

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2021-523479  
(P2021-523479A)

(43) 公表日 令和3年9月2日(2021.9.2)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G 1 6 B 30/20 (2019.01)</b>	G 1 6 B 30/20	4 B O 2 9
<b>C 1 2 Q 1/68 (2018.01)</b>	C 1 2 Q 1/68	4 B O 6 3
<b>C 1 2 M 1/34 (2006.01)</b>	C 1 2 M 1/34	Z

審査請求 未請求 予備審査請求 未請求 (全 59 頁)

(21) 出願番号 特願2020-564123 (P2020-564123)  
 (86) (22) 出願日 令和1年5月13日 (2019.5.13)  
 (85) 翻訳文提出日 令和3年1月5日 (2021.1.5)  
 (86) 国際出願番号 PCT/US2019/032065  
 (87) 国際公開番号 W02019/222120  
 (87) 国際公開日 令和1年11月21日 (2019.11.21)  
 (31) 優先権主張番号 62/671, 884  
 (32) 優先日 平成30年5月15日 (2018.5.15)  
 (33) 優先権主張国・地域又は機関 米国 (US)  
 (31) 優先権主張番号 62/671, 260  
 (32) 優先日 平成30年5月14日 (2018.5.14)  
 (33) 優先権主張国・地域又は機関 米国 (US)

(71) 出願人 516144164  
 クアンタム-エスアイ インコーポレイテッド  
 QUANTUM-S I INCORPORATED  
 アメリカ合衆国 06437 コネチカット州 ギルフォード オールド ウィットフィールド ストリート 530  
 (74) 代理人 100105957  
 弁理士 恩田 誠  
 (74) 代理人 100068755  
 弁理士 恩田 博宣  
 (74) 代理人 100142907  
 弁理士 本田 淳

最終頁に続く

(54) 【発明の名称】 機械学習可能な生物学的ポリマーアセンブリ

(57) 【要約】

本明細書には、高分子の生物学的ポリマーアセンブリを生成するための機械学習技術が記載されている。例えば、システムは、機械学習技術を使用して、生物のDNAのゲノムアセンブリ、生物のDNAの一部の遺伝子配列、またはタンパク質のアミノ酸配列を生成し得る。システムは、シーケンシングデバイスによって生成された生物学的ポリマー配列および配列から生成されたアセンブリにアクセスし得る。システムは、配列およびアセンブリを使用して機械学習モデルへの入力を生成し得る。システムは、入力を機械学習モデルに提供して、対応する出力を取得し得る。システムは、対応する出力を使用して、アセンブリ内の位置において生物学的ポリマーを同定し、次にアセンブリ内の位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、更新されたアセンブリを取得し得る。

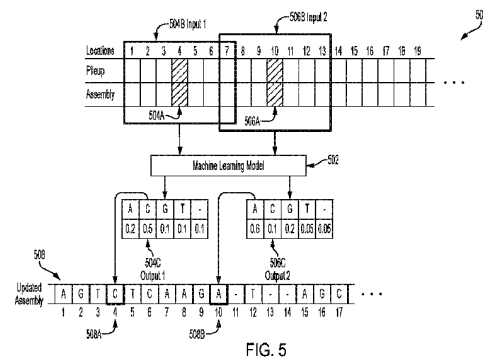


FIG. 5

**【特許請求の範囲】****【請求項 1】**

高分子の生物学的ポリマーアセンブリを生成する方法であって、  
少なくとも1つのコンピュータハードウェアプロセッサを使用して、  
複数の生物学的ポリマー配列と、個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリとにアクセスするステップと、

前記複数の生物学的ポリマー配列および前記アセンブリを使用して、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップと、

前記第1の入力を前記トレーニングされた深層学習モデルに提供して、第1の複数のアセンブリ位置の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がその位置に存在する1つまたは複数の尤度を示す対応する第1の出力を取得するステップと、

前記トレーニングされた深層学習モデルの前記第1の出力を使用して、前記第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップと、

前記第1の複数のアセンブリ位置において同定された生物学的ポリマーを示すように前記アセンブリを更新して、更新されたアセンブリを取得するステップとを実行するステップを含む方法。

**【請求項 2】**

前記高分子がタンパク質を含み、前記複数の生物学的ポリマー配列が複数のアミノ酸配列を含み、前記アセンブリが個々のアセンブリ位置におけるアミノ酸を示す、請求項1に記載の方法。

**【請求項 3】**

前記高分子が核酸を含み、前記複数の生物学的ポリマー配列が複数のヌクレオチド配列を含み、前記アセンブリが個々のアセンブリ位置におけるヌクレオチドを示す、請求項1または任意の他の先行する請求項に記載の方法。

**【請求項 4】**

前記アセンブリは、前記第1の複数のアセンブリ位置のうちの第1のアセンブリ位置における第1のヌクレオチドを示し、

前記第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップは、前記第1のアセンブリ位置において第2のヌクレオチドを同定することを含み、

前記アセンブリの更新することは、前記第1のアセンブリ位置において前記第2のヌクレオチドを示すように前記アセンブリを更新することを含む、請求項3または任意の他の先行する請求項に記載の方法。

**【請求項 5】**

前記アセンブリを更新して、前記更新されたアセンブリを取得した後、

前記複数のヌクレオチド配列を前記更新されたアセンブリに整列させるステップと、

前記複数のヌクレオチド配列および前記更新されたアセンブリを使用して、前記トレーニングされた深層学習モデルに提供される第2の入力を生成するステップと、

前記第2の入力を前記トレーニングされた深層学習モデルに提供して、第2の複数のアセンブリ位置の各々に関して、1つまたは複数の個々のヌクレオチドの各々がその位置に存在する1つまたは複数の尤度を示す対応する第2の出力を取得するステップと、

前記トレーニングされた深層学習モデルの前記第2の出力に基づいて、前記第2の複数のアセンブリ位置におけるヌクレオチドを同定するステップと、

前記第2の複数のアセンブリ位置において同定されたヌクレオチドを示すように前記更新されたアセンブリを更新して、第2の更新されたアセンブリを取得するステップとをさらに含む、請求項3または任意の他の先行する請求項に記載の方法。

**【請求項 6】**

前記複数のヌクレオチド配列を前記アセンブリに整列させることをさらに含む、請求項3または任意の他の先行する請求項に記載の方法。

**【請求項 7】**

前記複数のヌクレオチド配列が少なくとも9個のヌクレオチド配列を含む、請求項6ま

10

20

30

40

50

たは任意の他の先行する請求項に記載の方法。

【請求項 8】

前記トレーニングされた深層学習モデルへの前記第 1 の入力を生成するステップは、前記第 1 の複数のアセンブリ位置を選択すること、選択された第 1 の複数のアセンブリ位置に基づいて前記第 1 の入力を生成することを含む、請求項 3 または任意の他の先行する請求項に記載の方法。

【請求項 9】

前記アセンブリ内の前記第 1 の複数の位置を選択することは、前記アセンブリが前記第 1 の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定すること、  
決定された尤度を使用して、前記第 1 の複数のアセンブリ位置を選択することを含む、請求項 8 または任意の他の先行する請求項に記載の方法。

10

【請求項 10】

前記トレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、前記複数のヌクレオチド配列の個々の 1 つを前記アセンブリと比較することを含む、請求項 3 または任意の他の先行する請求項に記載の方法。

【請求項 11】

前記第 1 の複数のアセンブリ位置のうちの第 1 のアセンブリ位置におけるヌクレオチドを同定するために前記トレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、

20

前記第 1 のアセンブリ位置の近傍の 1 つまたは複数のアセンブリ位置の各々における複数のヌクレオチドの各々に関して、

ヌクレオチドがその位置にあることを示す複数のヌクレオチド配列の数を示すカウントを決定すること、

前記アセンブリがその位置においてヌクレオチドを示しているかどうかに基づいて参照値を決定すること、

前記カウントと前記参照値との差異を示すエラー値を決定すること、

前記参照値および前記エラー値を前記第 1 の入力に含ませることを含む、請求項 3 または任意の他の先行する請求項に記載の方法。

【請求項 12】

30

前記アセンブリがその位置においてヌクレオチドを示しているかどうかに基づいて前記参照値を決定することは、

前記アセンブリがその位置においてヌクレオチドを示している場合、前記参照値が第 1 の値であると決定すること、

前記アセンブリがその位置においてヌクレオチドを示していない場合、前記参照値が第 2 の値であると決定することを含む、請求項 11 または任意の他の先行する請求項に記載の方法。

【請求項 13】

前記第 1 の値は、前記複数のヌクレオチド配列の数であり、

前記第 2 の値は 0 である、請求項 12 または任意の他の先行する請求項に記載の方法。

40

【請求項 14】

前記トレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、複数の列を有するデータ構造に値を配置することを含み、

第 1 の列は、第 1 のアセンブリ位置において複数のヌクレオチドに対して決定された参照値およびエラー値を保持し、

第 2 の列は、前記第 1 のアセンブリ位置の近傍にある 1 つまたは複数のアセンブリ位置のうちの第 2 のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する、請求項 11 または任意の他の先行する請求項に記載の方法。

【請求項 15】

前記第 1 のアセンブリ位置の近傍の 1 つまたは複数のアセンブリ位置が、前記第 1 のア

50

センブリ位置とは別個の少なくとも2つのアセンブリ位置を含む、請求項11または任意の他の先行する請求項に記載の方法。

【請求項16】

1つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する1つまたは複数の尤度が、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置において存在する尤度を含み、

前記第1の複数のアセンブリ位置における生物学的ポリマーを同定することは、第1のヌクレオチドが第1の位置に存在する尤度が複数のヌクレオチドのうちの第2のヌクレオチドが第1のアセンブリ位置に存在する尤度よりも大きいことを決定することによって前記第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドが複数のヌクレオチドのうちの第1のヌクレオチドであることを決定することを含む、請求項3または任意の他の先行する請求項に記載の方法。

10

【請求項17】

前記複数のヌクレオチド配列から前記アセンブリを生成することをさらに含む、請求項3または任意の他の先行する請求項に記載の方法。

【請求項18】

前記複数のヌクレオチド配列から前記アセンブリを生成することは、前記アセンブリとなる複数のヌクレオチド配列からコンセンサス配列を決定することを含む、請求項17または任意の他の先行する請求項に記載の方法。

【請求項19】

前記複数のヌクレオチド配列から前記アセンブリを生成することは、オーバーラップ・レイアウト・コンセンサス(OLC)アルゴリズムを前記複数のヌクレオチド配列に適用することを含む、請求項17または任意の他の先行する請求項に記載の方法。

20

【請求項20】

参照高分子のシーケンシングから取得された生物学的ポリマー配列と、前記参照高分子の所定のアセンブリとを含むトレーニングデータにアクセスするステップと、

前記トレーニングデータを使用して深層学習モデルをトレーニングして、トレーニングされた深層学習モデルを取得するステップとをさらに含む、請求項1または任意の他の先行する請求項に記載の方法。

【請求項21】

前記参照高分子は前記高分子とは異なる、請求項20または任意の他の先行する請求項に記載の方法。

30

【請求項22】

深層学習モデルが畳み込みニューラルネットワーク(CNN)を含む、請求項1または任意の他の先行する請求項に記載の方法。

【請求項23】

高分子の生物学的ポリマーアセンブリを生成するためのシステムであって、

少なくとも1つのコンピュータハードウェアプロセッサと、

命令を格納する少なくとも1つの非一時的なコンピュータ可読記憶媒体とを備え、前記命令は、前記少なくとも1つのコンピュータハードウェアプロセッサによる実行時に、前記少なくとも1つのコンピュータハードウェアプロセッサに、

40

複数の生物学的ポリマー配列と、個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリとにアクセスするステップと、

前記複数の生物学的ポリマー配列および前記アセンブリを使用して、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップと、

前記第1の入力を前記トレーニングされた深層学習モデルに提供して、第1の複数のアセンブリ位置の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がその位置に存在する1つまたは複数の尤度を示す対応する第1の出力を取得するステップと、

前記トレーニングされた深層学習モデルの前記第1の出力を使用して、前記第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップと、

50

前記第 1 の複数のアセンブリ位置において同定された生物学的ポリマーを示すように前記アセンブリを更新して、更新されたアセンブリを取得するステップとを実行させる、システム。

【請求項 2 4】

前記高分子がタンパク質を含み、前記複数の生物学的ポリマー配列が複数のアミノ酸配列を含み、前記アセンブリが個々のアセンブリ位置におけるアミノ酸を示す、請求項 2 3 に記載のシステム。

【請求項 2 5】

前記高分子が核酸を含み、前記複数の生物学的ポリマー配列が複数のヌクレオチド配列を含み、前記アセンブリが個々のアセンブリ位置におけるヌクレオチドを示す、請求項 2 3 または任意の他の先行する請求項に記載のシステム。

10

【請求項 2 6】

前記アセンブリは、前記第 1 の複数のアセンブリ位置のうちの第 1 のアセンブリ位置における第 1 のヌクレオチドを示し、

前記第 1 の複数のアセンブリ位置における生物学的ポリマーを同定するステップは、前記第 1 のアセンブリ位置において第 2 のヌクレオチドを同定することを含み、

前記アセンブリを更新することは、前記第 1 のアセンブリ位置において前記第 2 のヌクレオチドを示すように前記アセンブリを更新することを含む、請求項 2 5 または任意の他の先行する請求項に記載のシステム。

【請求項 2 7】

前記命令は、前記アセンブリを更新して、前記更新されたアセンブリを取得した後、前記少なくとも 1 つのコンピュータハードウェアプロセッサに、

前記複数のヌクレオチド配列を前記更新されたアセンブリに整列させるステップと、

前記複数のヌクレオチド配列および前記更新されたアセンブリを使用して、前記トレーニングされた深層学習モデルに提供される第 2 の入力を生成するステップと、

前記第 2 の入力を前記トレーニングされた深層学習モデルに提供して、第 2 の複数のアセンブリ位置の各々に関して、1 つまたは複数の個々のヌクレオチドの各々がその位置に存在する 1 つまたは複数の尤度を示す対応する第 2 の出力を取得するステップと、

前記トレーニングされた深層学習モデルの前記第 2 の出力に基づいて、前記第 2 の複数のアセンブリ位置におけるヌクレオチドを同定するステップと、

20

30

前記第 2 の複数のアセンブリ位置において同定されたヌクレオチドを示すように前記更新されたアセンブリを更新して、第 2 の更新されたアセンブリを取得するステップとをさらに実行させる、請求項 2 5 または任意の他の先行する請求項に記載のシステム。

【請求項 2 8】

前記命令は、前記少なくとも 1 つのコンピュータハードウェアプロセッサに、前記複数のヌクレオチド配列を前記アセンブリに整列させるステップを実行させる請求項 2 5 または任意の他の先行する請求項に記載のシステム。

【請求項 2 9】

前記複数のヌクレオチド配列が少なくとも 9 個のヌクレオチド配列を含む、請求項 2 8 または任意の他の先行する請求項に記載のシステム。

40

【請求項 3 0】

前記トレーニングされた深層学習モデルへの前記第 1 の入力を生成するステップは、前記第 1 の複数のアセンブリ位置を選択すること、

選択された第 1 の複数のアセンブリ位置に基づいて前記第 1 の入力を生成することを含む、請求項 2 5 または任意の他の先行する請求項に記載のシステム。

【請求項 3 1】

前記アセンブリ内の前記第 1 の複数の位置を選択することは、

前記アセンブリが前記第 1 の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定すること、

決定された尤度を使用して、前記第 1 の複数のアセンブリ位置を選択することを含む、

50

請求項 30 または任意の他の先行する請求項に記載のシステム。

【請求項 32】

前記トレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、前記複数のヌクレオチド配列の個々の 1 つを前記アセンブリと比較することを含む、請求項 25 または任意の他の先行する請求項に記載のシステム。

【請求項 33】

前記第 1 の複数のアセンブリ位置のうちの第 1 のアセンブリ位置におけるヌクレオチドを同定するためにトレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、

前記第 1 のアセンブリ位置の近傍の 1 つまたは複数のアセンブリ位置の各々における複数のヌクレオチドの各々に関して、

ヌクレオチドがその位置にあることを示す複数のヌクレオチド配列の数を示すカウントを決定すること、

前記アセンブリがその位置においてヌクレオチドを示しているかどうかに基づいて参照値を決定すること、

前記カウントと前記参照値との差異を示すエラー値を決定すること、

前記参照値および前記エラー値を前記第 1 の入力に含ませることを含む、請求項 25 または任意の他の先行する請求項に記載のシステム。

【請求項 34】

前記アセンブリがその位置においてヌクレオチドを示すかどうかに基づいて前記参照値を決定することは、

前記アセンブリがその位置においてヌクレオチドを示している場合、前記参照値が第 1 の値であると決定すること、

前記アセンブリがその位置においてヌクレオチドを示していない場合、前記参照値が第 2 の値であると決定することを含む、請求項 33 または任意の他の先行する請求項に記載のシステム。

【請求項 35】

前記第 1 の値は、前記複数のヌクレオチド配列の数であり、

前記第 2 の値は 0 である、請求項 34 または他の先行する請求項に記載のシステム。

【請求項 36】

前記トレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、複数の列を有するデータ構造に値を配置することを含み、

第 1 の列は、第 1 のアセンブリ位置において複数のヌクレオチドに対して決定された参照値およびエラー値を保持し、

第 2 の列は、前記第 1 のアセンブリ位置の近傍にある 1 つまたは複数のアセンブリ位置のうちの第 2 のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する、請求項 33 または任意の他の先行する請求項に記載のシステム。

【請求項 37】

前記第 1 のアセンブリ位置の近傍の 1 つまたは複数のアセンブリ位置が、前記第 1 のアセンブリ位置とは別の少なくとも 2 つのアセンブリ位置を含む、請求項 33 または任意の他の先行する請求項に記載のシステム。

【請求項 38】

1 つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する 1 つまたは複数の尤度が、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置において存在する尤度を含み、

前記第 1 の複数のアセンブリ位置における生物学的ポリマーを同定することは、第 1 のヌクレオチドが第 1 の位置に存在する尤度が複数のヌクレオチドのうちの第 2 のヌクレオチドが第 1 のアセンブリ位置に存在する尤度よりも大きいことを決定することによって前記第 1 の複数のアセンブリ位置のうちの第 1 のアセンブリ位置におけるヌクレオチドが複

10

20

30

40

50

数のヌクレオチドのうちの第1のヌクレオチドであることを決定することを含む、請求項25または任意の他の先行する請求項に記載のシステム。

【請求項39】

前記命令は、前記少なくとも1つのコンピュータハードウェアプロセッサに、前記複数のヌクレオチド配列から前記アセンブリを生成することを実行させる、請求項25または任意の他の先行する請求項に記載のシステム。

【請求項40】

前記複数のヌクレオチド配列から前記アセンブリを生成することは、前記アセンブリとなる複数のヌクレオチド配列からコンセンサ配列を決定することを含む、請求項39または任意の他の先行する請求項に記載のシステム。

10

【請求項41】

前記複数のヌクレオチド配列から前記アセンブリを生成することは、オーバーラップ・レイアウト・コンセンサ(OLC)アルゴリズムを前記複数のヌクレオチド配列に適用することを含む、請求項39または任意の他の先行する請求項に記載のシステム。

【請求項42】

前記命令は、前記少なくとも1つのコンピュータハードウェアプロセッサに参照高分子のシーケンシングから取得された生物学的ポリマー配列と、前記参照高分子の所定のアセンブリとを含むトレーニングデータにアクセスするステップと、前記トレーニングデータを使用して深層学習モデルをトレーニングして、トレーニングされた深層学習モデルを取得するステップとをさらに実行させる、請求項23または任意の他の先行する請求項に記載のシステム。

20

【請求項43】

前記参照高分子は前記高分子とは異なる、請求項42または任意の他の先行する請求項に記載の方法。

【請求項44】

深層学習モデルが畳み込みニューラルネットワーク(CNN)を含む、請求項23または任意の他の先行する請求項に記載のシステム。

【請求項45】

命令を格納する少なくとも1つの非一時的なコンピュータ可読記憶媒体であって、前記命令は、少なくとも1つのコンピュータハードウェアプロセッサによる実行時に、前記少なくとも1つのコンピュータハードウェアプロセッサに高分子の生物学的ポリマーアセンブリを生成する方法を実行させ、前記方法は、

30

複数の生物学的ポリマー配列と、個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリとにアクセスするステップと、

前記複数の生物学的ポリマー配列および前記アセンブリを使用して、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップと、

前記第1の入力を前記トレーニングされた深層学習モデルに提供して、第1の複数のアセンブリ位置の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がその位置に存在する1つまたは複数の尤度を示す対応する第1の出力を取得するステップと、

前記トレーニングされた深層学習モデルの前記第1の出力を使用して、前記第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップと、

40

前記第1の複数のアセンブリ位置において同定された生物学的ポリマーを示すように前記アセンブリを更新して、更新されたアセンブリを取得するステップとを含む、少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項46】

前記高分子がタンパク質を含み、前記複数の生物学的ポリマー配列が複数のアミノ酸配列を含み、前記アセンブリが個々のアセンブリ位置におけるアミノ酸を示す、請求項45に記載の少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項47】

