、出現開始位置「-」及び出現回数「0」のように示されている。

【 0 0 6 8 】

次に、上位コンピュータ 10 は、各キーについて、その出現回数が所定の許容値を超えているか否かを判断する (S 3 0 6)。本開示において、所定の許容値は、階層的インデックスにおいて同じキーが重複して存在し得ることを許容する値である。本例では、所定の許容値は「1」としている。つまり、所定の許容値が「1」であれば、階層的インデックスにおいて各キーは唯一の存在となる。また、所定の許容値が大きいほど、階層的インデックスにおいて重複したキーが存在する可能性が高くなる一方、階層的インデックスの作成は高速化される。上位コンピュータ 10 は、出現回数が所定の許容値を超えているキーがあると判断する場合 (S 3 0 6 の Yes)、そのキーに対して追加キーを追加する (S 3 0 8)。

【 0 0 6 9 】

追加キーは、元の文字列における該キーに続く 1 以上の文字である。本例では、追加キーは 1 文字としている。追加キーの追加により得られる部分文字列は、新たなキーとみなされる。以下では、追加キーが追加された新たなキーを元のキーと区別するために「拡張キー」と称し、その配列を拡張キー配列と称する場合がある。追加キーが追加されることにより、各キーどうしが異なるものとして識別されることになる。図 4 B (a) は、元のキーに 1 個の追加キーが追加された拡張キーからなる拡張キー配列の一例を示している。なお、配列番号「12」の拡張キー「GA」については、元の文字列において「A」に続く文字がないため、終端文字列として例えば「\$」を割り当てている。また、配列番号「13」、「17」及び「18」のキーについては、その出現回数が 1 回であるため、追加キーは追加されない。

10

20

30

40

50

【 0 0 7 0 】

次に、上位コンピュータ10は、各キー（すなわち、拡張キー）を、図4B（b）に示すように、該追加キーに従って例えば昇順に更にソートする（S308）。この場合、元のキーのソート順が優先される。同図中、例えば、配列番号「2」及び「3」のキー配列「AT」については、追加キーによるソートで、その順序が入れ替わっていることがわかる。続いて、上位コンピュータ10は、処理S306に戻り、全てのキーの出現回数が所定の許容値を超えなくなるまで上記の処理を繰り返す。

【 0 0 7 1 】

すなわち、上位コンピュータ10は、キー配列における各キーの出現開始位置及び出現回数を同定し（S305）、出現回数が所定の許容値を超えているキーがないと判断する場合（S306のNo）、所望の階層的インデックスが作成されたため、処理を終了する。

10

【 0 0 7 2 】

図5は、図4に示される参照文字列に基づく階層的インデックスの作成過程におけるキーの出現開始位置及び出現回数を説明するための図である。例えば、図5（a）において出現回数が2以上であるキーには、追加キーが追加され（図4B（a））、各キー（拡張キー）は、追加キーに従ってソートされる（図4B（b））。これにより、図5（b）に示されるように、現在のキー配列における各キーの出現開始位置及び出現回数が同定される。同図に示す例では、拡張キー「CGA」及び「TAC」の出現回数が2となっている。したがって、これらの拡張キーのそれぞれについて、同様に、追加キーが追加されソートされる（図4C（a）及び（b））。なお、元の配列番号「17」（ハッシュ値でソート後の配列番号「8」）のキー「CG」については、キー「A」に続く文字がないため、終端文字列として例えば「\$」が割り当てられている。これにより、図5（c）に示されるように、該キーの出現開始位置及び出現回数が同定される。以上により、全ての拡張キーは、その出現回数が「1」となったため、インデックスの作成処理が終了する。

20

【 0 0 7 3 】

一例として、キー「TA」について考える。キー「TA」は、配列番号「14」～「16」にあることから（図4A（c）参照）、図6（a）に示すように、その出現開始位置は「14」、出現回数は「3」となる。次に、キー「TA」に追加キーが追加されソートされることより（図4B（b）参照）、図6（b）に示すように、追加キー「C」を含む拡張キー「TAC」については、その出現開始位置は「14」、出現回数は「2」となる一方、追加キー「T」を含む拡張キー「TAT」については、その出現開始位置は「16」、出現回数は「1」となる。したがって、出現回数が「2」のキー「TAC」について、更に追加キー「C」及び「G」がそれぞれ追加され、これにより、図6（c）に示すように、キー「TAG」の出現回数は「1」となる。このようにして、拡張キー配列は、階層的なツリー構造として把握される。

30

【 0 0 7 4 】

以上のようにして、上位コンピュータ10は、例えばヒトゲノム参照配列に基づいて、階層的インデックスを作成する。このような階層的インデックスは、ハッシュ値に従ってソートされているため、特定のキー（部分文字列）に関連する階層的インデックスの部分的なデータ配列構造は、メインメモリにおける特定のアドレス領域に集約的に展開され（データのシーケンシャル化）、これにより、メインメモリに対するランダムアクセスの回数を大幅に減らすことができるようになる。

40

【 0 0 7 5 】

なお、上記では、簡単化のため、極めて短い文字列を例にして説明したが、例えば、ヒトゲノムを扱う場合には、部分文字列のキー長を10～20程度、追加キーのキー長を2～8、所定の許容値を5～40とすることが好ましく、部分文字列のキー長を10～15程度、追加キーのキー長を2～4、所定の許容値を10～20とすることがより好ましい。

【 0 0 7 6 】

50

図7は、本発明の一実施形態に係るコンピュータシステムによるクエリ文字列に基づく参照文字列の探索処理の一例を説明するフローチャートである。すなわち、図7は、図2に示したクエリ文字列に基づく参照文字列の探索処理(S202)の詳細を示している。かかる探索処理により、クエリ文字列における所定のキーが参照文字列にマッピングされ、参照文字列における一致文字列及びその出現開始位置が同定される。なお、以下では、上位コンピュータ10による探索処理が説明されるが、並列的に動作する下位コンピュータ20による探索処理も同様である。

【0077】

同図に示すように、上位コンピュータ10は、データベース30から階層的インデックスを読み出して、メモリ上に展開し、記憶する(S701)。上位コンピュータ10は、高速化の観点から、階層的インデックスを構成するデータをメインメモリ上に連続的に展開し、記憶する。ここで、連続的に展開とは、データが同一バンクにおける連続的なメモリアドレスに配置されることを含む。

10

【0078】

次に、上位コンピュータ10は、参照文字列に対してマッピングを行うためのクエリ文字列を受信する(S702)。例えば、参照文字列がヒトゲノム参照配列であれば、クエリ文字列は、シーケンサから読み出される塩基対のデータ片である。以下では、理解容易のため、クエリ文字列は「TACC」であるものとして説明する。

【0079】

次に、上位コンピュータ10は、受信したクエリ文字列について、所定のキー長に従って各キーを連続的に切り出す(S703)。本例では、所定のキー長の初期値は「2」であるものとする。したがって、切り出される各キーは、「TA」、「AC」、「CC」、及び「C\$」となる。また、本例では、各キーのサンプリング間隔の初期値は1であるものとする。キーのサンプリング間隔とは、該キーに従って階層的インデックスを順に参照する開始位置の間隔である。なお、ヒトゲノムの解析であれば、サンプリング間隔の初期値は、例えば3~5程度であり得る。後述するように、サンプリング間隔は、キーの長さに応じて伸長される。

20

【0080】

続いて、上位コンピュータ10は、切り出した各キーについて、所定のハッシュ関数を用いてハッシュ値を算出する(S704)。上述したように、ハッシュ関数は、4種類の文字「A」、「C」、「G」、及び「T」にそれぞれ割り当てた「0」~「3」の数値により、4進数で表現した値を出力する関数として定義される。例えば、キー「TA」について、ハッシュ値は、 $h(TA) = 3 \times 4 + 0 = 12$ となる。