前記高分子が核酸を含み、前記複数の生物学的ポリマー配列が複数のヌクレオチド配列

50

を含み、前記アセンブリが個々のアセンブリ位置におけるヌクレオチドを示す、請求項 4 5 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【請求項 4 8】

前記アセンブリは、前記第 1 の複数のアセンブリ位置のうちの第 1 のアセンブリ位置における第 1 のヌクレオチドを示し、

前記第 1 の複数のアセンブリ位置における生物学的ポリマーを同定するステップは、前記第 1 のアセンブリ位置において第 2 のヌクレオチドを同定することを含み、

前記アセンブリを更新することは、前記第 1 のアセンブリ位置において前記第 2 のヌクレオチドを示すように前記アセンブリを更新することを含む、請求項 4 7 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

10

【請求項 4 9】

前記方法は、前記アセンブリを更新して、前記更新されたアセンブリを取得した後、前記複数のヌクレオチド配列を前記更新されたアセンブリに整列させるステップと、前記複数のヌクレオチド配列および前記更新されたアセンブリを使用して、前記トレーニングされた深層学習モデルに提供される第 2 の入力を生成するステップと、

前記第 2 の入力を前記トレーニングされた深層学習モデルに提供して、第 2 の複数のアセンブリ位置の各々に関して、1 つまたは複数の個々のヌクレオチドの各々がその位置に存在する 1 つまたは複数の尤度を示す対応する第 2 の出力を取得するステップと、

前記トレーニングされた深層学習モデルの前記第 2 の出力に基づいて、前記第 2 の複数のアセンブリ位置におけるヌクレオチドを同定するステップと、

前記第 2 の複数のアセンブリ位置において同定されたヌクレオチドを示すように前記更新されたアセンブリを更新して、第 2 の更新されたアセンブリを取得するステップとをさらに含む、請求項 4 7 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

20

【請求項 5 0】

前記方法が、前記複数のヌクレオチド配列を前記アセンブリに整列させるステップをさらに含む、請求項 4 7 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【請求項 5 1】

前記複数のヌクレオチド配列が少なくとも 9 個のヌクレオチド配列を含む、請求項 5 0 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

30

【請求項 5 2】

前記トレーニングされた深層学習モデルへの前記第 1 の入力を生成することは、前記第 1 の複数のアセンブリ位置を選択すること、選択された第 1 の複数のアセンブリ位置に基づいて前記第 1 の入力を生成することを含む、請求項 4 7 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【請求項 5 3】

前記アセンブリ内の前記第 1 の複数の位置を選択することは、前記アセンブリが前記第 1 の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定すること、

決定された尤度を使用して、前記第 1 の複数のアセンブリ位置を選択することを含む、請求項 5 2 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

40

【請求項 5 4】

前記トレーニングされた深層学習モデルに提供される前記第 1 の入力を生成するステップは、前記複数のヌクレオチド配列の個々の 1 つを前記アセンブリと比較することを含む、請求項 4 7 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコン

50

コンピュータ可読記憶媒体。

【請求項 55】

前記第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドを同定するためにトレーニングされた深層学習モデルに提供される前記第1の入力を生成するステップは、

前記第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置の各々における複数のヌクレオチドの各々に関して、

ヌクレオチドがその位置にあることを示す複数のヌクレオチド配列の数を示すカウントを決定すること、

前記アセンブリがその位置においてヌクレオチドを示しているかどうかに基づいて参照値を決定すること、

前記カウントと前記参照値との差異を示すエラー値を決定すること、

前記参照値および前記エラー値を前記第1の入力に含ませることを含む、請求項47または任意の他の先行する請求項に記載の少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項 56】

前記アセンブリがその位置においてヌクレオチドを示すかどうかに基づいて参照値を決定することは、

前記アセンブリがその位置においてヌクレオチドを示している場合、前記参照値が第1の値であると決定すること、

前記アセンブリがその位置においてヌクレオチドを示していない場合、前記参照値が第2の値であると決定することを含む、請求項55または任意の他の先行する請求項に記載の少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項 57】

前記第1の値は、前記複数のヌクレオチド配列の数であり、

前記第2の値は0である、請求項56または任意の他の先行する請求項に記載の少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項 58】

前記トレーニングされた深層学習モデルに提供される前記第1の入力を生成するステップは、複数の列を有するデータ構造に値を配置することを含み、

第1の列は、第1のアセンブリ位置において複数のヌクレオチドに対して決定された参照値およびエラー値を保持し、

第2の列は、前記第1のアセンブリ位置の近傍にある1つまたは複数のアセンブリ位置のうちの第2のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する、請求項55または任意の他の先行する請求項に記載の少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項 59】

前記第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置が、前記第1のアセンブリ位置とは別の少なくとも2つのアセンブリ位置を含む、請求項55または任意の他の先行する請求項に記載の少なくとも1つの非一時的なコンピュータ可読記憶媒体。

【請求項 60】

1つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する1つまたは複数の尤度が、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置において存在する尤度を含み、

前記第1の複数のアセンブリ位置における生物学的ポリマーを同定することは、第1のヌクレオチドが第1の位置に存在する尤度が複数のヌクレオチドのうちの第2のヌクレオチドが第1のアセンブリ位置に存在する尤度よりも大きいことを決定することによって前記第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドが複数のヌクレオチドのうちの第1のヌクレオチドであることを決定することを含む、請求項47または任意の他の先行する請求項に記載の少なくとも1つの非一時的なコンピュータ

10

20

30

40

50

可読記憶媒体。

【請求項 6 1】

前記方法が、前記複数のヌクレオチド配列から前記アセンブリを生成するステップをさらに含む、請求項 4 7 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【請求項 6 2】

前記複数のヌクレオチド配列から前記アセンブリを生成することは、前記アセンブリとなる複数のヌクレオチド配列からコンセンサス配列を決定することを含む、請求項 6 1 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

10

【請求項 6 3】

前記複数のヌクレオチド配列から前記アセンブリを生成することは、オーバーラップ・レイアウト・コンセンサス (OLC) アルゴリズムを前記複数のヌクレオチド配列に適用することを含む、請求項 6 1 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【請求項 6 4】

前記方法が

参照高分子のシーケンシングから取得された生物学的ポリマー配列と、前記参照高分子の所定のアセンブリとを含むトレーニングデータにアクセスするステップと、

前記トレーニングデータを使用して深層学習モデルをトレーニングして、トレーニングされた深層学習モデルを取得するステップとをさらに含む、請求項 4 5 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

20

【請求項 6 5】

前記参照高分子は前記高分子とは異なる、請求項 6 4 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【請求項 6 6】

深層学習モデルが畳み込みニューラルネットワーク (CNN) を含む、請求項 4 5 または任意の他の先行する請求項に記載の少なくとも 1 つの非一時的なコンピュータ可読記憶媒体。

【発明の詳細な説明】

30

【技術分野】

【0001】

本開示は、高分子 (例えば、核酸またはタンパク質) の生物学的ポリマー (例えば、ゲノムアセンブリ、ヌクレオチド配列、またはタンパク質配列) のアセンブリを生成することに関する。

【背景技術】

【0002】

シーケンシングデバイスは、アセンブリを生成するために使用することができるシーケンシングデータを生成し得る。一例として、シーケンシングデータは、ゲノムを (全体的または部分的に) 組み立てるために使用することができる生物学的サンプルからの DNA のヌクレオチド配列を含み得る。別の例として、シーケンシングデータは、タンパク質配列を (全体的または部分的に) 組み立てるために使用することができるアミノ酸配列を含み得る。

40

【発明の概要】

【0003】

一態様によれば、高分子の生物学的ポリマーアセンブリを生成する方法が提供される。方法は、少なくとも 1 つのコンピュータハードウェアプロセッサを使用して、複数の生物学的ポリマー配列と、個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリとにアクセスするステップと、複数の生物学的ポリマー配列およびアセンブリを使用して、トレーニングされた深層学習モデルに提供される第 1 の入力を生成するステップと、

50

第1の入力をトレーニングされた深層学習モデルに提供して、第1の複数のアセンブリ位置の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がその位置に存在する1つまたは複数の尤度を示す対応する第1の出力を取得するステップと、トレーニングされた深層学習モデルの第1の出力を使用して、第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップと、第1の複数のアセンブリ位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、更新されたアセンブリを取得するステップとを含む。

**【0004】**

一実施形態によれば、高分子はタンパク質を含み、複数の生物学的ポリマー配列は複数のアミノ酸配列を含み、アセンブリは個々のアセンブリ位置におけるアミノ酸を示す。

10

一実施形態によれば、高分子は核酸を含み、複数の生物学的ポリマー配列は複数のヌクレオチド配列を含み、アセンブリは個々のアセンブリ位置におけるヌクレオチドを示す。

**【0005】**

一実施形態によれば、アセンブリは、第1の複数のアセンブリ位置のうちの第1のアセンブリ位置における第1のヌクレオチドを示し、第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップは、第1のアセンブリ位置において第2のヌクレオチドを同定することを含み、アセンブリを更新するステップは、第1のアセンブリ位置における第2のヌクレオチドを示すようにアセンブリを更新することを含む。

**【0006】**

一実施形態によれば、方法は、アセンブリを更新して、更新されたアセンブリを取得した後、複数のヌクレオチド配列を更新されたアセンブリに整列させるステップと、複数のヌクレオチド配列および更新されたアセンブリを使用して、トレーニングされた深層学習モデルに提供される第2の入力を生成するステップと、第2の入力をトレーニングされた深層学習モデルに提供して、第2の複数のアセンブリ位置の各々に関して、1つまたは複数の個々のヌクレオチドの各々がその位置に存在する1つまたは複数の尤度を示す対応する第2の出力を取得するステップと、トレーニングされた深層学習モデルの第2の出力に基づいて、第2の複数のアセンブリ位置におけるヌクレオチドを同定するステップと、第2の複数のアセンブリ位置において同定されたヌクレオチドを示すように更新されたアセンブリを更新して、第2の更新されたアセンブリを取得するステップとを含む。

20

**【0007】**

一実施形態によれば、方法は、複数のヌクレオチド配列をアセンブリに整列させるステップをさらに含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも5個のヌクレオチド配列を含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも9個のヌクレオチド配列を含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも10個のヌクレオチド配列を含む。

30

**【0008】**

一実施形態によれば、トレーニングされた深層学習モデルへの第1の入力を生成するステップは、第1の複数のアセンブリ位置を選択すること、選択された第1の複数のアセンブリ位置に基づいて第1の入力を生成することを含む。一実施形態によれば、アセンブリ内の第1の複数のアセンブリ位置を選択することは、アセンブリが第1の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定すること、および決定された尤度を使用して、第1の複数のアセンブリ位置を選択することを含む。

40

**【0009】**

一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップは、複数のヌクレオチド配列の個々の1つをアセンブリと比較することを含む。一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成して、第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドを同定することは、第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置における複数のヌクレオチドの各々に関して、ヌクレオチドがその位置にあることを示す複数のヌクレオチド配列の数を示すカウントを決定すること、アセンブリがその位置に

50

においてヌクレオチドを示しているかどうかに基づいて参照値を決定すること、カウントと参照値との間の差異を示すエラー値を決定すること、第1の入力に参照値およびエラー値を含ませることを含む。

【0010】

一実施形態によれば、アセンブリがその位置においてヌクレオチドを示すかどうかに基づいて参照値を決定することは、アセンブリがその位置においてヌクレオチドを示している場合、参照値が第1の値であると決定すること、アセンブリがその位置においてヌクレオチドを示していない場合、参照値が第2の値であると決定することを含む。一実施形態によれば、第1の値は、複数のヌクレオチド配列の数であり、第2の値は0である。

【0011】

一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップは、複数の列を有するデータ構造に値を配置することを含み、第1の列は、第1のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持し、第2の列は、第1のアセンブリ位置の近傍にある1つまたは複数のアセンブリ位置のうちの第2のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する。一実施形態によれば、第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置は、第1のアセンブリ位置とは別の少なくとも2つのアセンブリ位置を含む。

【0012】

一実施形態によれば、1つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する1つまたは複数の尤度は、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置に存在する尤度を含み、第1の複数のアセンブリ位置における生物学的ポリマーを同定することは、第1のヌクレオチドが第1の位置に存在する尤度が複数のヌクレオチドのうちの第2のヌクレオチドが第1のアセンブリ位置に存在する尤度よりも大きいことを決定することによって第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドが複数のヌクレオチドのうちの第1のヌクレオチドであることを同定することを含む。

【0013】

一実施形態によれば、方法は、複数のヌクレオチド配列からアセンブリを生成するステップをさらに含む。一実施形態によれば、複数のヌクレオチド配列からアセンブリを生成するステップは、アセンブリとなる複数のヌクレオチド配列からコンセンサス配列を決定することを含む。一実施形態によれば、複数のヌクレオチド配列からアセンブリを生成するステップは、オーバーラップ・レイアウト・コンセンサス(OLC)アルゴリズムを複数のヌクレオチド配列に適用することを含む。

【0014】

一実施形態によれば、方法は、参照高分子のシーケンシングから取得された生物学的ポリマー配列と、参照高分子の所定のアセンブリとを含むトレーニングデータにアクセスするステップと、トレーニングデータを使用して深層学習モデルをトレーニングして、トレーニングされた深層学習モデルを取得するステップとをさらに含む。一実施形態によれば、参照高分子は、高分子とは異なる。一実施形態によれば、深層学習モデルは、畳み込みニューラルネットワーク(CNN)を含む。

【0015】

別の態様によれば、高分子の生物学的ポリマーアセンブリを生成するためのシステムが提供される。システムは、少なくとも1つのコンピュータハードウェアプロセッサと、命令を格納する少なくとも1つの非一時的なコンピュータ可読記憶媒体とを備え、命令は、少なくとも1つのコンピュータハードウェアプロセッサによる実行時に、少なくとも1つのコンピュータハードウェアプロセッサに、複数の生物学的ポリマー配列と、個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリとにアクセスするステップと、複数の生物学的ポリマー配列およびアセンブリを使用して、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップと、第1の入力をトレーニングされた

10

20

30

40

50

深層学習モデルに提供して、第1の複数のアセンブリ位置の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がその位置に存在する1つまたは複数の尤度を示す対応する第1の出力を取得するステップと、トレーニングされた深層学習モデルの第1の出力を使用して、第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップと、第1の複数のアセンブリ位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、更新されたアセンブリを取得するステップとを実行させる。

【0016】

一実施形態によれば、高分子はタンパク質を含み、複数の生物学的ポリマー配列は複数のアミノ酸配列を含み、アセンブリは個々のアセンブリ位置におけるアミノ酸を示す。

一実施形態によれば、高分子は核酸を含み、複数の生物学的ポリマー配列は複数のヌクレオチド配列を含み、アセンブリは個々のアセンブリ位置におけるヌクレオチドを示す。

【0017】

一実施形態によれば、アセンブリは、第1の複数のアセンブリ位置のうちの第1のアセンブリ位置における第1のヌクレオチドを示し、第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップは、第1のアセンブリ位置において第2のヌクレオチドを同定することを含み、アセンブリを更新するステップは、第1のアセンブリ位置における第2のヌクレオチドを示すようにアセンブリを更新することを含む。

【0018】

一実施形態によれば、命令はさらに、少なくとも1つのコンピュータハードウェアプロセッサに、アセンブリを更新して更新されたアセンブリを取得した後、複数のヌクレオチド配列を更新されたアセンブリに整列させるステップと、複数のヌクレオチド配列および更新されたアセンブリを使用して、トレーニングされた深層学習モデルに提供される第2の入力を生成するステップと、第2の入力をトレーニングされた深層学習モデルに提供して、第2の複数のアセンブリ位置の各々に関して、1つまたは複数の個々のヌクレオチドの各々がその位置に存在する1つまたは複数の尤度を示す対応する第2の出力を取得するステップと、トレーニングされた深層学習モデルの第2の出力に基づいて、第2の複数のアセンブリ位置におけるヌクレオチドを同定するステップと、第2の複数のアセンブリ位置において同定されたヌクレオチドを示すように更新されたアセンブリを更新して、第2の更新されたアセンブリを取得するステップとを実行させる。

【0019】

一実施形態によれば、命令はさらに、少なくとも1つのコンピュータハードウェアプロセッサに、複数のヌクレオチド配列をアセンブリに整列させるステップを実行させる。一実施形態によれば、複数のヌクレオチド配列は、少なくとも5個のヌクレオチド配列を含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも9個のヌクレオチド配列を含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも10個のヌクレオチド配列を含む。

【0020】

一実施形態によれば、トレーニングされた深層学習モデルへの第1の入力を生成するステップは、第1の複数のアセンブリ位置を選択すること、選択された第1の複数のアセンブリ位置に基づいて第1の入力を生成することを含む。一実施形態によれば、アセンブリ内の第1の複数の位置を選択することは、アセンブリが第1の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定すること、および決定された尤度を使用して、第1の複数のアセンブリ位置を選択することを含む。

【0021】

一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップは、複数のヌクレオチド配列の個々の1つをアセンブリと比較することを含む。一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成して、第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドを同定することは、第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置における複数のヌクレオチドの各々に関して、ヌクレオチドがその位置にあることを示

10

20

30

40

50

す複数のヌクレオチド配列の数を示すカウントを決定すること、アセンブリがその位置においてヌクレオチドを示しているかどうかに基づいて参照値を決定すること、カウントと参照値との間の差異を示すエラー値を決定すること、第1の入力に参照値およびエラー値を含ませることを含む。一実施形態によれば、アセンブリがその位置においてヌクレオチドを示すかどうかに基づいて参照値を決定することは、アセンブリがその位置においてヌクレオチドを示している場合、参照値が第1の値であると決定すること、アセンブリがその位置においてヌクレオチドを示していない場合、参照値が第2の値であると決定することを含む。一実施形態によれば、第1の値は、複数のヌクレオチド配列の数であり、第2の値は、0である。一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップは、複数の列を有するデータ構造に値を配置することを含み、第1の列は、第1のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持し、第2の列は、第1のアセンブリ位置の近傍にある1つまたは複数のアセンブリ位置のうちの第2のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する。一実施形態によれば、第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置は、第1のアセンブリ位置とは別の少なくとも2つのアセンブリ位置を含む。

10

20

30

40

50

**【0022】**

一実施形態によれば、1つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する1つまたは複数の尤度は、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置に存在する尤度を含み、第1の複数のアセンブリ位置における生物学的ポリマーを同定することは、第1のヌクレオチドが第1の位置に存在する尤度が複数のヌクレオチドのうちの第2のヌクレオチドが第1のアセンブリ位置に存在する尤度よりも大きいことを決定することによって第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドが複数のヌクレオチドのうちの第1のヌクレオチドであることを同定することを含む。

**【0023】**

一実施形態によれば、命令はさらに、少なくとも1つのコンピュータハードウェアプロセッサに、複数のヌクレオチド配列からアセンブリを生成するステップを実行させる。一実施形態によれば、複数のヌクレオチド配列からアセンブリを生成するステップは、アセンブリとなる複数のヌクレオチド配列からコンセンサス配列を決定することを含む。一実施形態によれば、複数のヌクレオチド配列からアセンブリを生成するステップは、オーバーラップ・レイアウト・コンセンサス(OLC)アルゴリズムを複数のヌクレオチド配列に適用することを含む。

**【0024】**

一実施形態によれば、命令はさらに、少なくとも1つのコンピュータハードウェアプロセッサに、参照高分子および参照高分子の所定のアセンブリのシーケンシングから取得された生物学的ポリマー配列を含むトレーニングデータにアクセスするステップと、トレーニングデータを使用して深層学習モデルをトレーニングし、トレーニングされた深層学習モデルを取得するステップとを実行させる。一実施形態によれば、参照高分子は高分子とは異なる。一実施形態によれば、深層学習モデルは、畳み込みニューラルネットワーク(CNN)を含む。

**【0025】**

別の態様によれば、非一時的なコンピュータ可読記憶媒体が提供される。非一時的なコンピュータ可読記憶媒体は、少なくとも1つのコンピュータハードウェアプロセッサによる実行時に、少なくとも1つのコンピュータハードウェアプロセッサに高分子の生物学的ポリマーアセンブリを生成する方法を実行させる命令を格納する。方法は、複数の生物学的ポリマー配列と、個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリとにアクセスするステップと、複数の生物学的ポリマー配列およびアセンブリを使用して、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップと、第1の入力をトレーニングされた深層学習モデルに提供して、第1の複数のアセンブリ位置

の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がその位置に存在する1つまたは複数の尤度を示す対応する第1の出力を取得するステップと、トレーニングされた深層学習モデルの第1の出力を使用して、第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップと、第1の複数のアセンブリ位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、更新されたアセンブリを取得するステップとを含む。

【0026】

一実施形態によれば、高分子はタンパク質を含み、複数の生物学的ポリマー配列は複数のアミノ酸配列を含み、アセンブリは個々のアセンブリ位置におけるアミノ酸を示す。

一実施形態によれば、高分子は核酸を含み、複数の生物学的ポリマー配列は複数のヌクレオチド配列を含み、アセンブリは個々のアセンブリ位置におけるヌクレオチドを示す。

【0027】

一実施形態によれば、アセンブリは、第1の複数のアセンブリ位置のうちの第1のアセンブリ位置における第1のヌクレオチドを示し、第1の複数のアセンブリ位置における生物学的ポリマーを同定するステップは、第1のアセンブリ位置において第2のヌクレオチドを同定することを含み、アセンブリを更新するステップは、第1のアセンブリ位置における第2のヌクレオチドを示すようにアセンブリを更新することを含む。

【0028】

一実施形態によれば、方法は、アセンブリを更新して、更新されたアセンブリを取得した後、複数のヌクレオチド配列を更新されたアセンブリに整列させるステップと、複数のヌクレオチド配列および更新されたアセンブリを使用して、トレーニングされた深層学習モデルに提供される第2の入力を生成するステップと、第2の入力をトレーニングされた深層学習モデルに提供して、第2の複数のアセンブリ位置の各々に関して、1つまたは複数の個々のヌクレオチドの各々がその位置に存在する1つまたは複数の尤度を示す対応する第2の出力を取得するステップと、トレーニングされた深層学習モデルの第2の出力に基づいて、第2の複数のアセンブリ位置におけるヌクレオチドを同定するステップと、第2の複数のアセンブリ位置において同定されたヌクレオチドを示すように更新されたアセンブリを更新して、第2の更新されたアセンブリを取得するステップとを含む。

【0029】

一実施形態によれば、方法は、複数のヌクレオチド配列をアセンブリに整列させるステップをさらに含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも5個のヌクレオチド配列を含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも9個のヌクレオチド配列を含む。一実施形態によれば、複数のヌクレオチド配列は、少なくとも10個のヌクレオチド配列を含む。

【0030】

一実施形態によれば、トレーニングされた深層学習モデルへの第1の入力を生成するステップは、第1の複数のアセンブリ位置を選択すること、選択された第1の複数のアセンブリ位置に基づいて第1の入力を生成することを含む。一実施形態によれば、アセンブリ内の第1の複数の位置を選択することは、アセンブリが第1の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定すること、および決定された尤度を使用して、第1の複数のアセンブリ位置を選択することを含む。

【0031】

一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップは、複数のヌクレオチド配列の個々の1つをアセンブリと比較することを含む。一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成して、第1の複数のアセンブリ位置のうちの第1のアセンブリ位置におけるヌクレオチドを同定することは、第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置における複数のヌクレオチドの各々に関して、ヌクレオチドがその位置にあることを示す複数のヌクレオチド配列の数を示すカウントを決定すること、アセンブリがその位置においてヌクレオチドを示しているかどうかに基づいて参照値を決定すること、カウントと

10

20

30

40

50

参照値との間の差異を示すエラー値を決定すること、第1の入力に参照値およびエラー値を含ませることを含む。一実施形態によれば、アセンブリがその位置においてヌクレオチドを示すかどうかに基づいて参照値を決定することは、アセンブリがその位置においてヌクレオチドを示している場合、参照値が第1の値であると決定すること、アセンブリがその位置においてヌクレオチドを示していない場合、参照値が第2の値であると決定することを含む。一実施形態によれば、第1の値は、複数のヌクレオチド配列の数であり、第2の値は、0である。一実施形態によれば、トレーニングされた深層学習モデルに提供される第1の入力を生成するステップは、複数の列を有するデータ構造に値を配置することを含み、第1の列は、第1のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持し、第2の列は、第1のアセンブリ位置の近傍にある1つまたは複数のアセンブリ位置のうち第2のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する。一実施形態によれば、第1のアセンブリ位置の近傍の1つまたは複数のアセンブリ位置は、第1のアセンブリ位置とは別の少なくとも2つのアセンブリ位置を含む。

#### 【0032】

一実施形態によれば、1つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する1つまたは複数の尤度は、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置に存在する尤度を含み、第1の複数のアセンブリ位置における生物学的ポリマーを同定することは、第1のヌクレオチドが第1の位置に存在する尤度が複数のヌクレオチドのうち第2のヌクレオチドが第1のアセンブリ位置に存在する尤度よりも大きいことを決定することによって第1の複数のアセンブリ位置のうち第1のアセンブリ位置におけるヌクレオチドが複数のヌクレオチドのうち第1のヌクレオチドであることを同定することを含む。

#### 【0033】

一実施形態によれば、方法は、複数のヌクレオチド配列からアセンブリを生成するステップをさらに含む。一実施形態によれば、複数のヌクレオチド配列からアセンブリを生成するステップは、アセンブリとなる複数のヌクレオチド配列からコンセンサス配列を決定することを含む。一実施形態によれば、複数のヌクレオチド配列からアセンブリを生成するステップは、オーバーラップ・レイアウト・コンセンサス(OLC)アルゴリズムを複数のヌクレオチド配列に適用することを含む。

#### 【0034】

一実施形態によれば、方法は、参照高分子のシーケンシングから取得された生物学的ポリマー配列と、参照高分子の所定のアセンブリとを含むトレーニングデータにアクセスするステップと、トレーニングデータを使用して深層学習モデルをトレーニングして、トレーニングされた深層学習モデルを取得するステップとをさらに含む。一実施形態によれば、参照高分子は、高分子とは異なる。一実施形態によれば、深層学習モデルは、畳み込みニューラルネットワーク(CNN)を含む。

#### 【図面の簡単な説明】

#### 【0035】

以下の図面を参照して、本出願の様々な態様および実施形態に関して説明する。図面は必ずしも一定の縮尺で描かれているわけではないことを理解されたい。複数の図面に表示されている構成要素は、表示されている全ての図面で同じ参照番号で示されている。

【図1A】本明細書に記載の技術のいくつかの実施形態による、本明細書に記載の技術の態様を実施し得るシステムを示す図である。

【図1B】本明細書に記載の技術のいくつかの実施形態による、本明細書に記載の技術の態様を実施し得るシステムを示す図である。

【図1C】本明細書に記載の技術のいくつかの実施形態による、本明細書に記載の技術の態様を実施し得るシステムを示す図である。

【図2A】本明細書に記載の技術のいくつかの実施形態による、アセンブリシステムの実施形態を示す図である。

10

20

30

40

50

【図 2 B】本明細書に記載の技術のいくつかの実施形態による、アセンブリシステムの実施形態を示す図である。

【図 2 C】本明細書に記載の技術のいくつかの実施形態による、アセンブリシステムの実施形態を示す図である。

【図 2 D】本明細書に記載の技術のいくつかの実施形態による、アセンブリシステムの実施形態を示す図である。

【図 3 A】本明細書に記載の技術のいくつかの実施形態による、生物学的ポリマーアセンブリを生成するための機械学習モデルをトレーニングするための例示的なプロセス 300 を示す図である。

【図 3 B】本明細書に記載の技術のいくつかの実施形態による、図 3 A のプロセスによって取得された機械学習モデルを使用して生物学的ポリマーアセンブリを生成するための例示的なプロセス 310 を示す図である。

【図 4 A】本明細書に記載の技術のいくつかの実施形態による、機械学習モデルへの入力を生成する例を示す図である。

【図 4 B】本明細書に記載の技術のいくつかの実施形態による、機械学習モデルへの入力を生成する例を示す図である。

【図 4 C】本明細書に記載の技術のいくつかの実施形態による、機械学習モデルへの入力を生成する例を示す図である。

【図 5】本明細書に記載の技術のいくつかの実施形態による、生物学的ポリマーアセンブリを更新する例を示す図である。

【図 6】本明細書に記載の技術のいくつかの実施形態による、生物学的ポリマーアセンブリを生成するために使用される例示的な畳み込みニューラルネットワーク (CNN) モデルの構造を示す図である。

【図 7】従来技術と比較した、本明細書に記載の技術のいくつかの実施形態により具体化されたアセンブリ技術の性能を示す図である。

【図 8】本明細書に記載の技術のいくつかの実施形態を実施する際に使用し得る例示的なコンピューティングデバイス 800 のブロック図である。

【発明を実施するための形態】

【0036】

高分子は、タンパク質またはタンパク質フラグメント、(任意のタイプの DNA の) DNA 分子またはフラグメント、または(任意のタイプの RNA の) RNA 分子またはフラグメントであり得る。生物学的ポリマーは、アミノ酸(例えば、高分子がタンパク質またはそのフラグメントである場合)、またはヌクレオチド(例えば、高分子が DNA、RNA、またはそのフラグメントである場合)であり得る。

【0037】

本発明者らは、機械学習技術を使用して高分子の生物学的ポリマーアセンブリを生成するシステムを開発した。例えば、本発明者らによって開発されたシステムは、機械学習技術を使用して、生物の DNA のゲノムアセンブリを生成するように構成され得る。別の例として、本発明者らによって開発されたシステムは、機械学習技術を使用してタンパク質のアミノ酸配列を生成するように構成され得る。

【0038】

いくつかの実施形態では、システムは、1つまたは複数の生物学的ポリマー配列(例えば、シーケンシングデバイスによって生成される)および配列から生成された初期アセンブリにアクセスし得る。アセンブリは、個々のアセンブリの位置において生物学的ポリマー(例えば、ヌクレオチド、アミノ酸)が存在することを示し得る。システムは、(1)配列と初期アセンブリとを使用して、機械学習モデルに提供される入力を生成し、(2)入力をトレーニング済みの機械学習モデルに提供して、対応する出力を取得し、(3)機械学習モデルから取得した出力を使用して初期アセンブリを更新し、更新されたアセンブリを取得することによって、初期アセンブリの生物学的ポリマーの表示のエラーを修正し得る。更新されたアセンブリは、初期アセンブリよりも生物学的ポリマーの表示にお

10

20

30

40

50

るエラーが少なくなり得る。

【0039】

いくつかの実施形態では、アセンブリは、複数の位置と、個々の位置における生物学的ポリマー（例えば、ヌクレオチドまたはアミノ酸）の表示とを含み得る。例として、アセンブリは、生物のゲノム内の位置におけるヌクレオチドを示すゲノムアセンブリであり得る。別の例として、アセンブリは、生物のDNAの一部のヌクレオチドの配列を示す遺伝子配列であり得る。別の例として、アセンブリは、タンパク質のアミノ酸配列（「タンパク質配列」とも呼ばれる）であり得る。生物学的ポリマーは、ヌクレオチド、アミノ酸、または他の任意のタイプの生物学的ポリマーであり得る。生物学的ポリマー配列は、本明細書では「配列」または「リード（read）」と呼ばれ得る。

10

【0040】

いくつかの従来 of 生物学的ポリマーアセンブリ技術は、シーケンシング技術を利用して高分子（例えば、DNA、RNA、またはタンパク質）の生物学的ポリマー配列を生成し、生成された配列を使用して高分子のアセンブリを生成し得る。例えば、シーケンシングデバイスは、生物のDNAサンプルからヌクレオチド配列を生成し得、その配列を使用して、生物のDNAのゲノムアセンブリを生成し得る。別の例として、シーケンシングデバイスは、タンパク質サンプルのアミノ酸配列を生成し得、その配列を使用して、タンパク質のより長いアミノ酸配列を組み立て得る。コンピューティングデバイスは、シーケンシングデバイスによって生成された配列にアセンブリアルゴリズムを適用してアセンブリを生成し得る。例えば、コンピューティングデバイスは、DNAサンプルのヌクレオチド配列にオーバーラップ・レイアウト・コンセンサス（OLC）アセンブリアルゴリズムを適用して、生物のゲノムアセンブリまたはその一部を生成し得る。

20

【0041】

核酸サンプルからヌクレオチド配列を生成するために使用されるシーケンシング技術の1つのタイプは、1000個未満のヌクレオチドのヌクレオチド配列（即ち、「ショートリード」）を生成する第2世代シーケンシング（「ショートリードシーケンシング」としても知られる）である。シーケンシング技術は、1000個以上のヌクレオチドのヌクレオチド配列（即ち、「ロングリード」）を生成し、かつ第2世代シーケンシングよりもアセンブリの大きな部分を提供する第3世代シーケンシング（「ロングリードシーケンシング」とも呼ばれる）に進化した。しかしながら、本発明者らは、第3世代シーケンシングは第2世代シーケンシングよりも精度が低く、その結果、ロングリードから生成されたアセンブリはショートリードから生成されたアセンブリよりも精度が低いことを認識した。本発明者らはまた、アセンブリの精度を向上するための従来 of エラー訂正技術は、計算コストおよび時間がかかることを認識した。従って、本発明者らは、（1）第3世代シーケンシングから生成されたアセンブリの精度を向上させ、（2）従来 of エラー訂正技術よりも効率的であるアセンブリのエラーを修正するための機械学習技術を開発した。

30

【0042】

本明細書に記載のいくつかの実施形態は、発明者がアセンブリの生成に関して認識した上記の問題の全てに対処する。しかしながら、本明細書に記載される全ての実施形態がこれらの問題の全てに対処するわけではないことを理解されたい。本明細書に記載の技術の実施形態は、生物学的ポリマーアセンブリの上記の問題に対処する以外の目的に使用し得ることも理解されたい。一例として、本明細書に記載の技術の実施形態を使用して、アミノ酸配列から生成されたタンパク質配列の精度を向上し得る。別の例として、本明細書に記載の技術の実施形態を使用して、ショートリードから生成されたアセンブリの精度を向上し得る。

40

【0043】

いくつかの実施形態では、システムは、（1）個々のアセンブリ位置に存在する生物学的ポリマーを示すアセンブリ（例えば、複数の生物学的ポリマー配列から生成される）にアクセスし、（2）複数の生物学的ポリマー配列およびアセンブリを使用して、トレーニ

50

ングされた深層学習モデルに提供される第1の入力を生成し、(3)第1の入力をトレーニングされた深層学習モデルに提供して、第1の複数のアセンブリ位置の各々に関して、1つまたは複数の個々の生物学的ポリマーの各々がそのアセンブリ位置に存在する1つまたは複数の尤度(例えば、確率)を示す対応する第1の出力を取得し、(4)トレーニングされた深層学習モデルの第1の出力を使用して、第1の複数のアセンブリ位置における生物学的ポリマーを同定し、(5)第1の複数のアセンブリ位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、更新されたアセンブリを取得するように構成されている。いくつかの実施形態では、システムは、複数の生物学的ポリマー配列をアセンブリに整列させるように構成され得る。

**【0044】**

いくつかの実施形態では、高分子はタンパク質であり得、複数の生物学的ポリマー配列は複数のアミノ酸配列であり得、アセンブリは個々のアセンブリ位置におけるアミノ酸を示す。いくつかの実施形態において、高分子は、核酸(例えば、DNA、RNA)であり得、複数の生物学的配列は、複数のヌクレオチド配列であり得、アセンブリは、個々のアセンブリ位置におけるヌクレオチドを示す。

**【0045】**

いくつかの実施形態では、アセンブリは、複数のアセンブリ位置のうちの第1のアセンブリ位置における第1のヌクレオチド(例えば、アデニン)を示す。第1の複数のアセンブリ位置における生物学的ポリマーを同定することは、第1のアセンブリ位置において第1のヌクレオチドとは異なる第2のヌクレオチド(例えば、チミン)を同定することを含み、アセンブリを更新することは、第1のアセンブリ位置における第2のヌクレオチド(例えば、チミン)を示すようにアセンブリを更新することを含む。

**【0046】**

いくつかの実施形態では、システムは、複数の更新の反復を実行するように構成され得る。システムは、アセンブリを更新して、更新されたアセンブリを取得した後、(1)複数のヌクレオチド配列を更新されたアセンブリに整列させ、(2)複数のヌクレオチド配列および更新されたアセンブリを使用して、トレーニングされた深層学習モデルに提供される第2の入力を生成し、(3)第2の入力をトレーニングされた深層学習モデルに提供して、第2の複数のアセンブリ位置の各々に関して、1つまたは複数の個々のヌクレオチドの各々がそのアセンブリ位置に存在する1つまたは複数の尤度(例えば、確率)を示す対応する第2の出力を取得し、(4)トレーニングされた深層学習モデルの第2の出力に基づいて、第2の複数のアセンブリ位置におけるヌクレオチドを同定し、(5)第2の複数のアセンブリ位置において同定されたヌクレオチドを示すように更新されたアセンブリを更新して、第2の更新されたアセンブリを取得するように構成され得る。

**【0047】**

いくつかの実施形態では、システムは、(1)第1の複数のアセンブリ位置を選択し、(2)選択された第1の複数のアセンブリ位置に基づいて第1の入力を生成することによって、トレーニングされた深層学習モデルへの第1の入力を生成するように構成され得る。いくつかの実施形態では、システムは、(1)アセンブリが第1の複数のアセンブリ位置においてヌクレオチドを不正確に示す尤度を決定し、(2)決定された尤度を使用して、第1の複数のアセンブリ位置を選択することによって、第1の複数のアセンブリ位置を選択するように構成され得る。

**【0048】**

いくつかの実施形態では、システムは、(例えば、1つまたは複数の特徴の値を決定するために)複数のヌクレオチド配列の個々の1つをアセンブリと比較することによって、トレーニングされた深層学習モデルに提供される第1の入力を生成するように構成され得る。いくつかの実施形態では、システムは、第1の入力の近傍にある1つまたは複数のアセンブリ位置の各々における複数のヌクレオチドの各々に関して、(1)ヌクレオチドがそのアセンブリ位置にあることを示す複数のヌクレオチド配列の数を示すカウントを決定し、(2)アセンブリがそのアセンブリ位置においてヌクレオチドを示しているかどうか

10

20

30

40

50

に基づいて参照値を決定し、(3) カウントと基準値との間の差異を示すエラー値を決定し、(4) 第1の入力に基準値およびエラー値を含ませることによって、第1の複数のアセンブリ位置の第1のアセンブリ位置におけるヌクレオチドを同定するための第1の入力を生成するように構成され得る。いくつかの実施形態では、システムは、アセンブリがそのアセンブリ位置においてヌクレオチドを示すかどうかに基づいて、(1) アセンブリがそのアセンブリ位置においてヌクレオチドを示している場合、参照値が第1の値(例えば、複数のヌクレオチド配列の数)であると決定し、(2) アセンブリがそのアセンブリ位置においてヌクレオチドを示していない場合、参照値が第2の値(例えば、0)であると決定することにより、参照値を決定するように構成され得る。いくつかの実施形態では、システムは、3個、4個、5個、6個、7個、8個、9個、10個、15個、20個、25個、30個、35個、40個、45個、または50個の位置の近傍を使用するように構成され得る。

10

**【0049】**

いくつかの実施形態では、システムは、行/列を有するデータ構造に値を配置することによって、第1のアセンブリ位置におけるヌクレオチドを同定するための第1の入力を生成するように構成され得、(1) 第1の行/列は、第1のアセンブリ位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持し、(2) 第2の行/列は、第1のアセンブリ位置の近傍にある第2の位置において複数のヌクレオチドに関して決定された参照値およびエラー値を保持する。

20

**【0050】**

いくつかの実施形態では、1つまたは複数の個々の生物学的ポリマーの各々がアセンブリ位置に存在する1つまたは複数の尤度は、複数のヌクレオチドの各々に関して、ヌクレオチドがアセンブリ位置において存在する尤度(例えば、確率)を含む。システムは、第1の複数のアセンブリ位置のうち第1のアセンブリ位置におけるヌクレオチドが複数のヌクレオチドのうち第1のヌクレオチドであることを同定することによって、アセンブリ内の第1の複数のアセンブリ位置における生物学的ポリマーを同定するように構成され得る。システムは、第1のヌクレオチドが第1のアセンブリ位置に存在する尤度が、複数のヌクレオチドのうち第2のヌクレオチドが第1のアセンブリ位置に存在する尤度よりも大きいことを決定することによって、第1のアセンブリ位置におけるヌクレオチドが第1のヌクレオチドであることを同定し得る。

30

**【0051】**

いくつかの実施形態では、システムは、複数のヌクレオチド配列からアセンブリ(例えば、初期アセンブリ)を生成するように構成され得る。いくつかの実施形態では、システムは、アセンブリとなる複数のヌクレオチド配列からコンセンサス配列を決定することによって(例えば、多数決を取ることによって)アセンブリを生成するように構成され得る。いくつかの実施形態では、システムは、オーバーラップ・レイアウト・コンセンサス(OLC)アルゴリズムを複数のヌクレオチド配列に適用することによって、複数のヌクレオチド配列からアセンブリを生成するように構成され得る。いくつかの実施形態では、システムは、(1) 参照高分子のシーケンシングから取得された生物学的ポリマー配列と、参照高分子の所定の生物学的ポリマーアセンブリとを含むトレーニングデータにアクセスし、(2) トレーニングデータを使用して深層学習モデル(畳み込みニューラルネットワークまたは再帰型ニューラルネットワークなど)をトレーニングして、トレーニングされた深層学習モデルを取得するように構成されている。いくつかの実施形態では、深層学習モデルをトレーニングするために使用される参照高分子は、アセンブリが生成されている高分子とは異なり得る。

40

**【0052】**

上記で導入され、以下でより詳細に説明される技術は、技術が特定の実施形態に限定されないことから、多数の方法のいずれかで実施され得ることを理解されたい。実施形態の詳細の例は、説明のみを目的として本明細書に提供されている。さらに、本明細書に記載の技術の態様は、特定の技術または技術の組み合わせの使用に限定されないことから、本