30

【0081】

次に、上位コンピュータ10は、算出したハッシュ値に従って、階層的インデックスを参照し(S705)、参照文字列における各キーの出現開始位置及び出現回数を同定する(S706)。例えば、キー「TA」であれば、ハッシュ値に従って、階層的インデックスにおける出現開始位置は14、出現回数は3であることが同定される(図6(a)参照)。

【0082】

次に、上位コンピュータ10は、各キーについて、その出現回数が所定の回数しきい値を超えているか否かを判断する(S707)。上述したように、本例では、所定の許容値は1としている。上位コンピュータ10は、キーの出現回数が所定の許容値を超えていると判断する場合(S707のYes)、続いて、現在のキー長が所定の上限値を超えているか否かを判断する(S708)。所定の上限値は、例えば、ヒトゲノムの解析であれば、20程度であり得るが、これに限られない。

40

【0083】

上位コンピュータ10は、各キーについて、現在のキー長が所定の上限値を超えていないと判断する場合(S708のNo)、該キーに対して追加キーを追加することにより、拡張キーを作成し(S709)、S706の処理に戻る。

50

【 0 0 8 4 】

例えば、上位コンピュータ10は、キー「T A」に対して追加キー「C」を追加して、その出現回数を調べ（図6（b）参照）、更に、追加キーが追加されたキー（拡張キー）に対して追加キー「C」を追加して、その出現回数を調べ（図6（c）参照）、出現回数が1になるまで繰り返す。

【 0 0 8 5 】

一方、上位コンピュータ10は、各キーについて、現在のキー長が所定の上限値を超えていると判断する場合（S708のYes）、該キーのサンプリング間隔を伸長（大きく）し（S710）、S705の処理に戻る。なお、サンプリング間隔の伸長は、キーについて、所定回数（例えば一回）だけ行われるようにしても良い。

10

【 0 0 8 6 】

また、上位コンピュータ10は、現在のキー（すなわち、拡張キー）について、その出現回数が所定の許容値を超えていないと判断する場合（S707のNo）、同定された参照文字列におけるキー（一致文字列）及びその出現開始位置を出力する（S711）。

【 0 0 8 7 】

以上のようにして、上位コンピュータ10は、階層的インデックスを用いて、参照文字列の中からクエリ文字列の少なくとも一部を探し出すことができる。とりわけ、本実施形態によれば、階層的インデックスの各キーはハッシュ値に従ってソートされているため、クエリ文字列に従った探索は、階層的インデックスにおける部分的・連続的なデータ配列構造に対して行われことになり、メインメモリに対するランダムアクセスの回数を大幅に減らすことができるようになる。

20

【 0 0 8 8 】

また、階層的インデックスの探索において、現在のキー（拡張キー）の長さが一定長以上になった場合に出現開始位置のサンプリング間隔を伸長するので、探索の回数が効率的に削減され、探索の高速化を図ることができる。一方で、サンプリング間隔の伸長は、本来同定されるべきキーを見逃す確率が高くなる可能性があるが、キーの長さが一定長以上になった場合にサンプリング間隔を長くしているので、このような見逃しの発生を効果的に抑制している。

【 0 0 8 9 】

次に、動的計画法を用いたアラインメントについて説明する。すなわち、上位コンピュータ10は、参照文字列における同定した出現開始位置近傍の文字列（被照合文字列）に対してクエリ文字列（照合文字列）がどれくらい変異しているか（変異度）を推定するために、動的計画法を用いたアラインメント処理を実行する。

30

【 0 0 9 0 】

アラインメントとは、2つの配列（文字列）をその要素どうしで置換、挿入及び欠損を許容しつつ比較して、定義されたスコア/ペナルティに従って変異度（類似度）を算出する手法である。例えば、文字列X（すなわち、参照文字列に基づく被照合文字列）と文字列Y（すなわち、クエリ文字列に基づく照合文字列）との比較において、文字列Yの一部の文字が置換され、挿入され、又は欠損する場合に、文字列Yは文字列Xに一致しないと判断される。つまり、文字列X及びYのそれぞれにおける各位置の文字どうしの関係は、一致する場合（一致）、一致しない場合（不一致）、及び一方の文字が存在しない場合（ギャップ）のいずれかであるといえる。ここで、一致する場合のスコアを例えば「+2」、不一致の場合のスコアを例えば「-1」、及びギャップの場合のスコアを例えば「-2」と定義する。そして、文字列X及び文字列Yの各要素（文字）どうしを比較して、これらのスコアを用いて各要素の変異度が算出される。本開示では、変異度の算出のために、グローバルアラインメントの一例であるNeedleman-Wunschアルゴリズムをベースにした改良アラインメントアルゴリズムが用いられる。以下、Needleman-Wunschアルゴリズムの基本的な考え方を説明し、更に、本発明に適用される改良アラインメントアルゴリズムを説明する。

40

【 0 0 9 1 】

50

例えば、文字列 $X = x_1 x_2 \dots x_i$ と文字列 $Y = y_1 y_2 \dots y_j$ との比較において、文字 x_i と文字 y_j との変異度 $F(i, j)$ は以下のように定義される。

【数 1】

$$F(i, j) = \max \left\{ \begin{array}{l} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{array} \right\} \quad \dots \text{式 1}$$

ここで、 \max は与えられた式の値の中から最大値を出力する関数、 s はスコア関数（一致： $s = 2$ 、不一致： $s = -1$ 、ギャップ： $s = -2$ ）、 d はギャップによるペナルティ（ $d = 2$ ）である。

10

【0092】

以下では、理解容易のため、被照合文字列「GCCTCGCT」と照合文字列「GCCATTC A」との間での動的計画法を用いたアラインメントを説明する。

【0093】

図 8 A は、本例における比較対象の文字列どうしを配列したアラインメント表である。表中、「」は空文字であり、 $F(0, 0) = 0$ とする。また、式 1 において、変異度 $F(i, j)$ の引数が負の値になる場合、範囲外であるため、計算は省略される（出力を `Null` とする。）。アラインメント表は、ある種のデータ構造としてメモリ上に展開され、プロセッサの利用に供される。

20

【0094】

まず、 $j = 0$ の場合において、 $F(1, 0)$ については、式 1 より、 \max 関数内のそれぞれは、

$$\begin{aligned} F(0, -1) + s(x_1, y_0) &= \text{Null} \\ F(0, 0) - d &= 0 - 2 = -2 \\ F(1, -1) - d &= \text{Null} \end{aligned}$$

であるから、

$$F(1, 0) = -2$$

となる。

30

【0095】

次に、 $F(2, 0)$ は、

$$F(2, 0) = F(1, 0) - d = -4$$

となる。同様にして、

$$\begin{aligned} F(n, 0) &= -nd \\ F(0, n) &= -nd \end{aligned}$$

となるため、アラインメント表は図 8 B に示すようになる。

【0096】

次に、 $j = 1$ の場合においても、同様に算出される。 $F(1, 1)$ については、 \max 関数内のそれぞれは、

40

$$\begin{aligned} F(0, 0) + s(x_1, y_1) &= 0 + 2 = 2 \\ F(0, 1) - d &= -2 - 2 = -4 \\ F(1, 0) - d &= -2 - 2 = -4 \end{aligned}$$

であり、これにより、

$$F(1, 1) = 2$$

となり、アラインメント表は図 8 C に示すようになる。

【0097】

以上の計算を同様に繰り返すことにより、図 8 D に示すようなアラインメント表が作成されることになる。作成されたアラインメント表において、要素位置（6, 7）の変異度 $F(6, 7)$ が最大値「7」を有している。したがって、同図に示すように、要素位置（

50

6, 7) から要素位置 (1, 1) までバックトラックがなされる。表中、右下斜め方向への矢印は文字の一致を示し、右方向への矢印は欠損を示し、下方向への矢印は挿入を示している。これにより、近似文字列「GCC<A>T[T]G」が導出されることになる。ただし、記号「< >」は、文字間への挿入を表し、記号「[]」は置換を表すものとする。つまり、参照文字列「GCCTCG」であるところ、クエリ文字列は、参照文字列の「C」と「T」の間に「A」が挿入され、参照文字列の「TCG」の「C」が「T」に置換されていることがわかる。

【0098】

以上のように、Needleman-Wunschアルゴリズムに従って、アラインメント表における変異度が最大値である要素位置を特定することにより、そこから近似文字列を導出することができる。しかしながら、Needleman-Wunschアルゴリズムではアラインメント表の全ての要素の変異度を算出するため、計算量が膨大となり、時間がかかっていた。そこで、本開示では、以下のような改良アラインメントアルゴリズムを提案し、これにより、計算量を削減し、処理の高速化を図っている。

【0099】

すなわち、本発明に適用される改良アラインメントアルゴリズムは、概略的には、アラインメント表における対角線上に位置する要素を中心とする所定の幅を有する計算領域を定め、該計算領域内の要素についてのみ変異度を算出することによりその最大値を決定し、該最大値が所定の条件を満たす場合に、該最大値に基づく要素から近似文字列を導出することを含む。最大値が所定の条件を満たさない場合には、計算領域が拡大され、同様に、変異度が算出されることによりその最大値を決定し、該最大値が所定の条件を満たすまで繰り返される。

【0100】

図9は、本発明の一実施形態に係るコンピュータシステムによる動的計画法を用いたアラインメント処理の一例を説明するフローチャートである。すなわち、図9は、図2に示したアラインメント処理の(S204)の詳細を示している。なお、以下では、上位コンピュータ10による一のクエリ文字列に基づくアラインメント処理が説明されるが、並列的に動作する下位コンピュータ20による他のクエリ文字列に基づくアラインメント処理も同様である。

【0101】

同図に示すように、上位コンピュータ10は、被照合文字列と照合文字列とからなる文字列ペアに基づいてアラインメント表を作成する(S901)。上述したように、アラインメント表は、ある種のデータ構造としてメモリ上に展開される。

【0102】

次に、上位コンピュータ10は、計算領域の幅 m (ただし、 m は正数)を初期値に設定する(S902)。幅 m の初期値は、例えば、マッピングにより得られた一致文字列の長さであり得るが、これに限られない。また、幅 m は、アラインメント表の対角線上の要素位置を中心とすることから、奇数の値に設定されるが、これに限られるものではない。これにより、アラインメント表の対角線上の要素位置を中心とする幅 m の計算領域が決定される。

【0103】

続いて、上位コンピュータ10は、計算領域の境界を画定する各要素に所定のダミー値を設定する(S903)。ダミー値は、変異度 F の値として十分に小さい値が選択される。例えば、ダミー値は、初期値の幅 m の例えば2~3倍程度の負の値であり、任意に設定することができる。

【0104】

次に、上位コンピュータ10は、計算領域内の各要素について、式1に従って変異度 F を算出する(S904)。続いて、上位コンピュータ10は、計算領域において最大値を有する最大変異度 F_{max} を決定し、その要素の位置を特定する(S905)。なお、最大変異度 F_{max} を持つ要素は、1つであるとは限らない。

10

20

30

40

50

【0105】

また、上位コンピュータ10は、アライメント表の行又は列における変異度Fの下限值 F_{Low} を算出する(S906)。下限値 F_{Low} は、該行又は列において、参照文字列とクエリ文字列とを比較して、m個の連続したギャップがあり、それ以外の部分は完全に又は実質的に一致したと仮定した場合の変異度Fの値である。すなわち、下限値 F_{Low} は、

$$F_{Low} = (\text{文字列の長さ} - m) \times s \quad \dots \text{式 2}$$

ただし、 $s = 2$ である。

で算出される。

【0106】

10

次に、上位コンピュータ10は、最大変異度 F_{max} と下限値 F_{Low} とを比較して、最大変異度 F_{max} が下限値 F_{Low} を超えているか否かを判断する(S907)。上位コンピュータ10は、最大変異度 F_{max} が下限値 F_{Low} を超えていないと判断する場合(S907のNo)、幅mを所定の大きさだけ拡幅する(S908)。例えば、 m は、 $m + 1$ とする(ただし、 m は文字列の文字数を超えないものとする)。

【0107】

上位コンピュータ10は、拡幅された幅mの計算領域に対して、同様に処理を行い、最大値 F_{max} が下限値 F_{Low} を超えるまで、上記処理を繰り返す。

【0108】

20

上位コンピュータ10は、最大値 F_{max} が下限値 F_{Low} を超えていると判断する場合(S907のYes)、該最大値を持つ要素の位置からバックトラックを行って、近似文字列を決定する(S909)。そして、上位コンピュータ10は、決定した近似文字列を出力する(S910)。

【0109】

例えば、被照合文字列「GGGATCCGATAATCGGTCCCCTAGG」(24文字)に対して照合文字列「GGGCATTCAACATAAGTCGGCCTG」(24文字)との間での、本発明に係るアラインメント法による変異度の算出例を説明する。なお、変異度Fの算出に式1を用いる点は、上述した例と同様である。なお、被照合文字列の長さ(24)と照合文字列の長さ(24)とは一致する必要はなく、典型的には、被照合文字列の長さの方が照合文字列の長さよりも長い。

30

【0110】

まず、上位コンピュータ10は、比較対象の文字列どうしを配列したアラインメント表を用意し、幅mの初期値を設定する。本例では、幅m(0)の初期値は7であるものとする。また、上位コンピュータ10は、幅mに従って規定される計算領域の境界部分の要素にダミー値を設定する。本例では、ダミー値は-20であるものとする。図10Aは、ダミー値が設定されたアラインメント表を示している。表中、ハッチングが描かれている要素が幅m(0) = 7での計算領域R(0)である。

【0111】

上位コンピュータ10は、計算領域R(0)内の各要素の変異度Fを式1に従って計算する。図10Bは、計算領域内の各要素の変異度Fが算出された状態を示している。上位コンピュータ10は、算出された変異度Fの中から、最大変異度 F_{max} を決定する。本例では、最大変異度 F_{max} は17である。図中、最大変異度 F_{max} が17である要素にはハッチングが示されている。

40

【0112】

続いて、上位コンピュータ10は、アラインメント表の行又は列における下限値 F_{Low} を算出する。本例では、下限値 F_{Low} は、

$$\begin{aligned} F_{Low} &= (24 - m) \times s \\ &= 17 \times 2 \\ &= 34 \end{aligned}$$

となる。

50

【 0 1 1 3 】

続いて、上位コンピュータ 10 は、最大変異度 F_{max} と下限値 F_{low} とを比較し、これにより、最大変異度 F_{max} が下限値 F_{low} を超えていないと判断するため、幅 m をだけ拡幅する。本例では、拡幅された幅 $m(1)$ を 15 に拡幅する。また、拡幅された幅 $m(1)$ の計算領域を $R(1)$ とする。

【 0 1 1 4 】

上位コンピュータ 10 は、同様にして、計算領域 $R(1)$ 内の各要素の変異度 F を式 1 に従って算出する。図 10 C は、計算領域内の各要素の変異度 F が算出された状態を示している。図中、計算領域 $R(1)$ に対して拡幅により追加された領域にハッチングが描かれている。これにより、最大変異度 F_{max} は 19 となる。また、このときの下限値 F_{low} は 18 となる。

10

【 0 1 1 5 】

したがって、上位コンピュータ 10 は、最大変異度 F_{max} が下限値 F_{low} を超えていると判断するため、最大変異度 $F_{max} = 19$ である要素からバックトラックし、これにより得られるパスに従って近似文字列を特定する。