50

明細書に開示される技術は、個別にまたは任意の適切な組み合わせで使用され得る。

【0053】

図1Aは、本明細書に記載の技術の態様を具体化し得るシステム100を示す。システム100は、1つまたは複数のシーケンシングデバイス102、アセンブリシステム104、モデルトレーニングシステム106、およびデータストア108Aを含み、これらの各々は、ネットワーク111に接続されている。

【0054】

いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、高分子の1つまたは複数のサンプル標本110のシーケンシングによってシーケンシングデータを生成するように構成され得る。例えば、サンプル標本110は、核酸（例えば、DNAおよび/またはRNA）、またはタンパク質（例えば、ペプチド）を含む生物学的サンプルであり得る。シーケンシングデータは、サンプル標本（単数または複数）110の生物学的ポリマー配列を含み得る。生物学的ポリマー配列は、高分子サンプル中に存在する生物学的ポリマーの順序および位置を示す英数字記号の配列として表され得る。いくつかの実施形態では、生物学的ポリマー配列は、生物学的サンプルのシーケンシングから生成されたヌクレオチド配列であり得る。例として、ヌクレオチド配列は、（1）アデニンを表す「A」、（2）シトシンを表す「C」、（3）グアニンを表す「G」、（4）チミンを表す「T」、（5）ウラシルを表す「U」、（6）配列内の位置にヌクレオチドが存在しないことを表す「-」を使用し得る。いくつかの実施形態では、生物学的ポリマー配列は、タンパク質サンプル（例えば、ペプチド）のシーケンシングから生成されたアミノ酸配列であり得る。一例として、アミノ酸配列は、タンパク質に存在し得る個々の異なるアミノ酸を表すために異なる英数字を使用する英数字配列であり得る。

10

20

【0055】

いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、核酸サンプル（例えば、DNAサンプル）のシーケンシングからヌクレオチド配列を生成するように構成され得る。いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、合成によって核酸サンプルをシーケンシングするように構成され得る。シーケンシングデバイス（単数または複数）102は、ヌクレオチドが、シーケンシングされている核酸に相補的である核酸の新たに合成された鎖に取り込まれるときに、ヌクレオチドを同定するように構成され得る。シーケンシング中に、重合酵素（例えば、DNAポリメラーゼ）は、ターゲット核酸分子のプライミング位置（「プライマー」と呼ばれる）に結合（例えば、付着）して、重合酵素の作用を介してヌクレオチドをプライマーに取り込み得る。シーケンシングデバイス（単数または複数）102は、取り込まれている各ヌクレオチドを検出するように構成され得る。いくつかの実施形態において、ヌクレオチドは、励起にตอบสนองして発光する個々の発光分子（例えば、フルオロフォア）と結合され得る。発光分子は、個々のヌクレオチドと結合している発光分子が取り込まれているときに励起され得る。シーケンシングデバイス（単数または複数）102は、発光を検出するための1つまたは複数のセンサを含み得る。各タイプのヌクレオチドは、個々のタイプの発光分子と結合され得る。シーケンシングデバイス（単数または複数）102は、検出された発光に基づいて発光分子のタイプを同定することによって、取り込まれているヌクレオチドを同定し得る。例えば、シーケンシングデバイス（単数または複数）102は、発光強度、寿命、波長、または他の特性を使用して、異なる発光分子を区別し得る。いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、ヌクレオチドの取り込み中に生成された電気信号を検出して、取り込まれているヌクレオチドを同定するように構成され得る。シーケンシングデバイス（単数または複数）102は、電気信号を検出し、それらを使用して取り込まれているヌクレオチドを同定するためのセンサ（単数または複数）を含み得る。

30

40

【0056】

いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、本明細書に記載されているものとは異なる技術を使用して核酸をシーケンシングするよう

50

に構成され得る。いくつかの実施形態は、本明細書に記載の核酸シーケンシングの特定の技術に限定されない。

【0057】

いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、タンパク質サンプル（例えば、ペプチド）のシーケンシングからアミノ酸配列を生成するように構成され得る。いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、個々のアミノ酸に選択的に結合する試薬を使用してタンパク質サンプルをシーケンシングするように構成され得る。試薬は、他のタイプのアミノ酸よりも1つまたは複数のタイプのアミノ酸に選択的に結合し得る。いくつかの実施形態において、試薬は、個々の発光分子と結合され得る。発光分子は、発光分子と結合されている試薬とアミノ酸との間の相互作用に反応して励起され得る。いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、発光分子の発光を検出することによってアミノ酸を同定するように構成され得る。シーケンシングデバイス102は、発光を検出するための1つまたは複数のセンサを含み得る。いくつかの実施形態において、各タイプのアミノ酸は、個々のタイプの発光分子と結合され得る。シーケンシングデバイス（単数または複数）102は、検出された発光に基づいて発光分子のタイプを同定することによってアミノ酸を同定し得る。一例として、シーケンシングデバイス（単数または複数）102は、発光強度、寿命、波長、または他の特性を使用して、異なる発光分子を区別し得る。いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、試薬とアミノ酸との間の結合相互作用の間に生成される電気信号を検出するように構成され得る。シーケンシングデバイス（単数または複数）102は、電気信号を検出するためのセンサ（単数または複数）を含み得、かつその信号を使用して、個々の結合相互作用に参与するアミノ酸を同定し得る。

10

20

【0058】

いくつかの実施形態では、シーケンシングデバイス（単数または複数）102は、本明細書に記載されているものとは異なる技術を使用してタンパク質をシーケンシングするように構成され得る。いくつかの実施形態は、本明細書に記載のタンパク質のシーケンシングの特定の技術に限定されない。

【0059】

図1Aの実施形態に示されるように、シーケンシングデバイス（単数または複数）102は、デバイス（単数または複数）102によって生成されたシーケンシングデータを、格納のためにデータストア108Aに送信するように構成され得る。シーケンシングデータは、高分子サンプルのシーケンシングから生成された配列を含み得る。シーケンシングデータは、1つまたは複数の他のシステムによって使用され得る。一例として、シーケンシングデータは、高分子のアセンブリを生成するためにアセンブリシステム104によって使用され得る。別の例として、シーケンシングデータは、アセンブリシステム104によって使用されるための機械学習モデルをトレーニングするためのトレーニングデータとして、モデルトレーニングシステム106によって使用され得る。シーケンシングデータの使用例が本明細書に記載される。

30

【0060】

いくつかの実施形態では、アセンブリシステム104は、シーケンシングデバイス（単数または複数）102によって生成されたシーケンシングデータを使用してアセンブリ112を生成するように構成されたコンピューティングデバイスであり得る。アセンブリシステム104は、アセンブリシステム104がアセンブリを生成するために使用する機械学習モデル104Aを含む。いくつかの実施形態では、機械学習モデル104Aは、モデルトレーニングシステム106から得られるトレーニングされた機械学習モデルであり得る。アセンブリシステム104によって使用され得る機械学習モデルの例は、本明細書に記載されている。

40

【0061】

いくつかの実施形態では、アセンブリシステム104は、初期アセンブリを更新するこ

50

とによってアセンブリ 112 を生成するように構成され得る。初期アセンブリは、従来のアセンブリアルゴリズムをシーケンシングデータに適用することで取得され得る。いくつかの実施形態では、アセンブリシステム 104 は、初期アセンブリを生成するように構成され得る。アセンブリシステム 104 は、シーケンシングデバイス（単数または複数）102 から取得されたシーケンシングデータにアセンブリアルゴリズムを適用することによって初期アセンブリを生成するように構成され得る。一例として、アセンブリシステム 104 は、オーバーラップ・レイアウト・コンセンサス（OLC：Overlap Layout Consensus）アセンブリまたはド・ブラウン・グラフ（DBG：De Bruijn Graph）アセンブリを、データストア 108 A からのシーケンシングデータ（例えば、ヌクレオチド配列）に適用して、初期アセンブリを生成し得る。いくつかの実施形態では、アセンブリシステム 104 は、アセンブリシステム 104 とは別のシステムによって生成された初期アセンブリを取得するように構成され得る。一例として、アセンブリシステム 104 は、シーケンシングデバイス（単数または複数）102 によって生成されたシーケンシングデータにアセンブリアルゴリズムを適用したアセンブリシステム 104 とは別のコンピューティングデバイスによって生成された初期アセンブリを受信し得る。

10

20

30

40

50

#### 【0062】

いくつかの実施形態では、アセンブリシステム 104 は、トレーニングされた機械学習モデル 104 A を使用して、アセンブリ（例えば、アセンブリアルゴリズムの適用から取得された初期アセンブリ）を更新または改良するように構成され得る。アセンブリシステム 104 は、アセンブリ内の 1 つまたは複数のエラーを修正することによって、かつ/またはアセンブリ内の生物学的ポリマーの表示を確認することによって、アセンブリを更新するように構成され得る。いくつかの実施形態では、アセンブリシステム 104 は、（1）シーケンシングデータおよびアセンブリを使用して機械学習モデル 104 A への入力を生成すること、（2）生成された入力を機械学習モデル 104 A に提供して、対応する出力を取得すること、（3）機械学習モデル 104 A から取得された出力を使用してアセンブリを更新することによってアセンブリを更新するように構成され得る。いくつかの実施形態では、機械学習モデル 104 A の出力は、アセンブリ内の複数の位置の各々に関して、1 つまたは複数の個々の生物学的ポリマー（例えば、ヌクレオチドまたはアミノ酸）の各々がアセンブリ内のその位置に存在する 1 つまたは複数の尤度を示し得る。一例として、出力は、位置の各々に関して、個々のヌクレオチドがその位置に存在する確率を示し得る。いくつかの実施形態では、アセンブリシステム 104 は、（1）機械学習モデル 104 A から取得された出力を使用して、アセンブリの位置における生物学的ポリマー（例えば、ヌクレオチドまたはアミノ酸）を同定し、（2）位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、更新されたアセンブリを取得するように構成され得る。本明細書では、機械学習モデルを使用してアセンブリを更新するための例示的な技術に関して説明している。

#### 【0063】

いくつかの実施形態では、アセンブリシステム 104 は、更新される（例えば、修正または確認される）べきアセンブリ内の位置を識別するように構成され得る。アセンブリシステム 104 は、選択された位置を使用して機械学習モデル 104 A への入力を生成するように構成され得る。いくつかの実施形態では、アセンブリシステム 104 は、（1）個々のアセンブリの位置における生物学的ポリマーの表示が不正確である尤度を決定すること、および（2）決定された尤度に基づいて修正されるべき位置を選択することによって更新されるべき位置を識別するように構成され得る。いくつかの実施形態では、アセンブリシステム 104 は、個々の位置に示される生物学的ポリマーが不正確である尤度を示す数値を決定し、尤度値に基づいて更新されるべき位置を選択するように構成され得る。一例として、アセンブリシステム 104 は、閾値よりも大きな不正確である尤度を有する位置を選択し得る。

#### 【0064】

いくつかの実施形態では、アセンブリシステム104は、アセンブリ内の位置に関する特徴値を決定することによって、機械学習モデル104Aへの入力を生成するように構成され得る。アセンブリシステム104は、アセンブリおよびアセンブリが生成された配列を使用して特徴値を決定するように構成され得る。例示的な特徴を本明細書において記載する。いくつかの実施形態では、アセンブリシステム104は、複数の位置の各々に関して機械学習モデル104Aへの入力を生成するように構成され得る。各位置に関して、アセンブリシステム104は、特徴値を決定し、機械学習モデル104Aへの入力として特徴値を提供して、対応する出力を取得するように構成され得る。アセンブリシステム104は、位置に関して提供された入力に対応する出力を使用して、その位置に示された生物学的ポリマーを修正するか、またはその位置において示された生物学的ポリマーが正確であることを確認するように構成され得る。いくつかの実施形態では、複数の位置は、アセンブリ内の全ての位置であり得る。いくつかの実施形態では、複数の位置は、アセンブリ内の一部の位置であり得る。

10

20

30

40

50

**【0065】**

一部の位置が更新される実施形態では、アセンブリシステム104は、一部の位置を選択するように構成され得る。アセンブリシステム104は、(1)アセンブリが複数の位置において生物学的ポリマーを不正確に示す尤度を決定すること、(2)尤度を使用して、複数の位置から一部の位置を選択することを含むいくつかの方法で一部の位置を選択するように構成され得る。例えば、アセンブリシステム104は、(1)閾値の尤度を超える尤度を有する位置を特定し、(2)特定された位置を一部の位置として選択し得る。

**【0066】**

いくつかの実施形態では、アセンブリシステム104は、位置の近傍の1つまたは複数の位置において決定された特徴値を使用して修正されるべき位置に関する入力を生成するように構成され得る。選択された位置に関して、機械学習モデル104Aは、アセンブリ内の周囲の位置からのコンテキスト情報を利用して、選択された位置に関する出力を生成し得る。いくつかの実施形態では、近傍の位置は、(1)選択された位置、および(2)選択された位置の周囲の一組の位置を含み得る。一例として、近傍は、機械学習モデル104Aが出力を生成することになる選択された位置を中心とする複数の位置のウィンドウであり得る。アセンブリシステム104は、5個の位置、10個の位置、15個の位置、20個の位置、25個の位置、30個の位置、35個の位置、40個の位置、45個の位置、および/または50個の位置のウィンドウを使用し得る。

**【0067】**

いくつかの実施形態では、アセンブリシステム104は、最終的なアセンブリ112を生成するために複数の更新の反復を実行するように構成され得る。一例として、アセンブリシステム104は、(1)初期アセンブリで1回目の反復を実行して、第1の更新されたアセンブリを取得し、(2)第1の更新されたアセンブリに対して2回目の反復を実行して、第2の更新されたアセンブリを取得し得る。いくつかの実施形態では、アセンブリシステム104は、更新を反復して実行するように構成され得る。アセンブリシステム104は、条件が満たされるまで更新の反復を実行するように構成され得る。例示的な条件が本明細書において記載されている。

**【0068】**

いくつかの実施形態では、モデルトレーニングシステム106は、データストア108Aに格納されたデータにアクセスし、アクセスされたデータを使用して、アセンブリを生成する際に使用するための機械学習モデルをトレーニングするように構成されたコンピューティングデバイスであり得る。いくつかの実施形態では、モデルトレーニングシステム106は、異なるアセンブリシステムに対する別個の機械学習モデルをトレーニングするように構成され得る。個々のアセンブリシステム用にトレーニングされた機械学習モデルは、アセンブリシステムの固有の特性に合わせて調整され得る。一例として、モデルトレーニングシステム106は、(1)第1のアセンブリシステム用の第1の機械学習モデルをトレーニングし、(2)第2のアセンブリシステム用の第2の機械学習モデルをトレー

ニングするように構成され得る。アセンブリシステムの各々に対する個別の機械学習モデルは、個々のアセンブリシステムの固有のエラープロファイルに合わせて調整され得る。例えば、異なるアセンブリシステムは、初期アセンブリを生成するために異なるアセンブリアルゴリズムを採用し得、各アセンブリシステム用にトレーニングされた機械学習モデルは、アセンブリアルゴリズムのエラープロファイルに合わせて調整され得る。

【0069】

いくつかの実施形態では、モデルトレーニングシステム106は、単一のトレーニングされた機械学習モデルを複数のアセンブリシステムに提供するように構成され得る。一例として、モデルトレーニングシステム106は、複数のアセンブリシステムからのアセンブリを集約して、単一の機械学習モデルをトレーニングし得る。複数のアセンブリシステムで採用されているアセンブリ技術における変動に起因するモデルの変動を軽減するために、単一の機械学習モデルが複数のアセンブリシステムに対して正規化され得る。いくつかの実施形態では、モデルトレーニングシステム106は、複数のシーケンシングデバイスに対して単一のトレーニングされた機械学習モデルを提供するように構成され得る。一例として、モデルトレーニングシステム106は、複数のシーケンシングデバイスからのシーケンシングデータを集約し、単一の機械学習モデルをトレーニングし得る。単一の機械学習モデルは、デバイスの変動に起因するモデルの変動を軽減するために、複数のシーケンシングデバイスに対して正規化され得る。

10

【0070】

いくつかの実施形態では、モデルトレーニングシステム106は、(1)1つまたは複数の参照高分子(例えば、DNA、RNA、タンパク質)のシーケンシングから取得された生物学的ポリマー配列と、(2)参照高分子(単数または複数)の1つまたは複数の所定のアセンブリを含むトレーニングデータを使用することによって機械学習モデルをトレーニングするように構成され得る。いくつかの実施形態では、モデルトレーニングシステム106は、所定のアセンブリ内の生物学的ポリマーの表示を、機械学習モデルをトレーニングするためのラベルとして使用するように構成され得る。ラベルは、アセンブリの位置における正確な表示または所望の表示を表し得る。一例として、トレーニングデータは、生物のDNAサンプルのシーケンシングから所得されるヌクレオチド配列、および生物の所定のゲノムアセンブリを含み得る。この例では、モデルトレーニングシステム106は、所定のゲノムアセンブリ内のヌクレオチドの表示を、トレーニングデータに教師あり学習アルゴリズムを適用するためのラベルとして使用し得る。

20

30

【0071】

いくつかの実施形態では、モデルトレーニングシステム106は、外部データベースのトレーニングデータにアクセスするように構成され得る。一例として、モデルトレーニングシステム106は、(1)パシフィック・バイオサイエンス社(Pacific Biosciences)のRS II(バックバイオ(PacBio(登録商標)))データベースおよび/またはオックスフォード・ナノポア社(Oxford Nanopore)のMiniION(ONT)データベースのシーケンシングデータ、(2)米国国立バイオ技術情報センター(NCBI)の参照ゲノムデータベースの所定のゲノムアセンブリにアクセスし得る。別の例として、モデルトレーニングシステム106は、ユニットプロット(UnitProt)データベースおよび/またはヒト・プロテオーム・プロジェクト(HPP: Human Proteome Project)データベースからタンパク質シーケンシングデータおよび関連するプロテオームアセンブリにアクセスし得る。

40

【0072】

いくつかの実施形態では、モデルトレーニングシステム106は、ラベル付けされたトレーニングデータを使用して教師あり学習トレーニングアルゴリズムを適用することによって機械学習モデルをトレーニングするように構成され得る。一例として、モデルトレーニングシステム504は、確率的勾配降下法を使用することによって、深層学習モデル(例えば、ニューラルネットワーク)をトレーニングし得る。別の例として、モデルトレー

50

ニングシステム106は、コスト関数を最適化することによってサポートベクターマシン(SVM)の決定境界を特定するためにSVMをトレーニングし得る。一例として、モデルトレーニングシステム106は、(1)シーケンシングデータと、シーケンシングデータへのアセンブリアルゴリズムの適用により生成されたアセンブリとを使用して、機械学習モデルへの入力を生成し、(2)高分子の所定のアセンブリ(例えば、公開データベースからの)を使用して入力にラベルを付け、(3)生成された入力および対応するラベルに教師ありトレーニングアルゴリズムを適用し得る。

#### 【0073】

いくつかの実施形態では、モデルトレーニングシステム106は、教師なし学習アルゴリズムをトレーニングデータに適用することによって機械学習モデルをトレーニングするように構成され得る。一例として、モデルトレーニングシステム106は、k平均クラスタリングを実行することによって、クラスタリングモデルのクラスタを特定し得る。いくつかの実施形態では、モデルトレーニングシステム106は、(1)シーケンシングデータと、シーケンシングデータへのアセンブリアルゴリズムの適用により生成されたアセンブリとを使用して、機械学習モデルへの入力を生成し、(2)生成された入力に教師なし学習アルゴリズムを適用し得る。一例として、モデルトレーニングシステム106は、モデルの各クラスタが個々のヌクレオチドを表すクラスタリングモデルをトレーニングし得、クラスタ分類は、ゲノムアセンブリまたは遺伝子配列内のある位置におけるヌクレオチドを示し得る。別の例として、モデルトレーニングシステム106は、モデルの各クラスタが個々のアミノ酸を表すクラスタリングモデルをトレーニングし得、クラスタ分類は、タンパク質配列内のある位置におけるアミノ酸を示し得る。

10

20

#### 【0074】

いくつかの実施形態では、モデルトレーニングシステム106は、半教師あり学習アルゴリズムをトレーニングデータに適用することによって機械学習モデルをトレーニングするように構成され得る。いくつかの実施形態では、モデルトレーニングシステム106は、(1)教師なし学習アルゴリズム(例えば、クラスタリング)をトレーニングデータに適用することによって一組のラベル付けされていないトレーニングデータにラベルを付けること、および(2)ラベル付けされたトレーニングデータに教師あり学習アルゴリズムを適用することによって、半教師あり学習アルゴリズムをトレーニングデータに適用するように構成され得る。一例として、モデルトレーニングシステム106は、(1)シーケンシングデータと、シーケンシングデータへのアセンブリアルゴリズムの適用により生成されたアセンブリとを使用して、機械学習モデルへの入力を生成し、(2)生成された入力に教師なし学習アルゴリズムを適用して入力にラベルを付け、(3)ラベル付けされたトレーニングデータに教師あり学習アルゴリズムを適用し得る。