【 0 1 1 6 】

すなわち、図 10 C に示す例では、上位コンピュータ 10 は、最大変異度 $F_{max} = 19$ である要素 (18, 22) を起点 (現在の要素位置) として、要素 (0, 0) 方向に向けて隣接する要素のうち変異度 F の値が最も大きい要素を同定し、そこに遷移する。したがって、変異度 $F = 17$ を持つ要素 (17, 21) が現在の要素位置となる。このような遷移を要素 (0, 0) まで繰り返すことにより、最終的に、近似文字列が導出される。なお、バックトラックにより得られるパスは、1 つとは限られず、複数である場合がある。

20

【 0 1 1 7 】

より具体的には、図 10 C に示すアラインメント表において、バックトラックにより得られるパスは 6 通りあり、各パスに従う近似文字列は、以下のとおりとなる。

(a) 第 1 のパス: $G G G < C > A < T > T C < A A > C - A T A A < G > T C G [G] C C$

(b) 第 2 のパス: $G G G < C > A < T > T C < A > [A] [C] A T A A < G > T C G [G] C C$

(c) 第 3 のパス: $G G G < C > A T [T] C < A > [A] [C] A T A A < G > T C G [G] C C$

30

(d) 第 4 のパス: $G G G < C > A T [T] C [A] < A C > A T A A < G > T C G [G] C C$

(e) 第 5 のパス: $G G G < C > A < T > T C < A A > C - A T A A < G > T C G [G] C C$

(f) 第 6 のパス: $G G G < C > A T < T > C < A > [A] [C] A T A A < G > T C G [G] C C$

ただし、記号「 $< >$ 」は、文字間への挿入を表し、記号「 $[]$ 」は文字の置換を表し、記号「 $-$ 」は文字の欠損を表すものとする。

なお、理解容易のため、上記の各パスに従う近似文字列を参照文字列との対比において図 11 A 及び 11 B に示している。

40

【 0 1 1 8 】

このように、改良されたアラインメントアルゴリズムでは、アラインメント表の全ての要素について変異度を算出するのではなく、所定の幅を有する計算領域を定め、必要に応じた範囲で計算領域を拡大させながら変異度を算出していくので、計算量を削減し、これにより、処理の高速化を図ることができるようになる。

【 0 1 1 9 】

以上のように、本実施形態によれば、参照文字列に基づく階層的インデックスが作成された後、与えられたクエリ文字列に従って、該階層的インデックスを探索することにより一致文字列 (及びその長さ) が同定され、同定された一致文字列に基づく被照合文字列と

50

照合文字列とからなる文字列ペアに対して近似文字列照合がなされることにより、近似文字列が導出される。

【0120】

上記各実施形態は、本発明を説明するための例示であり、本発明をこれらの実施形態にのみ限定する趣旨ではない。本発明は、その要旨を逸脱しない限り、さまざまな形態で実施することができる。

【0121】

例えば、本明細書に開示される方法においては、その結果に矛盾が生じない限り、ステップ、動作又は機能を並行して又は異なる順に実施しても良い。説明されたステップ、動作及び機能は、単なる例として提供されており、ステップ、動作及び機能のうちの一つか

10

【0122】

は、発明の要旨を逸脱しない範囲で、省略でき、また、互いに結合させることで一つのものとしてもよく、また、他のステップ、動作又は機能を追加してもよい。

また、本明細書では、さまざまな実施形態が開示されているが、一の実施形態における特定のフィーチャ（技術的事項）を、適宜改良しながら、他の実施形態に追加し、又は該他の実施形態における特定のフィーチャと置換することができ、そのような形態も本発明の要旨に含まれる。

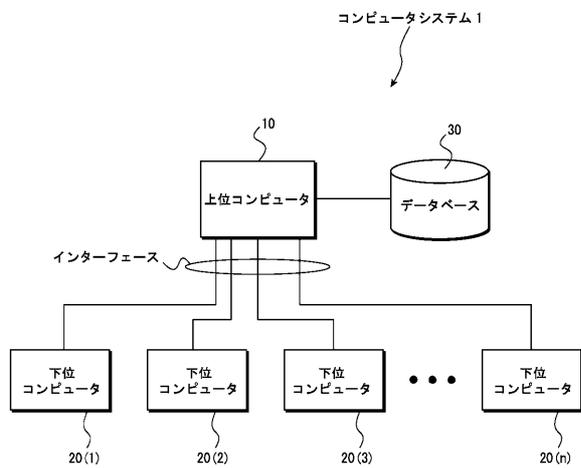
【符号の説明】

【0123】

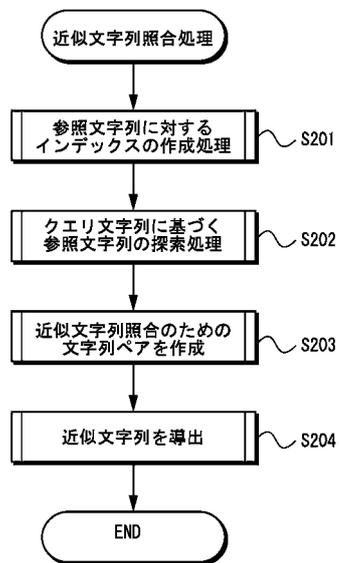
- 1 ... コンピュータシステム
- 10 ... 上位コンピュータ
- 20 ... 下位コンピュータ
- 30 ... データベース

20

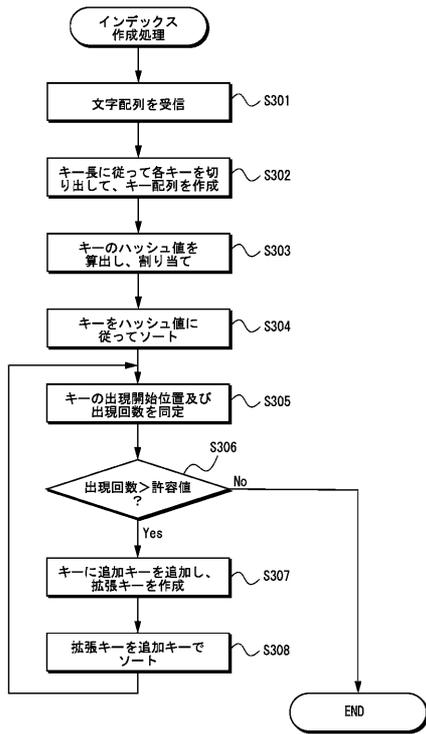
【図1】



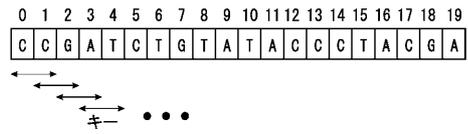
【図2】



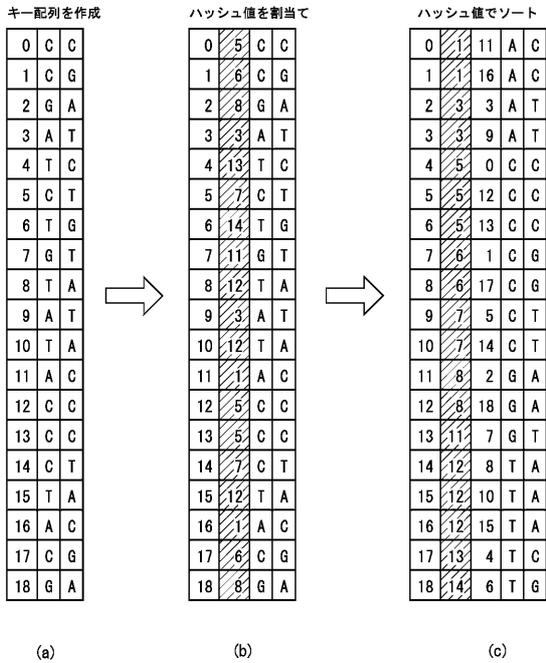
【 図 3 】



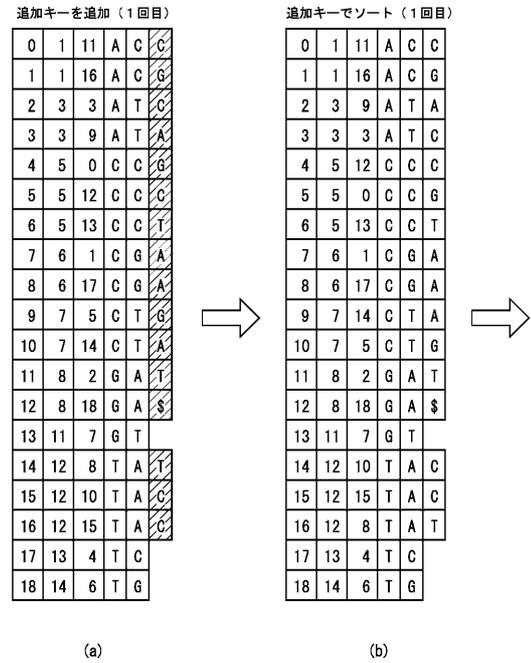
【 図 4 】



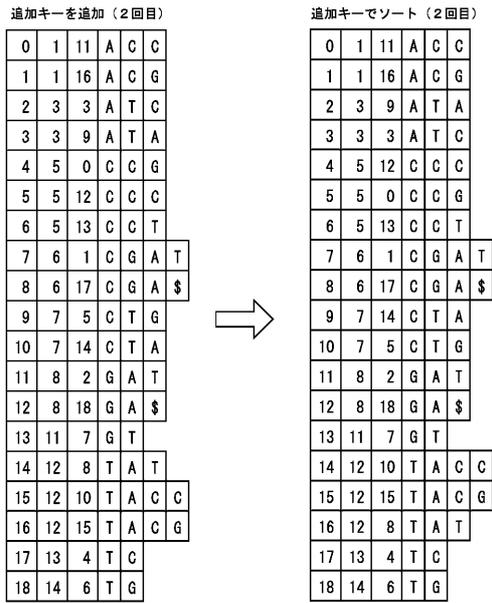
【 図 4 A 】



【 図 4 B 】



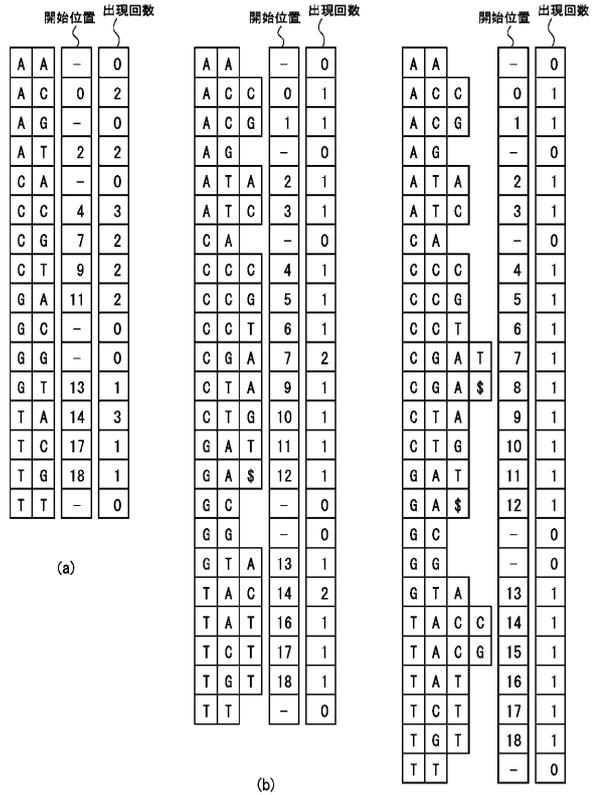
【図4C】



(a)

(b)

【図5】

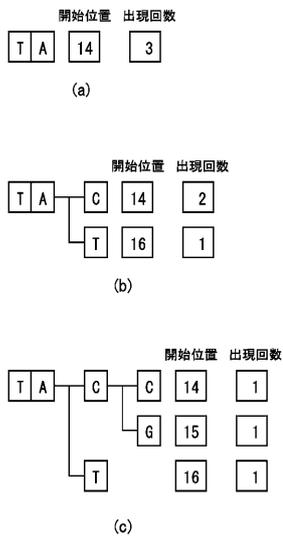


(a)

(b)

(c)

【図6】

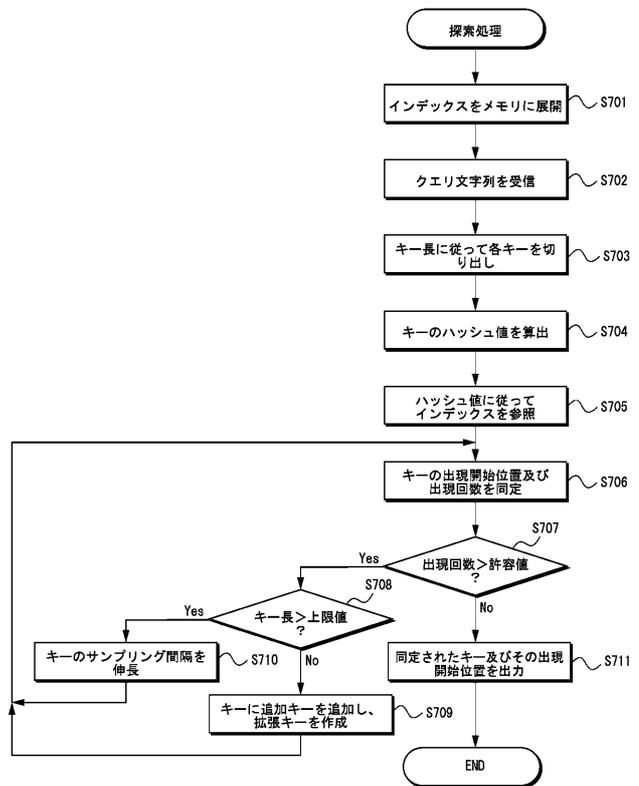


(a)

(b)

(c)

【図7】



【 図 8 A 】

参照文字列 (被照合文字列)

	i	0	1	2	3	4	5	6	7	8
j		φ	G	C	C	T	C	G	C	T
0	φ	0								
1	G									
2	C									
3	C									
4	A									
5	T									
6	T									
7	C									
8	A									

クエリ文字列 (照合文字列)

【 図 8 B 】

参照文字列 (被照合文字列)

	i	0	1	2	3	4	5	6	7	8
j		φ	G	C	C	T	C	G	C	T
0	φ	0	-2	-4	-6	-8	-10	-12	-14	-16
1	G	-2								
2	C	-4								
3	C	-6								
4	A	-8								
5	T	-10								
6	T	-12								
7	C	-14								
8	A	-16								

クエリ文字列 (照合文字列)

【 図 8 C 】

参照文字列 (被照合文字列)

	i	0	1	2	3	4	5	6	7	8
j		φ	G	C	C	T	C	G	C	T
0	φ	0	-2	-4	-6	-8	-10	-12	-14	-16
1	G	-2	2							
2	C	-4								
3	C	-6								
4	A	-8								
5	T	-10								
6	T	-12								
7	C	-14								
8	A	-16								

クエリ文字列 (照合文字列)