30

#### 【0075】

いくつかの実施形態では、機械学習モデルは、深層学習モデル(例えば、ニューラルネットワーク)を含み得る。いくつかの実施形態では、深層学習モデルは、畳み込みニューラルネットワーク(CNN: convolutional neural network)を含み得る。いくつかの実施形態では、深層学習モデルは、再帰型ニューラルネットワーク(RNN: recurrent neural network)、多層パーセプトロン、オートエンコーダ、および/またはCTC適合ニューラルネットワークモデルを含み得る。いくつかの実施形態では、機械学習モデルは、クラスタリングモデルを含み得る。一例として、クラスタリングモデルは、複数のクラスタを含み得、クラスタの各々は、生物学的ポリマー(例えば、ヌクレオチド、またはアミノ酸)に関連付けられている。

40

#### 【0076】

いくつかの実施形態では、モデルトレーニングシステム106は、複数のシーケンシングデバイスの各々に対する別個の機械学習モデルをトレーニングするように構成され得る。個々のシーケンシングデバイス用にトレーニングされた機械学習モデルは、シーケンシングデバイスの固有の特性に合わせて調整され得る。一例として、モデルトレーニングシステム106は、(1)第1のシーケンシングデバイス用の第1の機械学習モデ

50

ルをトレーニングし、(2)第2のシーケンシングデバイス用の第2の機械学習モデルをトレーニングし得る。個々のシーケンシングデバイス用にトレーニングされた機械学習モデルは、シーケンシングデバイスによって生成されたシーケンシングデータとともに使用するために最適化され得る。例えば、機械学習モデルは、シーケンシングデバイスによって使用される特定のシーケンシング技術(例えば、第三世代シーケンシング)のために最適化され得る。

**【0077】**

いくつかの実施形態では、モデルトレーニングシステム106は、以前にトレーニングされた機械学習モデルを定期的に更新するように構成され得る。いくつかの実施形態では、モデルトレーニングシステム106は、新たなトレーニングデータを使用して機械学習モデルの1つまたは複数のパラメータの値を更新することによって、以前にトレーニングされたモデルを更新するように構成され得る。いくつかの実施形態では、モデルトレーニングシステム106は、以前に取得されたトレーニングデータと新たなトレーニングデータとの組み合わせを使用して新たな機械学習モデルをトレーニングすることによって、機械学習モデルを更新するように構成され得る。

10

**【0078】**

いくつかの実施形態では、モデルトレーニングシステム106は、異なるタイプのイベントのいずれか1つに反応して機械学習モデルを更新するように構成され得る。例えば、いくつかの実施形態では、モデルトレーニングシステム106は、ユーザコマンドに反応して機械学習モデルを更新するように構成され得る。一例として、モデルトレーニングシステム106は、ユーザがトレーニングプロセスの実行を命令し得るユーザインターフェースを提供し得る。いくつかの実施形態では、モデルトレーニングシステム106は、例えば、ソフトウェアコマンドに反応して、機械学習モデルを自動的に(即ち、ユーザコマンドに反応することなく)更新するように構成され得る。別の例として、いくつかの実施形態では、モデルトレーニングシステム106は、1つまたは複数の条件の検出に反応して機械学習モデルを更新するように構成され得る。例えば、モデルトレーニングシステム106は、期間の満了を検出することに反応して、機械学習モデルを更新し得る。別の例として、モデルトレーニングシステム106は、閾値量(例えば、配列の数および/またはアセンブリの数)の新たなトレーニングデータを受信することに反応して、機械学習モデルを更新し得る。

20

30

**【0079】**

図1Aに示される例示的な実施形態では、モデルトレーニングシステム106は、アセンブリシステム104から分離されているが、いくつかの実施形態では、モデルトレーニングシステム106は、アセンブリシステム104の一部であり得る。図1Aに示される例示的な実施形態では、アセンブリシステム104は、シーケンシングデバイス(単数または複数)102から分離されているが、いくつかの実施形態では、アセンブリシステム104は、シーケンシングデバイスの構成要素であり得る。いくつかの実施形態では、シーケンシングデバイス102、モデルトレーニングシステム106、およびアセンブリシステム104は、各々、単一のシステムの構成要素であり得る。

**【0080】**

いくつかの実施形態では、データストア108Aは、データを格納するためのシステムであり得る。いくつかの実施形態では、データストア108Aは、1つまたは複数のコンピューティングデバイス(例えば、サーバ)によってホストされる1つまたは複数のデータベースを含み得る。いくつかの実施形態では、データストア108Aは、1つまたは複数の物理ストレージデバイスを含み得る。一例として、物理ストレージデバイス(単数または複数)は、1つまたは複数のソリッドステートドライブ、ハードディスクドライブ、フラッシュドライブ、および/または光学ドライブを含み得る。いくつかの実施形態では、データストア108Aは、データを格納する1つまたは複数のファイルを含み得る。一例として、データストア108Aは、データを格納する1つまたは複数のテキストファイルを含み得る。別の例として、データストア108Aは、1つまたは複数のXMLファイ

40

50

ルを含み得る。いくつかの実施形態では、データストア 108A は、コンピューティングデバイスのストレージ（例えば、ハードドライブ）であり得る。いくつかの実施形態では、データストア 108A は、クラウドストレージシステムであり得る。

#### 【0081】

いくつかの実施形態では、ネットワーク 111 は、無線ネットワーク、有線ネットワーク、またはそれらの任意の適切な組み合わせであり得る。一例として、ネットワーク 111 は、インターネットなどのワイドエリアネットワーク（WAN）であり得る。いくつかの実施形態では、ネットワーク 111 は、ローカルエリアネットワーク（LAN）であり得る。ローカルエリアネットワークは、シーケンシングデバイス（単数または複数）102、アセンブリシステム 104、モデルトレーニングシステム 106、およびデータストア 108A の間の有線接続および / または無線接続によって形成され得る。いくつかの実施形態は、本明細書に記載の特定のタイプのネットワークに限定されない。

10

#### 【0082】

図 1B は、遺伝子アセンブリを生成するように構成された場合の例示的なシステム 100 を示す。遺伝子アセンブリは、ゲノムアセンブリまたは遺伝子配列であり得る。例えば、出力されるアセンブリ 112 は、遺伝子アセンブリであり得る。シーケンシングデバイス（単数または複数）102 は、核酸サンプル 110 をシーケンシングしてヌクレオチド配列を生成するように構成され得る。一例として、シーケンシングデバイス（単数または複数）102 は、生物からの DNA サンプルをシーケンシングして、ヌクレオチド配列を生成し得る。シーケンシングデバイス（単数または複数）102 によって生成されたヌクレオチド配列は、データストア 108B に格納され得る。アセンブリシステム 104 は、機械学習モデル 104A を使用して遺伝子アセンブリを生成するように構成され得る。一例として、アセンブリシステム 104 は、（1）シーケンシングデバイス（単数または複数）102 によって生成されたヌクレオチド配列にアセンブリ技術（例えば、OLC）を適用することによって初期遺伝子アセンブリを取得し、（2）機械学習モデル 104A を使用して初期遺伝子アセンブリを更新して、遺伝子アセンブリ 112 を取得し得る。

20

#### 【0083】

図 1C は、タンパク質配列を生成するように構成された場合の例示的なシステム 100 を示す。例えば、出力されるアセンブリ 112 は、タンパク質配列であり得る。シーケンシングデバイス（単数または複数）102 は、タンパク質サンプル 110 をシーケンシングしてアミノ酸配列を生成するように構成され得る。一例として、シーケンシングデバイス（単数または複数）102 は、タンパク質からペプチドをシーケンシングして、アミノ酸配列を生成し得る。シーケンシングデバイス（単数または複数）102 によって生成されたアミノ酸配列は、データストア 108C に格納され得る。アセンブリシステム 104 は、機械学習モデル 104A を使用してタンパク質配列を生成するように構成され得る。一例として、タンパク質シーケンシングシステム 104 は、（1）シーケンシングデバイス（単数または複数）102 によって生成されたアミノ酸配列にアセンブリアルゴリズムを適用することによってタンパク質配列を取得し、（2）機械学習モデル 104A を使用してタンパク質配列を更新して、タンパク質配列を取得し得る。

30

40

#### 【0084】

図 2A は、本明細書に記載の技術のいくつかの実施形態による、アセンブリを生成するためのアセンブリシステム 200 を示す。アセンブリシステム 200 は、図 1A ~ 図 1C を参照して上記で説明したアセンブリシステム 104 であり得る。アセンブリシステム 200 は、シーケンシングデータ 202 を使用してアセンブリ 204 を生成するように構成されたコンピューティングデバイスであり得る。アセンブリシステム 200 は、特徴生成器 200A および機械学習モデル 200B を含む複数の構成要素を含む。アセンブリシステム 200C は、任意選択的に、アセンブラ 200C を含み得る。

#### 【0085】

いくつかの実施形態では、特徴生成器 200A は、機械学習モデルへの入力として提供

50

され得る1つまたは複数の特徴の値を決定するように構成され得る。特徴生成器200Aは、(1)配列データ202、および(2)アセンブリ(例えば、配列データ202へのアセンブリアルゴリズムの適用により得られる)から特徴(単数または複数)の値を決定するように構成され得る。配列データ202は、アセンブリを生成するためにアセンブリアルゴリズムによって使用される複数の配列を含み得る。いくつかの実施形態では、特徴生成器200Aは、配列の各々をアセンブリと比較することによって特徴(単数または複数)の値を決定するように構成され得る。いくつかの実施形態では、特徴生成器200Aは、配列をアセンブリの一部と整列させるように構成され得る。例えば、特徴生成器200Aは、配列をアセンブリ内の一組の位置に整列させ得、アセンブリ内の一組の位置における生物学的ポリマーの表示は、整列された配列から決定されたものである。特徴生成器200Aは、整列された配列を、アセンブリ内の一組の位置において示される生物学的ポリマー(例えば、ヌクレオチド、アミノ酸)と比較することによって、特徴(単数または複数)の値を決定するように構成され得る。特徴(単数または複数)の値を決定するための例示的な技術は、図4A~図4Cを参照して以下に説明される。

10

#### 【0086】

図2Aの実施形態に示されるように、特徴生成器200Aは、機械学習モデル200Bに提供される入力を生成するように構成され得る。いくつかの実施形態では、特徴生成器200Aは、アセンブリ内の複数の位置の各々に対して入力を生成するように構成され得る。いくつかの実施形態では、特徴生成器200Aは、複数の位置を選択し、選択された複数の位置を使用して入力を生成するように構成され得る。いくつかの実施形態では、特徴生成器200Aは、アセンブリが複数の位置において生物学的ポリマーを不正確に示す複数の尤度を決定し、決定された複数の尤度を使用して複数の位置を選択することによって複数の位置を選択するように構成され得る。いくつかの実施形態では、特徴生成器200Aは、アセンブリ内に示された生物学的ポリマーとは異なる生物学的ポリマーを特定する位置に整列された配列の数に基づいて、アセンブリが、ある位置において生物学的ポリマーを不正確に示す尤度を決定するように構成され得る。特徴生成器200Aは、尤度が閾値尤度を超えると決定されたときに、その位置に対する入力を生成するように構成され得る。

20

#### 【0087】

いくつかの実施形態では、特徴生成器200Aは、(1)ターゲット位置において同定される生物学的ポリマー、(2)ターゲット位置の近傍の1つまたは複数の他の位置において同定される生物学的ポリマーを使用して、アセンブリ内のターゲット位置に関して機械学習モデル200Bに提供される入力を生成するように構成され得る。いくつかの実施形態では、特徴生成器200Aは、ターゲット位置およびターゲット位置の近傍にある他の位置(単数または複数)における特徴値を決定するように構成され得る。近傍の他の位置(単数または複数)における特徴値は、ターゲット位置に関する出力を生成するために機械学習モデル200Aにコンテキスト情報を提供し得る。いくつかの実施形態では、近傍のサイズは、設定可能なパラメータであり得る。例えば、近傍のサイズは、ソフトウェアアプリケーションにおけるユーザ入力によって指定され得る。

30

#### 【0088】

いくつかの実施形態では、特徴生成器200Aは、ターゲット位置の近傍の位置において決定された特徴値を含むウィンドウとして入力を生成するように構成され得る。ターゲット位置の近傍は、ターゲット位置と、ターゲット位置のウィンドウ内の1つまたは複数の他の位置とを含み得る。いくつかの実施形態では、ウィンドウのサイズは、2個の位置、3個の位置、5個の位置、10個の位置、15個の位置、20個の位置、25個の位置、30個の位置、35個の位置、40個の位置、45個の位置、または50個の位置であり得る。いくつかの実施形態では、特徴生成器200Aは、60個の位置、70個の位置、80個の位置、90個の位置、または100個の位置の近傍のサイズを使用するように構成され得る。いくつかの実施形態では、ウィンドウは、ターゲット位置を中心に配置され得る。

40

50

## 【 0 0 8 9 】

いくつかの実施形態では、機械学習モデル 2 0 0 B は、図 1 A ~ 図 1 C を参照して上記で説明した機械学習モデル 1 0 4 A であり得る。図 1 A の実施形態に示されるように、機械学習モデル 2 0 0 B は、特徴生成器 2 0 0 A からの入力を受信するように構成され得る。機械学習モデル 2 0 0 B は、特徴生成器 2 0 0 A によって提供される個々の入力に対応する出力を生成するように構成され得る。機械学習モデル 2 0 0 B は、アセンブリ内の複数の位置における生物学的ポリマー（例えば、ヌクレオチドまたはアミノ酸）を同定するためにアセンブリシステム 2 0 0 によって使用される出力を生成するように構成され得る。いくつかの実施形態では、機械学習モデル 2 0 0 B は、位置に関して、複数の生物学的ポリマーの各々がその位置に存在する尤度を出力するように構成され得る。一例として、機械学習モデル 2 0 0 B は、複数のヌクレオチドの各々に関して、ヌクレオチドがその位置に存在する確率を出力し得る。別の例として、機械学習モデル 2 0 0 B は、複数のアミノ酸の各々に関して、アミノ酸がその位置に存在する確率を出力し得る。いくつかの実施形態では、アセンブリシステム 2 0 0 は、アセンブリ内のある位置における生物学的ポリマーを、機械学習モデル 2 0 0 B の出力によって示されるような、生物学的ポリマーのその位置において存在する尤度が最も高い生物学的ポリマーであると同定するように構成され得る。一例として、アセンブリシステム 2 0 0 は、複数のヌクレオチドの中から、その位置に存在する可能性が最も高いヌクレオチドを選択し得る。別の例として、アセンブリシステム 2 0 0 は、複数のアミノ酸の中から、その位置に存在する可能性が最も高いアミノ酸を選択し得る。

10

20

## 【 0 0 9 0 】

いくつかの実施形態では、アセンブリシステム 2 0 0 は、機械学習モデル 2 0 0 B から取得した出力を使用して、出力アセンブリ 2 0 4 を生成するように構成され得る。アセンブリシステム 2 0 0 は、機械学習モデル 2 0 0 B から取得された出力からアセンブリ内の位置において同定された生物学的ポリマーを使用してアセンブリを更新するように構成され得る。アセンブリシステム 2 0 0 は、アセンブリ内の位置において同定された生物学的ポリマーを示すようにアセンブリを更新して、出力アセンブリ 2 0 4 を取得するように構成され得る。一例として、アセンブリは、アセンブリ内の第 1 の位置においてアデニンを示し、アセンブリ内の第 2 の位置においてグアニンを示し得る。この例では、アセンブリシステム 2 0 0 は、( 1 ) 機械学習モデル 2 0 0 B から取得された出力を使用して、第 1 の位置におけるヌクレオチドがチミンであり、第 2 の位置におけるヌクレオチドがグアニンであることを同定し、( 2 ) アセンブリ内の第 1 の位置をチミンを示すように更新し、第 2 の位置において示されたヌクレオチドを変更せずに維持して、出力アセンブリ 2 0 4 を生成し得る。上記の例によって示されるように、アセンブリシステム 2 0 0 は、他の位置（単数または複数）における生物学的ポリマーの表示を変更せずに、機械学習モデル 2 0 0 B から取得された出力を使用して、アセンブリ内の位置（単数または複数）における生物学的ポリマーの表示を変更し得る。例えば、アセンブリシステム 2 0 0 は、アセンブリ内のある位置において同定された生物学的ポリマーが、アセンブリで示された生物学的ポリマーと一致することを決定して、更新されたアセンブリ内でその位置における表示を変更せずに維持し得る。

30

40

## 【 0 0 9 1 】

図 1 A の実施形態に示されるように、アセンブラ 2 0 0 C は、アセンブリを特徴生成器 2 0 0 A に提供するように構成され得る。いくつかの実施形態では、アセンブラ 2 0 0 C は、アセンブリアルゴリズムを（例えば、高分子サンプルのシーケンシングから受信される）配列データ 2 0 2 に適用することによって、特徴生成器 2 0 0 A に提供されるアセンブリを生成するように構成され得る。一例として、アセンブラ 2 0 0 C は、アセンブリアルゴリズムを、配列データ 2 0 2 に含まれるヌクレオチド配列に適用して、アセンブリを生成するように構成され得る。次に、アセンブリ内の位置における生物学的ポリマーを同定するための出力を取得するために機械学習モデル 2 0 0 B に提供される入力を生成するために、アセンブリが特徴生成器 2 0 0 A に提供され得る。アセンブラ 2 0 0 C によ

50

て生成されたアセンブリは、出力アセンブリ204を生成するために、機械学習モデル200Bから取得された出力を使用してアセンブリシステム200によって更新され得る。

【0092】

いくつかの実施形態では、アセンブラ200Cは、オーバーレイ・レイアウト・コンセンサス(OLC: overlay layout consensus)アルゴリズムを、配列データ202に含まれるヌクレオチド配列に適用して、アセンブリを生成するように構成され得る。シーケンシングデバイスは、核酸(単数または複数)を含む生物学的サンプルの複数のコピーをシーケンシングし得る。結果として、配列データ202は、アセンブリの各部分(例えば、一組の位置)に関して、アセンブリの一部に整列する複数の配列を含み得る。アセンブリ内の位置をカバーする配列の平均数は、配列の「カバレッジ」と呼ばれ得る。アセンブラ200Cは、(1)配列の重複領域に基づいて重複グラフを生成し、(2)重複グラフを使用して、アセンブリの個々の一部に整列する配列(「コンティグ(contigs)」とも呼ばれる)のレイアウトを生成し、(3)アセンブリの一部に整列する各組の配列に関して、アセンブリの一部を生成するために組内の配列のコンセンサスを取ることにによって、OLCアルゴリズムを配列に適用するように構成され得る。

10

【0093】

いくつかの実施形態では、アセンブラ200Cは、配列のペアを比較して、それらが生物学的ポリマー(例えば、ヌクレオチド)の1つまたは複数の同一の部分配列を含むかどうかを決定することによって、重複領域を有する配列を同定するように構成され得る。いくつかの実施形態では、アセンブラ200Cは、(1)少なくともヌクレオチドの閾値数(例えば、3、4、5、6、8、10、20、30、40、50、60、70、80、90、100、200、300、400、500)の同一の部分配列(単数または複数)を共有する配列のペアを重複配列として同定し、(2)各重複領域の長さ(即ち、ヌクレオチドの数)を決定し、(3)同定された重複配列および重複領域の長さに基づいて重複グラフを生成するように構成され得る。重複グラフは、重複する配列の個々のペアを接続する頂点およびエッジとしての配列を含み得る。決定された長さは、重複グラフにおけるエッジのラベルとして使用され得る。

20

【0094】

いくつかの実施形態では、アセンブラ200Cは、重複グラフを使用して配列を連結することによって、アセンブリの個々の一部に整列された複数組の配列のレイアウトを生成するように構成され得る。アセンブラ200Cは、配列を連結するために重複グラフを通るパスを発見するように構成され得る。一例として、アセンブラ200Cは、連結された配列を取得するためにヌクレオチドを表す一組の英数字を連結し得る。いくつかの実施形態では、アセンブラ200Cは、グリーディアルゴリズム(greedy algorithm)を重複グラフに適用して、連結された配列を同定し得る。一例として、アセンブラ200Cは、グリーディアルゴリズムを適用して、最短共通超文字列(shortest common superstring)を連結された配列として同定し得る。

30

【0095】

いくつかの実施形態では、アセンブラ200Cは、レイアウト配列を使用してアセンブリを生成するように構成され得る。いくつかの実施形態では、アセンブラ200Cは、各組がアセンブリの一部と整列する、複数の組のレイアウト配列を同定し得る。アセンブラ200Cは、アセンブリの一部と整列するレイアウト配列のコンセンサスを取ることにによって、アセンブリの一部を生成するように構成され得る。いくつかの実施形態では、アセンブラ200Cは、アセンブリの一部内のある位置における生物学的ポリマー(例えば、ヌクレオチド)が、アセンブリの一部に整列した配列の大多数がその位置にあることを示す生物学的ポリマーであると決定することによって、コンセンサスを取るように構成され得る。一例として、アセンブラ200Cは、ヌクレオチド配列の重複グラフを生成し、アセンブリ内の一組の4個の位置に対応する4個のヌクレオチド配列「TAGA」、「TAGA」、「TAGT」、「TAGA」、および「TAGC」を同定し得る。この例では、