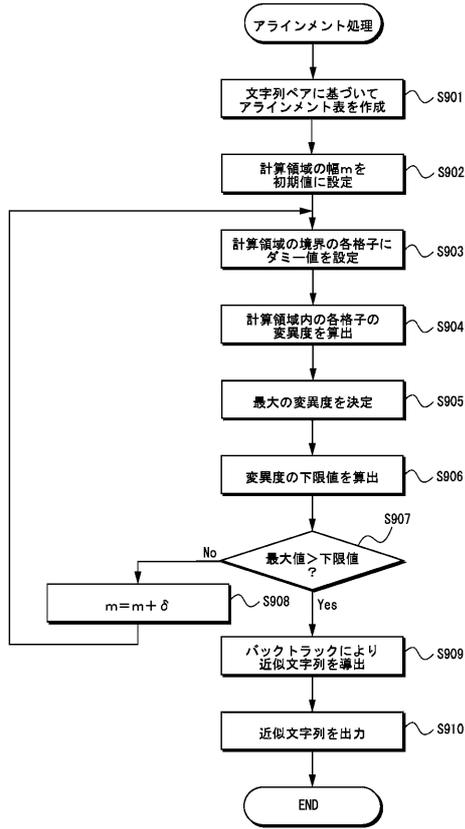
【 図 8 D 】

参照文字列 (被照合文字列)

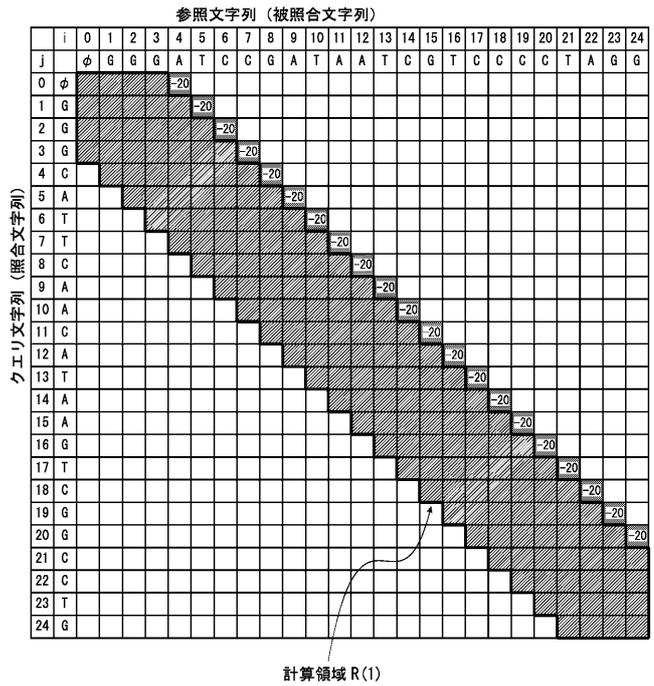
	i	0	1	2	3	4	5	6	7	8
j		φ	G	C	C	T	C	G	C	T
0	φ	0	-2	-4	-6	-8	-10	-12	-14	-16
1	G	-2	2	0	-2	-4	-6	-8	-10	-12
2	C	-4	0	4	2	0	-2	-4	-6	-8
3	C	-6	-2	2	6	4	2	0	-2	-4
4	A	-8	-4	0	4	5	3	1	-1	-3
5	T	-10	-6	-2	2	6	4	2	0	1
6	T	-12	-8	-4	0	4	5	3	1	2
7	C	-14	-10	-6	-2	2	3	7	5	3
8	A	-16	-12	-8	-4	0	1	5	6	4

クエリ文字列 (照合文字列)