40

50

アセンブラ 200C は、4 個のヌクレオチド配列の全てが最初の 3 個の位置が「TAG」であることを示し、ヌクレオチド配列の大多数が 4 番目の位置が「A」であることを示すので、4 個のヌクレオチド配列間のコンセンサスを「TAGA」と決定し得る。

【0096】

いくつかの実施形態では、アセンブリシステム 200 は、機械学習技術を使用して OLC アルゴリズムのコンセンサスステップを実行するように構成され得る。アセンブラ 200C がアセンブリを生成するために使用されるレイアウトを生成すると、システムは、レイアウトおよびレイアウトから取得されたコンセンサスアセンブリを使用して機械学習モデルへの入力を生成するように構成され得る。いくつかの実施形態では、アセンブリシステム 200 は、出力アセンブリ 204 を得るために、本明細書に記載の技術を使用してコンセンサスアセンブリを更新するように構成され得る。

10

【0097】

いくつかの実施形態では、アセンブラ 200C は、参照により本明細書に組み込まれる、ゲノミクス (Genomics)、第 95 巻、第 6 号、2010 年 6 月に公開された「次世代シーケンシングデータのためのアセンブリアルゴリズム (Assembly Algorithms for Next-Generation Sequencing Data)」に記載されたシーケンシングデータ 202 にアルゴリズムを適用するように構成され得る。いくつかの実施形態では、アセンブラ 200C は、OLC アルゴリズム以外のアセンブリアルゴリズムを配列データ 202 に適用してアセンブリを生成するように構成され得る。いくつかの実施形態では、アセンブラ 200C は、ド・ブラウン・グラフ (DBG) アセンブリを配列データ 202 に適用するように構成され得る。いくつかの実施形態は、特定のタイプのアセンブリアルゴリズムに限定されない。いくつかの実施形態では、アセンブラ 200C は、配列データ 202 を使用してアセンブリを生成するように構成されたソフトウェアアプリケーションを含み得る。一例として、システムは、HGAP アセンブラ、ファルコン (Falcon) アセンブラ、カヌ (Canu) アセンブラ、ヒンジ (Hinge) アセンブラ、ミニアスム (Miniasm) アセンブラ、またはフライ (Flye) アセンブラを含み得る。別の例として、システムは、SPADES アセンブリアプリケーション、レイ (Ray) アセンブリアプリケーション、ABYSS アセンブリアプリケーション、ALLPATHS-LG アセンブリアプリケーション、またはトリニティ (Trinity) アセンブリアプリケーションを含み得る。いくつかの実施形態は、特定のアセンブラに限定されない。

20

30

【0098】

図 2A の破線によって示されるように、いくつかの実施形態では、アセンブラ 200C は、アセンブリシステムに含まれなくてもよい。アセンブリシステム 200 は、別個のシステムからアセンブリを受信し、受信したアセンブリを更新して出力アセンブリ 204 を生成するように構成され得る。一例として、別個のコンピューティングデバイスは、アセンブリアルゴリズム (例えば、OLC) を配列データ 202 に適用して、アセンブリを生成し、生成されたアセンブリをアセンブリシステム 200 に送信し得る。

【0099】

図 2B は、図 2A を参照して上記のアセンブリシステム 200 の実施形態を示し、アセンブリシステム 200 は、機械学習モデル 200B から特徴生成器 200A へのフィードバック矢印によって示されるように、アセンブリに対する更新の複数の反復を実行するように構成される。いくつかの実施形態では、アセンブリシステム 200 は、第 1 の更新されたアセンブリを取得した後、機械学習モデル 200B への入力として提供され得る 1 つまたは複数の特徴の値を決定するように構成され得る。特徴生成器 200A は、(1) 配列データ 202 と、(2) アセンブリアルゴリズムの配列データ 202 への適用から取得された初期アセンブリを更新することから取得された第 1 の更新されたアセンブリとから特徴 (単数または複数) の値を決定するように構成され得る。特徴生成器 200A は、出力を得るために決定された特徴 (単数または複数) の値を機械学習モデル 200B への入力として提供するように構成され得る。アセンブリシステム 200 は、機械学習モデル 2

40

50

00Bからの出力を使用して、(1)第1の更新されたアセンブリ内の個々の位置における生物学的ポリマーを同定し、(2)個々の位置において同定された生物学的ポリマーを示すように第1の更新されたアセンブリを更新して、第2の更新されたアセンブリを取得するように構成され得る。第2の更新されたアセンブリは、アセンブリシステム200によって出力されたアセンブリ204であり得る。

#### 【0100】

いくつかの実施形態では、アセンブリシステム200は、条件が満たされるまで更新の反復を実行するように構成され得る。いくつかの実施形態では、アセンブリシステム104は、閾値回数の反復が実行されたシステムが判定するまで、更新の反復を実行するように構成され得る。いくつかの実施形態では、反復の閾値回数は、ユーザ入力(例えば、ソフトウェアコマンド、またはハードコードされた値)によって設定され得る。いくつかの実施形態では、アセンブリシステム104は、反復の閾値回数を決定するように構成され得る。一例として、アセンブリシステム200は、初期アセンブリを取得するために使用されたアセンブリ技術のタイプに基づいて、更新の反復の閾値回数を決定し得る。いくつかの実施形態では、アセンブリシステム200は、指定された停止基準が満たされるまで、アセンブリを反復して更新するように構成され得る。一例として、アセンブリシステム200は、(1)最新の更新の反復から取得された現在のアセンブリと前のアセンブリとの間の差異の数を決定し、(2)差異の数が差異の閾値数より少ない場合、および/または差異のパーセンテージが閾値パーセンテージより少ない場合、アセンブリの反復した更新を停止するように決定し得る。

10

20

#### 【0101】

図2Cは、図2Aを参照して上記のアセンブリシステム200の実施形態を示し、アセンブリシステム200は、特徴生成器200Aから機械学習モデル200Bへの複数の矢印によって示されるように、アセンブリの複数の位置を並列に修正するように構成される。図2Aを参照して説明したように、いくつかの実施形態では、特徴生成器200Aは、複数の位置の各々に関して、機械学習モデル200Bに提供される入力を生成するように構成され得る。図2Cの実施形態では、アセンブリシステム200は、アセンブリの複数の位置を並列に更新するように構成され得る。アセンブリシステム200は、(1)アセンブリ内の第1の位置を更新し、(2)アセンブリ内の第1の位置の更新を完了する前に、アセンブリ内の第2の位置の更新を開始するように構成され得る。いくつかの実施形態では、アセンブリシステム200は、複数の入力を並列に生成すること、かつ/または複数の個々の位置に対して生成された複数の入力を機械学習モデル200Bに並列に提供することによって、複数の位置を並列に更新するように構成され得る。一例として、特徴生成器200Aは、(1)機械学習モデル200Bへの第1の位置に関する第1の入力を生成および/または提供し、(2)機械学習モデル200Bから第1の入力に対応する出力を取得する前に、機械学習モデル200Bへの第2の位置に関する第2の入力を生成および/または提供し得る。

30

#### 【0102】

いくつかの実施形態では、図2Cのアセンブリシステム200は、アセンブリの複数の位置を並列に更新するように構成された複数のプロセッサを含むコンピューティングデバイスであり得る。いくつかの実施形態では、アセンブリシステム200は、マルチスレッドアプリケーションを使用するように構成され得、アプリケーションの各スレッドは、アセンブリ内の個々の位置を1つまたは複数の他のスレッドと並列に更新するように構成される。

40

#### 【0103】

図2Dは、図2Aを参照して上記のアセンブリシステム200の実施形態を示し、アセンブリシステム200は、(1)機械学習モデル200Bから特徴生成器200Aへの矢印によって示されるように、更新の複数の反復を実行し、(2)特徴生成器200Aから機械学習モデル200Bへの複数の矢印によって示されるように、アセンブリの複数の位置を並列に修正するように構成されている。いくつかの実施形態では、アセンブリシステ

50

ム 2 0 0 は、図 2 B を参照して上記のように複数の更新の反復を実行し、各更新サイクル中に、図 2 C を参照して上記のようにアセンブリ内の複数の位置を並列に更新するように構成され得る。

【 0 1 0 4 】

図 3 A は、本明細書に記載の技術のいくつかの実施形態による、生物学的ポリマーアセンブリを生成するために機械学習モデルをトレーニングするための例示的なプロセス 3 0 0 を示す。プロセス 3 0 0 は、任意の適切なコンピューティングデバイス（単数または複数）によって実行され得る。一例として、プロセス 3 0 0 は、図 1 A ~ 図 1 C を参照して説明されたモデルトレーニングシステム 1 0 6 によって実行され得る。プロセス 3 0 0 は、本明細書で説明される機械学習モデルをトレーニングするために実行され得る。一例として、プロセス 3 0 0 が、図 6 を参照して説明した畳み込みニューラルネットワーク（CNN）6 0 0 などの深層学習モデルをトレーニングするために実行され得る。

10

【 0 1 0 5 】

いくつかの実施形態では、機械学習モデルは、深層学習モデルであり得る。いくつかの実施形態では、深層学習モデルはニューラルネットワークであり得る。例として、機械学習モデルは、アセンブリ内の複数の位置における生物学的ポリマー（例えば、ヌクレオチド、アミノ酸）を同定する際に使用される出力を生成する畳み込みニューラルネットワーク（CNN）であり得る。別の例として、機械学習モデルは、CTC 適合ニューラルネットワークであり得る。いくつかの実施形態では、深層学習モデルの一部は、個別にトレーニングされ得る。一例として、深層学習モデルは、入力データを 1 つまたは複数の特徴（単数または複数）の値にエンコードする第 1 の部分と、特徴（単数または複数）の値を入力として受信して、1 つまたは複数の生物学的ポリマーを同定する出力を生成する第 2 の部分とを有し得る。

20

【 0 1 0 6 】

いくつかの実施形態では、機械学習モデルは、クラスタリングモデルであり得る。いくつかの実施形態では、モデルの各クラスタは、生物学的ポリマーに関連付けられ得る。例示的な例として、クラスタリングモデルは 5 つのクラスタを含み得、各クラスタは個々のヌクレオチドに関連付けられている。例えば、第 1 のクラスタはアデニンに関連付けられ得、第 2 のクラスタはシトシンに関連付けられ得、第 3 のクラスタはグアニンに関連付けられ得、第 4 のクラスタはチミンに関連付けられ得、第 5 のクラスタは、（例えば、アセンブリ内のある位置において）ヌクレオチドが存在しないことを示し得る。クラスタおよび関連する生物学的ポリマーの例示的な数は、例示の目的で本明細書に記載されている。

30

【 0 1 0 7 】

プロセス 3 0 0 は、ブロック 3 0 2 で開始し、プロセス 3 0 0 を実行するシステムは、1 つまたは複数の参照高分子（例えば、DNA、RNA、またはタンパク質）のシーケンシングによるシーケンシングデータにアクセスする。いくつかの実施形態では、システムは、参照高分子のシーケンシングによるシーケンシングデータにデータベースからアクセスするように構成され得る。一例として、システムは、細菌のシーケンシングにより取得されたシーケンシングデータに ONG データベースからアクセスし得る。シーケンシングデータは、高分子の 1 つまたは複数のサンプルをシーケンシングすることにより取得され得る。一例として、シーケンシングデータは、酵母の一種であるサッカロミセス・セレビシエ（*Saccharomyces cerevisiae*）の生物学的サンプルから取得され得る。別の例として、シーケンシングデータは、タンパク質のペプチドサンプルをシーケンシングすることから取得され得る。いくつかの実施形態では、シーケンシングデータは、核酸（例えば、DNA、RNA）を含む生物学的サンプルをシーケンシングすることから取得されたヌクレオチド配列を含み得る。いくつかの実施形態では、シーケンシングデータは、タンパク質サンプル（例えば、タンパク質からのペプチド）をシーケンシングすることから取得されたアミノ酸配列を含み得る。

40

【 0 1 0 8 】

いくつかの実施形態では、システムは、機械学習モデルが、ターゲットシーケンシ

50

グ技術によって生成されたシーケンシングデータから生成されたアセンブリの精度を向上させるようにトレーニングされ得るように、ターゲットシーケンシング技術によるシーケンシングデータにアクセスするように構成され得る。機械学習モデルは、機械学習モデルがターゲットシーケンシング技術の特徴的なエラーを修正するために最適化され得るように、ターゲットシーケンシング技術のエラープロファイルに関してトレーニングされ得る。いくつかの実施形態では、システムは、第三世代シーケンシングにより取得されたデータにアクセスするように構成され得る。いくつかの実施形態では、第三世代シーケンシングは、1分子リアルタイムシーケンシングであり得る。一例として、システムは、ヌクレオチドに結合された発光分子による発光を検出することによって核酸サンプルをシーケンシングするシステムから取得されたデータにアクセスし得る。別の例として、システムは、アミノ酸と選択的に相互作用する試薬に結合された発光分子による発光を検出することによってペプチドをシーケンシングするシステムから取得されたデータにアクセスし得る。いくつかの実施形態では、システムは、第2世代シーケンシングから取得されたデータにアクセスするように構成され得る。一例として、システムは、サンガー・シーケンシング (Sanger sequencing)、マキサムギルバート・シーケンシング (Maxam-Gilbert sequencing)、ショットガン・シーケンシング (shotgun sequencing)、パイロ・シーケンシング (pyrosequencing)、コンビナトリアル・プローブ・アンカー合成 (combinatorial probe anchor synthesis)、またはライゲーション (ligation) によるシーケンシングから取得された

10

20

#### 【0109】

次に、プロセス300はブロック304に移行し、システムは、ブロック302で取得されたシーケンシングデータの少なくとも一部から生成されたアセンブリにアクセスする。いくつかの実施形態では、システムは、アセンブリアルゴリズム (例えば、OLCアセンブリ、DBGアセンブリ) のシーケンシングデータへの適用により取得されたアセンブリにアクセスするように構成され得る。いくつかの実施形態では、システムは、アセンブリアルゴリズムをシーケンシングデータに適用することによってアセンブリにアクセスするように構成され得る。いくつかの実施形態では、システムは、1つまたは複数のアセンブリアルゴリズムのシーケンシングデータへの適用により生成された所定のアセンブリにアクセスするように構成され得る。一例として、アセンブリは、以前に別のコンピューティングデバイスによって実行され、データベースに格納されてもよい。例えば、シーケンシングデータが取得されたデータベースは、1つまたは複数のアセンブリアルゴリズムのシーケンシングデータへの適用により生成されたアセンブリをも格納し得る。

30

40

#### 【0110】

いくつかの実施形態では、システムは、ターゲットアセンブリ技術により生成されたアセンブリにアクセスするように構成され得、機械学習モデルは、ターゲットアセンブリ技術の特徴的なエラーを修正するようにトレーニングされ得る。機械学習モデルは、機械学習モデルがターゲットアセンブリ技術の特徴的なエラーを修正するために最適化され得るように、ターゲットアセンブリ技術のエラープロファイルに関してトレーニングされ得る。いくつかの実施形態では、システムは、特定のアセンブリアルゴリズムおよび/またはソフトウェアアプリケーションによって生成されたアセンブリにアクセスするように構成され得る。例として、システムは、カナ (Canu) アセンブラ、ミニアスム (Miniasm) アセンブラ、またはフライ (Flye) アセンブラによって生成されたアセンブ

50

りにアクセスし得る。いくつかの実施形態では、システムは、アセンブラのクラスから生成されたアセンブリにアクセスするように構成され得る。一例として、システムは、グリーディ・アルゴリズムアセンブラまたはグラフメソッド・アセンブラから生成されたアセンブリにアクセスし得る。いくつかの実施形態は、特定のアセンブリ技術に限定されない。

#### 【0111】

次に、プロセス300は、ブロック306に移行し、システムは、参照高分子（単数または複数）の1つまたは複数の所定のアセンブリにアクセスする。いくつかの実施形態では、参照高分子（単数または複数）の所定のアセンブリは、個々の高分子（単数または複数）に関する真のまたは正確なアセンブリを表し得る。従って、システムは、参照高分子（単数または複数）の所定のアセンブリを使用してトレーニングデータにラベルを付けるように構成され得る。一例として、システムは、NCBIデータベースから生物のDNAの参照ゲノムにアクセスし得る。この例では、システムは参照ゲノムを使用して、ゲノムアセンブリ内のヌクレオチドを同定するための機械学習モデルをトレーニングするための教師あり学習の実行の際に使用するラベルを決定し得る。別の例として、システムは、ユニットプロット（UnitProt）データベースからタンパク質の参照タンパク質配列にアクセスし、参照タンパク質配列を使用して、タンパク質配列内のアミノ酸を同定するための機械学習モデルをトレーニングするための教師あり学習の実行の際に使用するラベルを決定し得る。

10

#### 【0112】

次に、プロセス300はブロック308に移行し、システムは、ブロック302～308でアクセスされるデータを使用して機械学習モデルをトレーニングする。いくつかの実施形態では、システムは、（1）ブロック302においてアクセスされたシークエンシングデータおよびブロック304においてアクセスされたアセンブリを使用して、機械学習モデルへの入力を生成し、（2）ブロック306においてアクセスされた所定のアセンブリを使用して、生成された入力にラベルを付け、（3）ラベル付けされたトレーニングデータに教師あり学習アルゴリズムを適用するように構成され得る。いくつかの実施形態では、システムは、シークエンシングデータを使用して1つまたは複数の特徴の値を生成することによって、機械学習モデルへの入力を生成するように構成され得る。いくつかの実施形態では、システムは、アセンブリ内の各位置に対する特徴（単数または複数）の値を決定するように構成され得る。一例として、システムは、（1）個々のヌクレオチドに対するカウントを決定し、各カウントは、ヌクレオチドがその位置に存在することを示すヌクレオチド配列の数を示し、（2）カウントを使用して特徴（単数または複数）の値を決定することによって、位置に関する特徴の値を決定し得る。入力を生成して、入力にラベルを付けるための例示的な技術は、図4A～図4Cを参照して本明細書に記載されている。

20

30

#### 【0113】

いくつかの実施形態では、システムは、ラベル付けされたトレーニングデータを使用して深層学習モデルをトレーニングするように構成され得る。いくつかの実施形態では、システムは、ラベル付けされたトレーニングデータを使用して決定木モデルをトレーニングするように構成され得る。いくつかの実施形態では、システムは、ラベル付けされたトレーニングデータを使用してサポートベクターマシン（SVM：support vector machine）をトレーニングするように構成され得る。いくつかの実施形態では、システムは、ラベル付けされたトレーニングデータを使用してナイーブベイズ分類器（NBC：Naive Bayes classifier）をトレーニングするように構成され得る。

40

#### 【0114】

いくつかの実施形態では、システムは、確率的勾配降下法を使用することによって機械学習モデルをトレーニングするように構成され得る。システムは、目的関数を最適化するために機械学習モデルのパラメータを反復的に変更して、トレーニングされた機械学習モ

50

デルを取得し得る。例えば、システムは確率的勾配降下法を使用して、畳み込みネットワークのフィルタおよび/またはニューラルネットワークの重みをトレーニングし得る。

【0115】

いくつかの実施形態では、システムは、ラベル付けされたトレーニングデータを使用して教師ありトレーニングを実行するように構成され得る。いくつかの実施形態では、システムは、(1)機械学習モデルに生成された入力を提供して、対応する出力を取得し、(2)出力を使用してアセンブリ内の複数の位置に存在する生物学的ポリマーを同定し、(2)同定された生物学的ポリマーと参照アセンブリの複数の位置において示されている生物学的ポリマーとの間の差異に基づいて機械学習モデルをトレーニングすることによって機械学習モデルをトレーニングするように構成され得る。参照アセンブリ内のある位置において示される生物学的ポリマーは、個々の入力に関するラベルであり得る。差異は、機械学習モデルが、現在の組のパラメータで構成された場合に、ラベルを再現する際にどの程度良好に動作するかの尺度を提供し得る。例として、機械学習モデルのパラメータは、確率的勾配降下法および/またはモデルのトレーニングに適した他の反復最適化手法を使用して更新され得る。一例として、システムは、決定された差異に基づいてモデルの1つまたは複数のパラメータを更新するように構成され得る。

10

【0116】

いくつかの実施形態では、システムは、教師なしトレーニングアルゴリズムを一組のラベル付けされていないトレーニングデータに適用し得る。図3Aの実施形態は、ブロック306において参照高分子の所定のアセンブリにアクセスすることを含むが、いくつかの実施形態では、システムは、所定のアセンブリにアクセスすることなくトレーニングを実行するように構成され得る。これらの実施形態では、システムは、教師なしトレーニングアルゴリズムをトレーニングデータに適用して、機械学習モデルをトレーニングするように構成され得る。システムは、(1)シークエンシングデータと、シークエンシングデータから生成されたアセンブリとを使用してモデルへの入力を生成し、(2)生成された入りに教師なしトレーニングアルゴリズムを適用することによってモデルをトレーニングするように構成され得る。いくつかの実施形態では、機械学習モデルはクラスタリングモデルであり得、システムは、教師なし学習アルゴリズムをトレーニングデータに適用することによって、クラスタリングモデルのクラスタを識別するように構成され得る。各クラスタは、生物学的ポリマー(例えば、ヌクレオチドまたはアミノ酸)と関連付けられ得る。一例として、システムは、トレーニングデータを使用してk平均クラスタリングを実行して、クラスタ(例えば、クラスタ重心)を識別し得る。