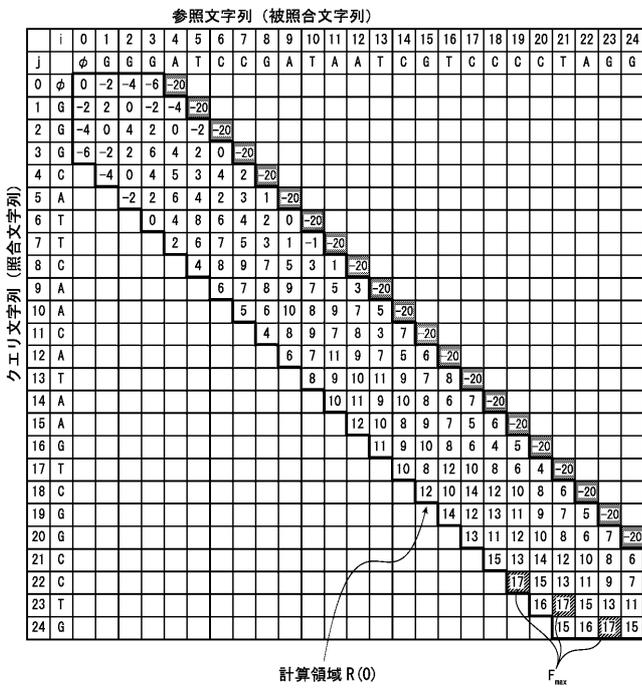
【図9】



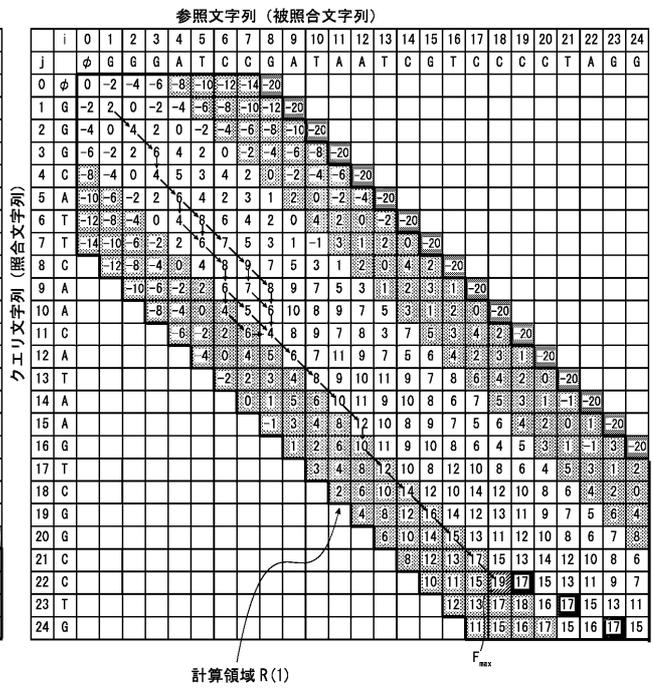
【図10A】



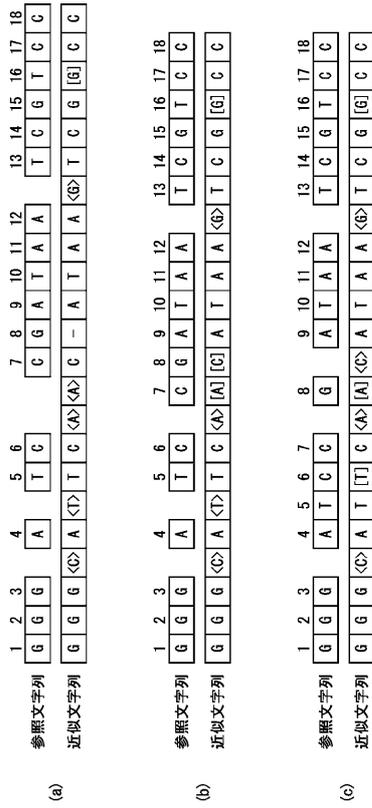
【図10B】



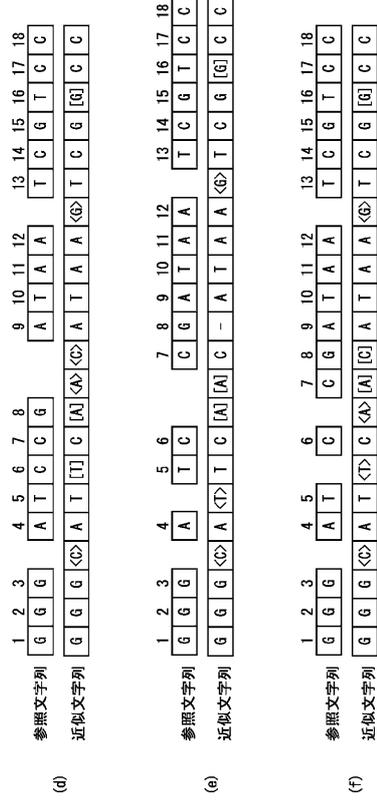
【図10C】



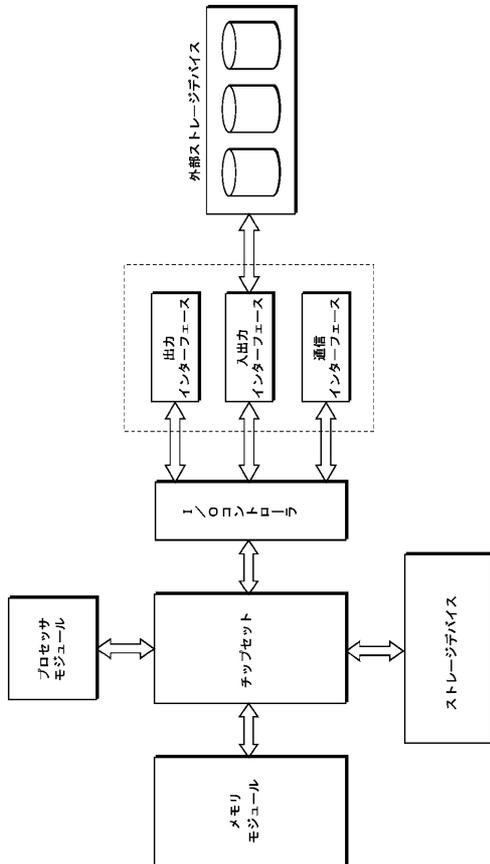
【 図 1 1 A 】



【 図 1 1 B 】



【 図 1 2 】



【手続補正書】

【提出日】令和5年1月26日(2023.1.26)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

コンピューティングデバイスに、クエリ文字列に基づいて参照文字列における近似文字列を検索するための方法を実行させるためのコンピュータプログラムであって、

前記方法は、

前記参照文字列に基づいて階層的インデックスを作成することと、

前記クエリ文字列の少なくとも一部と一致する前記参照文字列における部分文字列を同定するために、前記階層的インデックスを参照して、前記参照文字列に対する前記クエリ文字列のマッピングを行うことと、

前記のマッピングにより同定される前記部分文字列に基づいて作成される、前記参照文字列に基づく被照合文字列と前記クエリ文字列に基づく照合文字列とに基づいて、前記被照合文字列に近似する文字列を前記近似文字列として導出することと、を含み、

前記階層的インデックスを作成することは、

前記参照文字列から所定長の各第1のキーを切り出すことと、

切り出された前記各第1のキーについて、所定のハッシュ関数により該第1のキーに基づいて算出されるハッシュ値を割り当てた第1のキー配列を作成することと、

作成された前記第1のキー配列を更新することと、

更新された前記第1のキー配列を前記階層的インデックスとして出力することと、を含み、

前記第1のキー配列を更新することは、

前記第1のキー配列における前記各第1のキーについて、前記参照文字列における該第1のキーの出現回数を同定することと、

同定された前記第1のキーの前記出現回数に従って、該第1のキーに対して前記参照文字列における該第1のキーに続く少なくとも1以上の文字からなる第1の追加キーを追加することにより前記第1のキー配列を更新することと、を含む、
コンピュータプログラム。

【請求項2】

前記第1のキー配列を作成することは、前記ハッシュ値に従って前記第1のキー配列における前記各第1のキーをソートすることを含む、

請求項1に記載のコンピュータプログラム。

【請求項3】

前記第1のキー配列を更新することは、

前記同定された出現回数が所定の許容値を超えているか否かを判断することと、

前記同定された出現回数が前記所定の許容値を超えていると判断される場合に、前記第1のキーに対して前記参照文字列における該第1のキーに続く少なくとも1以上の文字からなる前記第1の追加キーを追加することにより新たな第1のキーを作成することと、

前記新たな第1のキーについて、前記参照文字列における該新たな第1のキーの出現回数を同定することと、を含む、

請求項1又は2に記載のコンピュータプログラム。

【請求項4】

前記第1のキー配列を更新することは、前記第1の追加キーに従って前記第1のキー配列における前記新たな第1のキーをソートすることを更に含む、

請求項1から3のいずれか一項に記載のコンピュータプログラム。

【請求項 5】

前記第 1 のキー配列を更新することは、前記同定された前記出現回数が所定の許容値を超えていないと判断されるまで、現在の前記第 1 のキーに新たな前記第 1 の追加キーを順次に追加することにより新たな前記第 1 のキーを作成することを含む、請求項 3 又は 4 に記載のコンピュータプログラム。

【請求項 6】

前記第 1 のキー配列を前記階層的インデックスとして出力することは、前記同定された前記出現回数が所定の許容値を超えていないと判断される場合に、現在の前記第 1 のキー配列を前記階層的インデックスとして出力することを含む、請求項 3 から 5 のいずれか一項に記載のコンピュータプログラム。

10

【請求項 7】

前記マッピングを行うことは、前記クエリ文字列から所定長の各第 2 のキーを切り出すことと、前記クエリ文字列から切り出された前記各第 2 のキーについて、前記所定のハッシュ関数により該第 2 のキーに基づいて算出されるハッシュ値を割り当てた第 2 のキー配列を作成することと、前記各第 2 のキーについて、前記ハッシュ値に従って、所定のサンプリング間隔で、前記階層的インデックスを参照し、該第 2 のキーの出現開始位置及び出現回数を同定することと、を含む、請求項 1 から 6 のいずれか一項に記載のコンピュータプログラム。

20

【請求項 8】

前記第 2 のキーの前記出現開始位置及び前記出現回数を同定することは、前記第 2 のキーの前記出現回数が前記所定の許容値を超えているか否かを判断することと、前記第 2 のキーの前記出現回数が前記所定の許容値を超えていると判断される場合に、前記第 2 のキーに対して前記クエリ文字列における該第 2 のキーに続く少なくとも 1 以上の文字からなる第 2 の追加キーを追加することにより新たな第 2 のキーを作成することと、前記第 2 のキーの前記出現回数が前記所定の許容値を超えていないと判断される場合に、同定された現在の前記第 2 のキーを前記部分文字列として出力するとともに該第 2 のキーの前記出現開始位置を出力することと、を含む、請求項 7 に記載のコンピュータプログラム。

30

【請求項 9】

前記第 2 のキーの前記出現開始位置及び前記出現回数を同定することは、前記第 2 のキーの前記同定された前記出現回数が前記所定の許容値を超えていないと判断されるまで、現在の前記第 2 のキーに新たな前記第 2 の追加キーを順次に追加して、前記新たな第 2 のキーを作成することを更に含む、請求項 8 に記載のコンピュータプログラム。

【請求項 10】

前記第 2 のキーの前記出現回数が前記所定の許容値を超えていると判断される場合に、該第 2 のキーの前記所定のサンプリング間隔を大きくする、請求項 8 又は 9 に記載のコンピュータプログラム。

40

【請求項 11】

前記近似文字列として導出することは、前記マッピングにより同定された前記部分文字列に基づく、前記被照合文字列と前記照合文字列とからなる文字列ペアを受信することと、前記文字列ペアに基づいて少なくとも 1 つの前記近似文字列を導出するために、所定のアラインメント処理を実行することと、導出された前記少なくとも 1 つの近似文字列を出力することと、を含む、請求項 8 から 10 のいずれか一項に記載のコンピュータプログラム。