20

30

【0117】

いくつかの実施形態では、システムは、半教師あり学習アルゴリズムをトレーニングデータに適用するように構成され得る。システムは、(1)教師なし学習アルゴリズム(例えば、クラスタリング)をトレーニングデータに適用することによって一組のラベル付けされていないトレーニングデータにラベルを付け、(2)ラベル付けされたトレーニングデータに教師あり学習アルゴリズムを適用し得る。一例として、システムは、シークエンシングデータから生成された入力およびシークエンシングデータから取得されたアセンブリにk平均クラスタリングを適用して、入力をクラスタリングし得る。次に、システムは、クラスタメンバーシップに基づく分類によって各入力にラベルを付け得る。次に、システムは、確率的勾配降下アルゴリズムおよび/または他の反復最適化手法をラベル付けされたデータに適用することによって、機械学習モデルをトレーニングし得る。

40

【0118】

ブロック308において機械学習モデルをトレーニングした後、プロセス300は終了する。いくつかの実施形態では、システムは、トレーニングされた機械学習モデルを格納するように構成され得る。システムは、機械学習モデルの1つまたは複数のトレーニングされたパラメータの値(単数または複数)を保存し得る。一例として、機械学習モデルは、1つまたは複数のニューラルネットワークを含み得、システムは、ニューラルネットワーク(単数または複数)のトレーニングされた重みの値を格納し得る。別の例として、機

50

械学習モデルは畳み込みニューラルネットワークを含み、システムは畳み込みニューラルネットワークの1つまたは複数のトレーニングされたフィルタを格納し得る。いくつかの実施形態では、システムは、アセンブリ（例えば、ゲノムアセンブリ、タンパク質配列、またはそれらの一部）を生成するために使用するためのトレーニングされた機械学習モデルを（例えば、アセンブリシステム104内に）格納するように構成され得る。

#### 【0119】

いくつかの実施形態では、システムは、新たなデータを取得し、新たなトレーニングデータを使用して機械学習モデルを更新するように構成され得る。いくつかの実施形態では、システムは、新たなトレーニングデータを使用して新たな機械学習モデルをトレーニングすることによって機械学習モデルを更新するように構成され得る。一例として、システムは、新たなトレーニングデータを使用して新たな機械学習モデルをトレーニングし得る。いくつかの実施形態では、システムは、新たなトレーニングデータを使用して機械学習モデルを再トレーニングして、機械学習モデルの1つまたは複数のパラメータを更新することによって機械学習モデルを更新するように構成され得る。一例として、モデルによって生成された出力（単数または複数）および対応する入力データは、以前に取得されたトレーニングデータとともにトレーニングデータとして使用され得る。いくつかの実施形態では、システムは、（例えば、図3Bを参照して以下に説明するプロセス310を実行することから得られる）アミノ酸を同定するデータおよび出力を使用して、トレーニングされた機械学習モデルを反復して更新するように構成され得る。一例として、システムは、第1のトレーニングされた機械学習モデル（例えば、教師モデル）に入力データを提供して、1つまたは複数のアミノ酸を同定する出力を取得するように構成され得る。次に、システムは、入力データおよび対応する出力を使用して機械学習モデルを再トレーニングして、第2のトレーニングされた機械学習モデル（例えば、学生モデル）を取得し得る。

10

20

#### 【0120】

いくつかの実施形態では、システムは、複数のシーケンシング技術の各々に関して別個の機械学習モデルをトレーニングするように構成され得る。機械学習モデルは、シーケンシング技術から取得したデータを使用して、個々のシーケンシング技術に関してトレーニングされ得る。機械学習モデルは、シーケンシング技術のエラープロファイルに関して調整され得る。いくつかの実施形態では、システムは、複数のアセンブリ技術の各々に関して別個の機械学習モデルをトレーニングするように構成され得る。機械学習モデルは、アセンブリ技術から取得したアセンブリを使用して、個々のアセンブリ技術に関してトレーニングされ得る。機械学習モデルは、アセンブリ技術のエラープロファイルに関して調整され得る。

30

#### 【0121】

いくつかの実施形態では、システムは、複数のシーケンシング技術に関して使用される一般化された機械学習モデルをトレーニングするように構成され得る。一般化された機械学習モデルは、複数のシーケンシング技術から集約されたデータを使用してトレーニングされ得る。いくつかの実施形態では、システムは、複数のアセンブリ技術に関して使用される一般化された機械学習モデルをトレーニングするように構成され得る。一般化された機械学習モデルは、複数のアセンブリ技術を使用して生成されたアセンブリを使用してトレーニングされ得る。

40

#### 【0122】

図3Bは、本明細書に記載の技術のいくつかの実施形態による、アセンブリ（例えば、ゲノムアセンブリ、遺伝子配列、タンパク質配列、またはそれらの一部）を生成するためのプロセス300から取得されたトレーニングされた機械学習モデルを使用するための例示的なプロセス310を示す。プロセス310は、任意の適切なコンピューティングデバイスによって実行され得る。一例として、プロセス310は、図1A～図1Cを参照して上記のアセンブリシステム104によって実行され得る。

#### 【0123】

プロセス310は、ブロック312で開始し、システムは、アセンブリを生成するため

50

に、シーケンシングデータに対するアセンブリアルゴリズム（例えば、OLCアセンブリまたはDBGアセンブリ）を実行する。一例として、システムは、DNAサンプルのシーケンシングから生成されたヌクレオチド配列に対してアセンブリアルゴリズムを適用し得る。別の例として、システムは、タンパク質からのペプチドサンプルのシーケンシングから生成されたアミノ酸配列にアセンブリアルゴリズムを適用し得る。システムは、図2A～図2Dのアセンブラ200Cを参照して、上記のようなアセンブリアルゴリズムを適用し得る。いくつかの実施形態では、システムは、アセンブリアプリケーションを含み得る。システムは、アセンブリアプリケーションを実行することによってアセンブリを生成するように構成され得る。アセンブリアプリケーションの例は、本明細書に記載されている。

10

**【0124】**

ブロック312の周囲の破線によって示されるように、いくつかの実施形態では、システムは、アセンブリアルゴリズムを実行しなくてもよい。システムは、別個のシステム（例えば、別個のコンピューティングデバイス）によって生成されたアセンブリを取得し、ブロック314～322のステップを実行して、取得されたアセンブリを更新し得る。

**【0125】**

次に、プロセス310は、ブロック312に移行し、システムがシーケンシングデータおよびアセンブリにアクセスする。いくつかの実施形態では、システムは、（例えば、ブロック312において）システムによって生成されたアセンブリにアクセスするように構成され得る。いくつかの実施形態では、システムは、別個のシステムによって生成されたアセンブリにアクセスするように構成され得る。一例として、システムは、システムとは別のコンピューティングデバイス上で実行されるソフトウェアアプリケーションによって生成されたアセンブリを受信し得る。いくつかの実施形態では、システムは、プロセス300でトレーニングされた機械学習モデルが更新するのに（例えば、エラーを修正するのに）最適化されたターゲットアセンブリ技術（例えば、アルゴリズムおよび/またはソフトウェアアプリケーション）から生成されたシーケンシングデータにアクセスするように構成され得る。例として、機械学習モデルは、カヌ（Canu）アセンブリアプリケーションから生成されたアセンブリでトレーニングされ、システムは、カヌアセンブリアプリケーションによって生成されたアセンブリにアクセスし得る。

20

**【0126】**

いくつかの実施形態では、システムは、アクセスされたアセンブリを生成するために使用された生物学的ポリマー配列を含むシーケンシングデータにアクセスするように構成され得る。一例として、アクセスされるシーケンシングデータは、ゲノムアセンブリまたは遺伝子配列を生成するためにアセンブリアルゴリズムが適用されたヌクレオチド配列を含み得る。別の例として、アクセスされるシーケンシングデータは、タンパク質配列を生成するためにアセンブリアルゴリズムが適用されたアミノ酸配列を含み得る。いくつかの実施形態では、システムは、プロセス300でトレーニングされた機械学習モデルが更新するのに最適化されたターゲットシーケンシング技術から生成されたシーケンシングデータにアクセスするように構成され得る。例として、機械学習モデルは、第三世代シーケンシングから生成されたシーケンシングデータでトレーニングされ得、システムは、第三世代シーケンシングから生成されたシーケンシングデータにアクセスし得る。

30

40

**【0127】**

次に、プロセス310は、ブロック316に移行し、システムは、シーケンシングデータおよびアセンブリを使用して、機械学習モデルに提供される入力を生成する。いくつかの実施形態では、システムは、アセンブリ内の個々の位置に関する入力を生成するように構成され得る。システムは、（1）シーケンシングデータからの配列をアセンブリ内の一組の位置に整列させ、（2）整列された配列の生物学的ポリマーを、アセンブリ内の位置に示される生物学的ポリマーと比較して、1つまたは複数の特徴の値を決定することによって、アセンブリ内の一組の位置に関する入力を生成するように構成し得る。いくつ

50

かの実施形態では、システムは、アセンブリ内の一組の位置における生物学的ポリマーを示すシーケンシングデータからの配列を同定することによって、アセンブリ内の一組の位置に配列を整列させるように構成され得る。一例として、アセンブリは、1から10,000のインデックスが付けられた位置を含み得、システムは、ヌクレオチド配列「TAGGTC」、「TAGTTC」、「TAGGCC」、「TAGGTC」が各々、アセンブリの5~10にインデックスが付けられた位置に整列することを決定し得る。この例では、システムは、各ヌクレオチド配列を、アセンブリ内の5~10にインデックスが付けられた位置において示された生物学的ポリマーと比較して、特徴(単数または複数)の値を決定し得る。特徴の例、および特徴の値の生成は、図4A~図4Cを参照して説明されている。

10

**【0128】**

いくつかの実施形態では、システムは、アセンブリ内の個々の位置に関する入力を作成するように構成され得る。システムは、機械学習モデルへの入力として提供する位置に関する入力を生成して、アセンブリ内の位置に存在する生物学的ポリマー(例えば、ヌクレオチド、アミノ酸)を同定するために使用され得る出力を取得するように構成され得る。いくつかの実施形態では、システムは、その位置における生物学的ポリマーの表示、およびその位置の近傍にある1つまたは複数の他の位置における生物学的ポリマーの表示に基づいて、アセンブリ内のある位置に関する入力を生成するように構成され得る。入力は、モデルが対応する出力を生成するために使用するアセンブリ内の位置の周囲のコンテキスト情報を機械学習モデルに提供し得る。システムは、その位置およびその位置の近傍の他の位置(単数または複数)における特徴(単数または複数)の値を決定することによって、その位置の近傍の位置における生物学的ポリマーの表示に基づいて、ある位置に関する入力を生成するように構成され得る。一例として、システムは、(1)位置を選択し、(2)選択された位置を中心とする近傍の位置を特定し、(3)選択された位置および近傍の位置の各々における特徴(単数または複数)の値である入力を生成し得る。

20

**【0129】**

いくつかの実施形態では、システムは、設定された近傍のサイズを使用するように構成され得る。本明細書において近傍のサイズの例が説明される。いくつかの実施形態では、システムによって使用される近傍の位置の数は、設定可能なパラメータであり得る。例えば、システムは、使用する近傍のサイズを指定するユーザ入力(例えば、ソフトウェアアプリケーションにおける)を受信し得る。いくつかの実施形態では、システムは、近傍のサイズを決定するように構成され得る。一例として、システムは、シーケンシングデータが生成されたシーケンシング技術および/またはアセンブリが生成されたアセンブリ技術に基づいて近傍のサイズを決定し得る。

30

**【0130】**

いくつかの実施形態では、システムは、(1)アセンブリ内の位置を選択し、(2)選択された位置に関する個々の入力を生成することによって機械学習モデルに提供される入力を生成するように構成され得る。いくつかの実施形態では、システムは、アセンブリがアセンブリ内の位置において生物学的ポリマーを不正確に示す尤度を決定し、決定された尤度を使用して入力を生成する位置を選択することによって、アセンブリ内の位置を選択するように構成され得る。一例として、システムは、アセンブリが位置において生物学的ポリマーを不正確に示す尤度が閾値尤度を超えるかどうかを決定し、尤度が閾値尤度を超える場合、その位置に関する入力を生成し得る。いくつかの実施形態では、システムは、生物学的ポリマーがその位置に存在することを示す整列された配列の数に基づいて、位置が生物学的ポリマーを不正確に示す尤度を決定するように構成され得る。システムは、生物学的ポリマーがその位置にあることを示す配列の数と配列の総数との間の差異である尤度を決定し得る。一例として、アセンブリは、一組の9個のヌクレオチド配列のからのコンセンサスに基づいて、アセンブリ内のある位置においてチミンを示し得、このとき、4個のヌクレオチド配列は、チミンがその位置に存在することを示し、2個のヌクレオチド配列は、グアニンがその位置に存在することを示し、3個のヌクレオチド配列は、アデニ

40

50

ンがその位置に存在することを示す。この例では、システムは、アセンブリが、アセンブリ内の位置にある生物学的ポリマーを、チミンを示すヌクレオチド配列の数(4)とヌクレオチド配列の総数(9)との間に差異があると不正確に示す尤度を決定して、5の値を取得し得る。システムは、5が閾値の差異(例えば、1、2、3、4)より大きいと判定し、その結果、位置に関する入力を生成し得る。

#### 【0131】

いくつかの実施形態では、システムは、1、2、3、4、5、6、7、8、9、または10の閾値の差異を使用するように構成され得る。いくつかの実施形態は、特定の閾値の差異に限定されない。いくつかの実施形態では、閾値の差異は、設定可能なパラメータであり得る。システムによって使用される閾値尤度は、システムがモデルに提供される入力を生成する位置の数に影響を与え得る。一例として、システムは、ソフトウェアアプリケーションへのユーザ入力として閾値の値を受信し得る。いくつかの実施形態では、システムは、設定された閾値尤度を使用し得る。一例として、閾値尤度の値がエンコードされ得る。いくつかの実施形態では、システムは、閾値尤度を自動的に決定するように構成され得る。一例として、システムは、アセンブリが生成されたアセンブリ技術および/またはシーケンシングデータが生成されたシーケンシング技術に基づいて閾値尤度を決定し得る。

#### 【0132】

いくつかの実施形態では、システムは、位置に関する入力を2次元行列として生成するように構成され得る。いくつかの実施形態では、マトリクスの各行/列は、アセンブリ内の個々の位置において決定された特徴(単数または複数)の値を指定し得る。いくつかの実施形態では、システムは、画像として入力を生成するように構成され得、画像のピクセルは、特徴(単数または複数)の値を保持する。一例として、画像の各行/列は、アセンブリ内の個々の位置において決定された特徴(単数または複数)の値を指定し得る。

#### 【0133】

次に、プロセス310は、ブロック318に移行し、システムは、対応する出力を取得するためにブロック316で生成された入力を機械学習モデルに提供する。いくつかの実施形態では、システムは、機械学習モデルへの別個の入力として、アセンブリ内の個々の位置に対して生成された入力を提供するように構成され得る。一例として、システムは、ターゲット位置に対応する出力を取得するために、機械学習モデルへの入力として、ターゲット位置およびその位置の近傍の位置において決定された一組の特徴値を提供し得る。いくつかの実施形態では、システムは、(例えば、図2C~図2Dを参照して上で説明したように)複数の位置に対して並列に生成された入力を提供するように構成され得る。一例として、システムは、(1)第1の位置に対して生成された第1の入力をモデルに提供し、(2)第1の入力に対応する第1の出力を取得する前に、第2の位置に対して生成された第2の入力をモデルに提供し得る。いくつかの実施形態では、システムは、複数の位置に対して生成された入力を順次提供するように構成され得る。例えば、システムは、(1)対応する第1の出力を取得するために、第1の位置に対して生成された第1の入力をモデルに提供し、(2)第1の出力を取得した後、対応する第2の出力を取得するために、第2の位置に対する第2の入力を提供し得る。

#### 【0134】

いくつかの実施形態では、機械学習モデルに提供される入力に対応する出力は、アセンブリ内の複数の位置の各々に関して、1つまたは複数の生物学的ポリマーの各々がその位置に存在する尤度を示し得る。一例として、出力は、ゲノムアセンブリ内の複数の位置の各々に関して、1つまたは複数のヌクレオチド(例えば、アデニン、グアニン、チミン、シトシン)の各々がその位置に存在する尤度(例えば、確率)を示し得る。別の例として、出力は、タンパク質配列内の複数の位置の各々に関して、1つまたは複数のアミノ酸の各々がその位置に存在する尤度を示し得る。いくつかの実施形態では、出力は、アセンブリ内のある位置に生物学的ポリマーが存在しない尤度を示し得る。一例として、システムは、「-」文字がアセンブリ内の位置における尤度を示し得る。

10

20

30

40

50

## 【 0 1 3 5 】

いくつかの実施形態では、モデルは、アセンブリ内の個々の位置に対応する出力を提供し得る。システムは、アセンブリ内のターゲット位置に対して生成された入力を提供し、ターゲット位置に存在する1つまたは複数の生物学的ポリマーの各々の尤度を示す対応する出力を取得し得る。一例として、システムは、ゲノムアセンブリ内の位置に対して生成された入力を提供し、一組の4つの可能性のあるヌクレオチド（例えば、アデニン、グアニン、チミン、シトシン）の各々がその位置に存在する尤度を示す対応する出力を取得し得る。例えば、尤度は、その位置に存在する各ヌクレオチドの確率値であり得る。

## 【 0 1 3 6 】

次に、プロセス310は、ブロック320に移行し、システムは、モデルから取得された出力を使用して、アセンブリ内の位置における生物学的ポリマーを同定する。いくつかの実施形態では、システムは、モデルに提供された対応する入力に回答してその位置に対して取得された出力を使用して、位置の各々に関して、その位置に存在する生物学的ポリマーを特定することによって、アセンブリ内の位置における生物学的ポリマーを特定するように構成され得る。モデルからの出力は、個々の位置に対応する複数組の出力値を含み得る。各組の出力値は、1つまたは複数の生物学的ポリマーの各々がアセンブリ内の個々の位置に存在する尤度を指定し得る。システムは、個々の位置においてその位置に存在する尤度が最も高い生物学的ポリマーである生物学的ポリマーを同定し得る。例として、アセンブリ内の第1の位置に関する一組の出力値は、アデニン(A) 0.1、シトシン(C) 0.6、グアニン(G) 0.1、チミン(T) 0.15、およびブランク(-) 0.05の組のその位置に関する尤度を示し得る。この例では、システムは、アセンブリ内の位置にあるシトシン(C)を同定し得る。いくつかの実施形態では、位置に関して生成された入力に対応するモデルからの出力は、その位置において生物学的ポリマーを指定する分類であり得る。一例として、モデルからの出力は、アデニン(A)、シトシン(C)、グアニン(G)、チミン(T)、またはブランク(-)の分類であり得る。

## 【 0 1 3 7 】

次に、プロセス310はブロック322に移行し、システムは、アセンブリを更新して、更新されたアセンブリを取得する。システムは、ブロック320において同定された生物学的ポリマーに基づいてアセンブリを更新するように構成され得る。いくつかの実施形態では、システムは、アセンブリ内の位置における生物学的ポリマーの表示を更新することによってアセンブリを更新するように構成され得る。いくつかの例では、ブロック320において位置に存在すると同定された生物学的ポリマーは、アセンブリ内の生物学的ポリマーの表示とは異なり得る。これらの例では、システムは、アセンブリ内の位置における生物学的ポリマーの表示を変更し得る。一例として、システムは、(1)モデルの出力を使用して、アデニン「A」の表示を有するアセンブリ内の第1の位置にチミン「T」が存在することを同定し、(2)アデニン「A」の以前の表示からチミン「T」を表示するようにアセンブリ内の第1の位置を変更し得る。いくつかの例では、ある位置に存在すると同定された生物学的ポリマーは、アセンブリ内のその位置における生物学的ポリマーの表示と同じであり得る。これらの例では、システムは、アセンブリ内のその位置における生物学的ポリマーの表示を変更しない。一例として、システムは、(1)モデルの出力を使用して、チミン「T」の表示を有するアセンブリ内の第1の位置においてチミン「T」が存在していることを同定し、(2)第1の位置の表示を変更せずに維持し得る。

## 【 0 1 3 8 】

いくつかの実施形態では、システムは、アセンブリ内の複数の位置を並列に更新するように構成され得る。一例として、システムは、(1)アセンブリ内の第1の位置の更新を開始し、(2)第1の位置における更新を完了する前に、アセンブリの第2の位置の更新を開始し得る。いくつかの実施形態では、システムは、アセンブリ内の位置を順次更新するように構成され得る。一例として、システムは、(1)アセンブリの第1の位置を更新し、(2)アセンブリの第1の位置における更新を完了した後、アセンブリの第2の位置を更新する。