50

【請求項 1 2】

前記所定のアラインメント処理を実行することは、
前記被照合文字列と前記照合文字列とに基づいて所定のアラインメント表を作成することと、
前記アラインメント表の対角線上の要素を中心にした幅 m を有する計算領域を設定することと、
設定された前記計算領域における各要素について、変異度を算出することと、
算出された前記変異度に基づいて、最大変異度を決定することと、
決定された前記最大変異度に基づいて、前記少なくとも 1 つの近似文字列を導出することを含む、
請求項 1 1 に記載のコンピュータプログラム。

10

【請求項 1 3】

前記所定のアラインメント処理を実行することは、
前記最大変異度と所定の下限値とを比較して、前記最大変異度が前記所定の下限値を超えているかを判断することと、
前記最大変異度が前記所定の下限値を超えていないと判断される場合に、新たな計算領域を設定するために、前記計算領域の前記幅 m を拡幅することと、
前記最大変異度が前記所定の下限値を超えていると判断される場合に、前記最大変異度を有する要素に基づいて、前記少なくとも 1 つの近似文字列を導出することと、を含み、
前記最大変異度が前記下限値を超えると判断されるまで、前記計算領域を拡幅することにより新たな計算領域を設定して前記変異度を算出することを繰り返す、
請求項 1 2 に記載のコンピュータプログラム。

20

【請求項 1 4】

前記所定の下限値は、所定の要素列に m 個の連続したギャップがあり、それ以外の部分は一致したと仮定した場合の変異度の値である、
請求項 1 3 に記載のコンピュータプログラム。

【請求項 1 5】

前記部分文字列に対して前記参照文字列における対応する所定の文字列を追加することにより前記被照合文字列を作成することと、
前記部分文字列に対して前記クエリ文字列における対応する所定の文字列を追加することにより前記照合文字列を作成することと、を更に含む、
請求項 1 1 から 1 4 のいずれか一項に記載のコンピュータプログラム。

30

【請求項 1 6】

コンピューティングデバイスに、クエリ文字列に基づいて参照文字列を探索するための階層的インデックスを作成する方法を実行させるためのコンピュータプログラムであって、

前記方法は、
前記参照文字列から所定長の各第 1 のキーを切り出すことと、
切り出された前記各第 1 のキーについて、所定のハッシュ関数により該第 1 のキーに基づいて算出されるハッシュ値を割り当てた第 1 のキー配列を作成することと、
作成された前記第 1 のキー配列を更新することと、
更新された前記第 1 のキー配列を前記階層的インデックスとして出力することと、を含み、

40

前記第 1 のキー配列を更新することは、
前記第 1 のキー配列における前記各第 1 のキーについて、前記参照文字列における該第 1 のキーの出現開始位置及び出現回数を同定することと、
同定された前記第 1 のキーの前記出現開始位置及び前記出現回数に従って、該第 1 のキーに対して前記参照文字列における該第 1 のキーに続く少なくとも 1 以上の文字からなる第 1 の追加キーを追加することにより前記第 1 のキー配列を更新することと、を含む、
コンピュータプログラム。

50

【請求項 17】

コンピューティングデバイスに、参照文字列に対してクエリ文字列のマッピングを行う方法を実行させるためのコンピュータプログラムであって、

前記方法は、

前記参照文字列から切り出される各第1のキー、及び前記各第1のキーについて、該第1のキーに基づいて所定のハッシュ関数により算出されるハッシュ値が割り当てられた第1のキー配列であって、前記参照文字列における前記第1のキーの出現開始位置及び出現回数に従って前記第1のキーに追加キーが追加された該第1のキー配列からなる階層的インデックスを読み出すことと、

前記クエリ文字列から所定のキー長を有する各第2のキーを切り出すことと、

前記クエリ文字列から切り出された前記各第2のキーについて、所定のハッシュ関数により該第2のキーに基づいて算出されるハッシュ値を割り当てた第2のキー配列を作成することと、

前記各第2のキーについて、前記ハッシュ値に従って、所定のサンプリング間隔で、前記階層的インデックスを参照し、前記第1のキー配列における前記第1のキーとの比較により、該第2のキーの出現開始位置及び出現回数を同定することと、

前記同定した出現回数が所定のしきい値を超えているか否かを判断することと、

前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記第2のキーに対して前記クエリ文字列における該第2のキーに続く少なくとも1以上の文字からなる追加キーを追加することにより新たな第2のキーを作成することと、

前記同定された前記出現回数が所定のしきい値を超えていないと判断される場合に、同定された現在の前記第2のキーの出現開始位置及び該キーを出力することと、を含み、

前記第2のキーの前記出現開始位置及び前記出現回数を同定することは、前記同定された前記出現回数が所定のしきい値を超えていないと判断されるまで、現在の前記第2のキーに新たな前記追加キーを順次に追加して、前記新たな第2のキーを作成することを含む、
コンピュータプログラム。

【請求項 18】

前記方法は、前記同定された前記出現回数が所定の許容値を超えていると判断される場合に、前記第2のキーの前記所定のサンプリング間隔を大きくするように構成される、請求項17に記載のコンピュータプログラム。

【請求項 19】

コンピューティングデバイスに、参照文字列における部分文字列とクエリ文字列との間の変異を所定のアラインメント処理により同定する方法を実行させるためのコンピュータプログラムであって、

前記方法は、

マッピングにより同定される前記部分文字列に基づいて作成される、前記参照文字列に基づく被照合文字列と前記クエリ文字列に基づく照合文字列とからなる文字列ペアを受信することと、

前記文字列ペアに基づいて、前記被照合文字列に近似する文字列を少なくとも1つの近似文字列として導出するために、所定のアラインメント処理を実行することと、

導出された前記少なくとも1つの近似文字列を出力することと、を含み、

前記所定のアラインメント処理を実行することは、

前記被照合文字列と前記照合文字列とに基づいて所定のアラインメント表を作成することと、

前記アラインメント表の対角線上の要素を中心にした幅mを有する計算領域を設定することと、

設定された前記計算領域における各要素について、変異度を算出することと、

算出された前記変異度に基づいて、最大変異度を決定することと、

決定された前記最大変異度に基づいて、前記少なくとも1つの近似文字列を導出するこ

10

20

30

40

50

とを含む、
コンピュータプログラム。

【請求項 20】

前記所定のアラインメント処理を実行することは、
前記最大変異度と所定の下限値とを比較して、前記最大変異度が前記所定の下限値を超えているかを判断することと、
前記最大変異度が前記所定の下限値を超えていないと判断される場合に、新たな計算領域を設定するために、前記計算領域の前記幅 m を拡幅することと、
前記最大変異度が前記所定の下限値を超えていると判断される場合に、前記最大変異度を有する要素に基づいて、前記少なくとも 1 つの近似文字列を導出することと、を含み、
前記最大変異度が前記下限値を超えると判断されるまで、前記計算領域を拡幅することにより新たな計算領域を設定して前記変異度を算出することを繰り返す、
請求項 19 に記載のコンピュータプログラム。

【請求項 21】

前記所定のアラインメント処理を実行することは、所定の要素列に m 個の連続したギャップがあり、それ以外の部分は一致したと仮定した場合の変異度の値を前記所定の下限値として設定することを更に含む、
請求項 20 に記載のコンピュータプログラム。

フロントページの続き

(72)発明者 牧野 淳一郎

神奈川県横浜市金沢区泥亀一丁目2番B-1312 先端加速システムズ株式会社内

(72)発明者 姫野 龍太郎

神奈川県横浜市金沢区泥亀一丁目2番B-1312 先端加速システムズ株式会社内

Fターム(参考) 5B175 DA01 DA10 FA01 HB01 KA04 KA08