10

20

30

40

50

## 【 0 1 3 9 】

いくつかの実施形態では、ブロック 3 2 2 においてアセンブリを更新して第 1 の更新されたアセンブリを取得した後、プロセス 3 1 0 は、ブロック 3 2 2 からブロック 3 1 6 への破線によって示されるように、ブロック 3 1 6 に戻ってもよい。いくつかの実施形態では、システムは、第 1 の更新されたアセンブリおよびシーケンシングデータを使用して機械学習モデルへの入力を生成するように構成され得る。一例として、システムは、シーケンシングデータの一組のヌクレオチド配列および第 1 の更新されたアセンブリを使用して、モデルへの入力を生成し得る。システムは、ヌクレオチド配列を第 1 の更新されたアセンブリの個々の位置に整列させて、上記のように機械学習モデルへの入力を生成し得る。次に、システムは、ブロック 3 1 6 から 3 2 2 における動作を実行して、第 2 の更新されたアセンブリを取得し得る。いくつかの実施形態では、アセンブリシステムは、条件が満たされるまで反復を実行するように構成され得る。

10

## 【 0 1 4 0 】

いくつかの実施形態では、システムは、閾値の反復回数が行われたとシステムが判定するまで、更新の反復を実行するように構成され得る。いくつかの実施形態では、反復の閾値回数は、ユーザ入力（例えば、ソフトウェアコマンド、またはハードコードされた値）によって設定され得る。いくつかの実施形態では、システムは、反復の閾値回数を決定するように構成され得る。一例として、システムは、初期アセンブリを取得するために使用されたアセンブリ技術のタイプに基づいて、更新の反復の閾値回数を決定し得る。いくつかの実施形態では、システムは、アセンブリが収束したことをシステムが検出するまで更新の反復を実行するように構成され得る。一例として、アセンブリシステムは、（ 1 ）最新の反復から取得された現在のアセンブリと前のアセンブリとの間の差異の数を決定し、（ 2 ）差異の数が差異の閾値数または差異のパーセンテージよりも少ない場合、更新の反復の実行を停止するように決定し得る。

20

## 【 0 1 4 1 】

いくつかの実施形態では、システムは、アセンブリへの単一の更新を実行するように構成され得、プロセス 3 1 0 は、アセンブリへの単一の更新を実行した後、ブロック 3 2 2 において終了し得る。更新されたアセンブリは、システムによって出力アセンブリとして出力され得る。一例として、システムは、出力アセンブリがブロック 3 1 4 においてアクセスされる初期アセンブリよりも正確であるように、アセンブリ内のエラーが修正されたゲノムアセンブリを出力し得る。別の例として、システムは、出力タンパク質配列がブロック 3 1 4 においてアクセスされる初期タンパク質配列よりも正確であるように、エラーが修正されたタンパク質配列を出力し得る。

30

## 【 0 1 4 2 】

いくつかの実施形態では、システムは、アセンブリの第 1 の部分に対して第 1 の数の更新の反復を実行し、アセンブリの第 2 の部分に対して第 2 の数の更新の反復を実行するように構成され得る。例として、システムは、（例えば、ブロック 3 1 6 ~ 3 2 2 で動作の複数の反復を実行することによって）ゲノムアセンブリの 1 ~ 1 0 0 のインデックスが付けられた位置を複数回更新し、（例えば、ブロック 3 1 6 ~ 3 2 2 で動作を 1 回実行することによって）ゲノムアセンブリの 1 0 1 ~ 2 0 0 のインデックスが付けられた位置を 1 回更新する。システムは、生物学的ポリマーを不正確に示し得る一部内の位置の数に基づいて、複数回更新するためのアセンブリの一部を決定するように構成され得る。一例として、システムは、（ 1 ）ウィンドウ位置（例えば、25 個、50 個、75 個、100 個、または 1000 個の位置）内で閾値の尤度を超える不正確な生物学的ポリマーの表示の尤度を有する位置の数を決定し、（ 2 ）数が位置の閾値数を越えたときに、ウィンドウ位置に対して更新サイクルを実行することを決定し得る。

40

## 【 0 1 4 3 】

図 4 A ~ 図 4 C は、本明細書に記載の技術のいくつかの実施形態による、機械学習モデルに提供される入力を生成する例を示す。

図 4 A は、ヌクレオチド配列 4 0 1（図 4 A において「パイルアップ」とラベル付けさ

50

れている)、ヌクレオチド配列 401 から生成された生物学的ポリマーのアセンブリ 402、およびアセンブリ内の個々の位置に関する生物学的ポリマーのラベル 404 を含むアレイ 400 を示す。一例として、図 4 A に示されるデータは、機械学習モデルをトレーニングするためのプロセス 300 を実行することから取得されたトレーニングデータであり得、(1) シークエンシングデータ 401 およびアセンブリ 402 は、ブロック 302 および 304 において取得され、(2) ラベル 404 は、ブロック 306 において取得される。別の例として、シークエンシングデータ 401 およびアセンブリ 402 は、トレーニングされた機械学習モデルを使用してアセンブリを生成するために、プロセス 310 のブロック 312 および / または 314 において取得され得る。

#### 【0144】

図 4 A の実施形態に示されるように、シークエンシングデータ 401 は、DNA をシークエンシングすることから生成されたヌクレオチド配列を含む。シークエンシングデータ 401 の各行はヌクレオチド配列である。図 4 A の例に示すように、ヌクレオチド配列は英数字の配列として表され、「A」はアデニンを表し、「C」はシトシンを表し、「G」はグアニンを表し、「T」はチミンを表し、「-」はその位置にヌクレオチドが存在しないことを表す。いくつかの実施形態は、個々のヌクレオチドまたはその欠如を表すための特定の組の英数字に限定されないことから、本明細書に記載の例示的な英数字は、例示の目的のためである。

#### 【0145】

図 4 A の実施形態では、アセンブリ 402 は、ヌクレオチド配列 401 から生成される。いくつかの実施形態では、アセンブリ 402 は、シークエンシングデータ 401 にアセンブリアルゴリズム(例えば、OLC アセンブリ)を適用することにより取得され得る。図 4 A の実施形態では、アセンブリ 402 は、ヌクレオチド配列のコンセンサスを取ることににより取得される。コンセンサスは、アセンブリ 402 内の各位置に関するヌクレオチド配列の多数決によって決定され、システムは、その位置に最大数のヌクレオチド配列によって示される生物学的ポリマーを同定する。システムは、複数のヌクレオチドの各々に関して、(1) (例えば、ヌクレオチドがその位置に存在することを示すことにより)ヌクレオチドを選出するヌクレオチド配列の数を決定し、(2) その位置において示される選出数が最も多いヌクレオチドを同定するように構成され得る。例として、強調表示された列 406 の位置に関して、(1) 4 個の配列はアデニンを示し、3 個の配列はシトシンを示し、2 個の配列はグアニンを示し、(2) アセンブリ 402 内の位置はアデニンを示す。別の例として、アセンブリ 402 の第 1 の位置に関して、全てのヌクレオチド配列はシトシンを示し、従って、アセンブリ 402 は、第 1 の位置においてシトシンを示す。

#### 【0146】

図 4 A の実施形態では、ラベル 404 は、アセンブリ 402 内の位置に対する所望の生物学的ポリマーを示し得る。いくつかの実施形態において、システムは、参照ゲノムからラベルを決定するように構成され得る。例えば、システムは、生物からの DNA サンプルをシークエンシングすることによりヌクレオチド配列を取得し、ヌクレオチド配列へのアセンブリアルゴリズムの適用によりアセンブリ 402 を取得し、生物の既知の参照ゲノムから(例えば、NCBI データベースから)ラベル 404 を取得し得る。ラベル 404 は、教師ありトレーニングのために使用され、かつ / または生成されたアセンブリの精度を決定するために使用される各位置に関する真のまたは正確な生物学的ポリマーの表示を表し得る。

#### 【0147】

図 4 B は、図 4 A に示されるデータ 400 から決定された値のアレイ 410 を示す。アレイ 410 は、アセンブリ 402 内の列 406 の位置に関する機械学習モデルへの入力生成の際の中間ステップを示す。アレイ 410 は、図 4 A のヌクレオチド配列を表す「パイルアップ」とラベル付けされた一組の行を含む。アセンブリ内の各位置に関して、システムは、複数のヌクレオチドの各々のカウントを決定する。カウントは、ヌクレオチドがアセンブリ内の位置にあることを示すヌクレオチド配列の数を示す。アレイ 410 の「パ

10

20

30

40

50

イルアップ」セクションの各エントリは、ヌクレオチドに関するカウントを保持する。例として、図 4 B における列 4 1 2 のカウントは、アデニンが 4、シトシンが 3、グアニンが 2、チミンが 0、ヌクレオチド無しが 0 である。別の例として、アレイ 4 1 0 の第 1 の列のカウントは、アデニンが 0、シトシンが 9、グアニンが 0、チミンが 0、ヌクレオチド無しが 0 である。

#### 【 0 1 4 8 】

アレイ 4 1 0 はさらに、図 4 B のアセンブリ 4 0 2 を表す、図 4 B において「アセンブリ」とラベル付けされた一組の行を含む。アセンブリ 4 0 2 内の各位置に関して、アレイ 4 1 0 は、その位置に示されたヌクレオチドから決定された列の値を含む。各位置に関して、システムは、複数のヌクレオチドの各々に参照値を割り当て得、参照値は、ヌクレオチドがアセンブリ内の位置において示されているかどうかを示す。一例として、図 4 B の 4 1 2 とラベル付けされた列において、アセンブリセクションは、( 1 ) アデニンはアセンブリ 4 0 2 内の対応する位置に示されているヌクレオチドであるため、アデニンに対する 9 の値を有し、( 2 ) 他のヌクレオチドの各々はアセンブリ 4 0 2 内の対応する位置に示されていないため、他のヌクレオチドの各々に対する 0 の値を有する。別の例として、アレイ 4 1 0 の第 1 の列において、アセンブリセクションは、( 1 ) シトシンはアセンブリ 4 0 2 内の対応する位置に示されているヌクレオチドであるため、シトシンに対する 9 の値を有し、( 2 ) 他のヌクレオチドの各々はアセンブリ 4 0 2 内の対応する位置に示されていないため、他のヌクレオチドの各々に対する 0 の値を有する。図 4 B の例に示されるように、いくつかの実施形態では、ヌクレオチドがアセンブリ位置に示されるときにアセンブリ位置においてヌクレオチドに割り当てられる参照値は、整列されたヌクレオチド配列の数に等しい(例えば、図 4 A の例では 9)。

10

20

#### 【 0 1 4 9 】

図 4 C は、図 4 B のアレイ 4 1 0 の値を使用して生成された特徴値のアレイ 4 2 0 を示す。いくつかの実施形態では、アレイ 4 2 0 は、対応する出力を得るために機械学習モデルへの入力として提供され得る。図 4 C の例では、アレイ 4 2 0 は、列 4 2 2 に対応するアセンブリ内の位置に関してモデルに提供される入力である。アレイ 4 2 0 は、列 4 2 2 に対応するターゲット位置において決定された特徴の値、およびターゲット位置の近傍における 2 4 個の位置に関して決定された特徴の値を含む。アレイ 4 2 0 は、ターゲット位置の左側にある 1 2 個の位置、およびターゲット位置の右側にある 1 2 個の位置に関する特徴の値を含む。

30

#### 【 0 1 5 0 】

アレイ 4 2 0 のパイルアップセクションにおいて、各列は、複数のヌクレオチドの各々に関するエラー値を指定する。列におけるヌクレオチドに関するエラー値は、( 1 ) ヌクレオチドが列に対応するアセンブリ 4 0 2 内の位置にあることを示すヌクレオチド配列の数と、( 2 ) アレイ 4 2 0 のアセンブリセクション内のヌクレオチドに割り当てられた参照値との間の差異を示す。例として、図 4 C の列 4 2 2 に関して、値は、( 1 ) アデニンが  $4 - 9 = -5$  であり、( 2 ) シトシンが  $3 - 0 = 3$  であり、( 3 ) グアニンが  $2 - 0 = 2$  であり、( 4 ) チミンが  $0 - 0 = 0$  であり、( 5 ) ブランクが  $0 - 0 = 0$  であるとして決定される。アレイ 4 2 0 のアセンブリセクションは、図 4 B のアレイ 4 1 0 のアセンブリセクションと同じであり得る。

40

#### 【 0 1 5 1 】

いくつかの実施形態では、アレイ 4 2 0 内のパイルアップの値は、アセンブリ 4 0 2 がある位置においてヌクレオチドを不正確に同定する尤度を示し得る。システムは、値を使用して機械学習モデルへの入力を生成する位置を選択し得る。図 4 C に示すように、パイルアップの非ゼロの値が強調表示されている。いくつかの実施形態では、システムは、ある位置におけるパイルアップ値が閾値を超えたときに、その位置に関して機械学習モデルに提供される入力を生成することを決定するように構成され得る。例えば、システムは、アデニンに関して決定された 5 の差異が 4 の閾値の差異を超えると決定することによって、列 4 2 2 に対応するアセンブリ 4 0 2 内の位置に関する入力を生成することを決定し得

50

る。閾値の差異の例は本明細書において説明されている。

【0152】

いくつかの実施形態では、アレイ420は、アセンブリ内の位置（例えば、列422に対応する位置）を更新するための機械学習モデルへの入力として提供され得る。システムは、機械学習モデルから取得した対応する出力を使用して、アセンブリ内の位置に存在するヌクレオチドを同定し、それに応じてアセンブリを更新し得る。いくつかの実施形態では、アレイ420は、モデルのトレーニングの一部として機械学習モデルに提供される複数の入力のうちの1つであり得る。システムは、機械学習モデルおよびラベル404から取得された対応する出力を使用して、機械学習モデルの1つまたは複数のパラメータへの調整を決定し得る。一例として、機械学習モデルはニューラルネットワークであり得、システムは、機械学習モデルの出力から同定されたヌクレオチドとラベルとの間の差異を使用して、ニューラルネットワークの重みに対する1つまたは複数の調整を決定し得る。

10

【0153】

図4Aの例示的な実施形態は、核酸に関連するデータを示しているが、いくつかの実施形態では、データは、タンパク質に関連し得る。例えば、配列401はアミノ酸配列であり得、アセンブリ402はタンパク質配列であり得、ラベル404はタンパク質配列中の各位置に関する参照アミノ酸であり得る。システムは、アミノ酸配列、タンパク質配列、および/またはラベルに基づいて、図4B～図4Cに示される値を決定し得る。

【0154】

図5は、本明細書に記載の技術のいくつかの実施形態による、アセンブリを更新するプロセスを示す。図5は、更新されたアセンブリ508を生成するために機械学習モデル502に提供されるアセンブリデータ500からの入力の生成を示す。アセンブリデータ500は、例えば、図4Cを参照して上記で説明したデータの形式であり得る。図示された更新のプロセスは、図1A～図1Cを参照して上記で説明されたアセンブリシステム104によって実行され得る。

20

【0155】

図5の実施形態に示されるように、システムは、更新されるべきアセンブリ内の位置504Aおよび506Aを選択する。一例として、システムは、(1)アセンブリがアセンブリ内の位置において生物学的ポリマー（例えば、ヌクレオチド、アミノ酸）を不正確に示す尤度を決定し、(2)位置504A、506Aにおける尤度が各々位置504A、506Aを選択するための閾値尤度を超えると決定することによって位置504A、506Aを選択し得る。システムが位置504A、506Aを選択すると、システムは、機械学習モデル502に提供される対応する入力を生成することを決定し得る。

30

【0156】

図5の実施形態に示されるように、システムは、位置504Aに対応する第1の入力504Bと、位置506Aに対応する第2の入力506Bとを生成する。システムは、図4A～図4Cを参照して上記のように入力504B、506Bの各々を生成し得る。例えば、システムは、(1)その位置を中心とする位置の近傍を選択し、(2)近傍の位置の各々において1つまたは複数の特徴の値を決定し、(3)特徴（単数または複数）の値を位置に関する入力として使用することによって、入力504B、506Bの各々を生成し得る。いくつかの実施形態では、システムは、特徴（単数または複数）の値をデータ構造に格納するように構成され得る。一例として、システムは、図4Cに示されるように、値を2次元アレイまたは画像内に格納し得る。

40

【0157】

図5の実施形態に示されるように、システムは、対応する出力を得るために、生成された入力504B、506Bの各々を機械学習モデル502への入力として提供する。出力504Cは、位置504Aに対して生成された入力504Bに対応し、出力506Cは、位置506Aから生成された入力506Bに対応する。いくつかの実施形態では、システムは、入力504B、506Bを機械学習モデル502に順次提供するように構成され得る。一例として、システムは、(1)入力504Bを機械学習モデル502に提供して、

50

対応する出力504Cを取得し、(2)出力504Cを取得した後、入力506Bを機械学習モデル502に提供して、対応する出力506Cを取得する。いくつかの実施形態では、システムは、入力504B、506Bを機械学習モデル502に並列に提供するように構成され得る。一例として、システムは、(1)入力504Bを機械学習モデル502に提供し、(2)入力504Bに対応する出力504Cを取得する前に、入力506Bを機械学習モデル502に提供する。

**【0158】**

図5の実施形態に示されるように、出力504C、506Cの各々は、1つまたは複数のヌクレオチドの各々がアセンブリ内の位置に存在する尤度を示す。図5の実施形態では、尤度は確率である。例として、出力504Cは、(1)4個の異なるヌクレオチドの各々に関して、ヌクレオチドが位置504Aに存在する確率と、(2)位置504Aにおいてヌクレオチドが存在しない確率(「-」文字によって表される)とを指定する。出力504Cにおいて、アデニンは0.2の確率を有し、シトシンは0.5の確率を有し、グアニンは0.1の確率を有し、チミンは0.1の確率を有し、ヌクレオチドが位置504Aにおいて存在しない確率は0.1である。別の例として、出力506Cは、(1)4個の異なるヌクレオチドの各々に関して、ヌクレオチドが位置506Aに存在する確率と、(2)位置506Aにおいてヌクレオチドが存在しない確率(「-」文字によって表される)とを指定する。この例では、アデニンは0.6の確率を有し、シトシンは0.1の確率を有し、グアニンは0.2の確率を有し、チミンは0.05の確率を有し、ヌクレオチドが位置504Aにおいて存在しない確率は0.05である。

10

20

**【0159】**

図5の実施形態に示されるように、システムは、機械学習モデル502から取得された出力を使用して、アセンブリ内の位置を更新して、更新されたアセンブリ508を取得する。いくつかの実施形態では、システムは、(1)機械学習モデルから取得した出力を使用して、位置において存在するヌクレオチドを同定し、(2)同定されたヌクレオチドを示すようにアセンブリ内の位置を更新して、更新されたアセンブリ508を取得することによってアセンブリを更新するように構成され得る。図5の例に示すように、システムは、(1)出力504Cを使用して、シトシンがその位置に存在する尤度が最も高いと判定し、(2)その位置においてシトシン「C」を示すように、更新されたアセンブリ508内の対応する位置508Aを設定することによって、初期アセンブリの位置504Aを更新する。別の例として、システムは、(1)出力506Cを使用して、アデニンがその位置に存在する尤度が最も高いと判定し、(2)アデニン「A」を示すように、更新されたアセンブリ508内の対応する位置508Bを設定することによって、初期アセンブリの位置506Aを更新する。いくつかの例では、システムは、(1)機械学習モデル502から取得した出力を使用して、ある位置において同定されたヌクレオチドが、その位置において既に示されていることを決定し、(2)更新されたアセンブリ508において位置における表示を変更せずに維持し得る。

30

**【0160】**

更新されたアセンブリ508は、初期アセンブリとは別に示されているが、いくつかの実施形態では、更新されたアセンブリ508は、初期アセンブリの更新されたバージョンであり得る。例えば、システムは、初期アセンブリをメモリに格納し、メモリ内の初期アセンブリの値を更新して、更新されたアセンブリ508を取得し得る。いくつかの実施形態では、システムは、更新されたアセンブリ508を、初期アセンブリとは別個のアセンブリとして生成し得る。例えば、システムは、初期アセンブリを第1のメモリ位置に格納し、更新されたアセンブリ508を別個のアセンブリとして第2のメモリ位置に格納し得る。

40

**【0161】**

いくつかの実施形態では、システムは、初期アセンブリ内の複数の位置において更新を順次実行するように構成され得る。一例として、システムは、(1)出力504Cを使用して、更新されたアセンブリ508内の位置508Aを更新し、(2)位置508Aにお

50

ける更新を完了した後、出力506Cを使用して、更新されたアセンブリ508内の位置508Bを更新する。いくつかの実施形態では、システムは、初期アセンブリ内の複数の位置において並列に更新を実行するように構成され得る。一例として、システムは、(1)出力504Cを使用して位置508Aの更新を開始し、(2)位置508Aにおける更新を完了する前に、出力506Cを使用して位置508Bの更新を開始する。

#### 【0162】

いくつかの実施形態では、システムは、アセンブリ内の個々の位置に関する入力を生成し、機械学習モデル502に入力を提供し、機械学習モデルからの出力を使用してアセンブリ内の複数の位置を並列に更新するプロセスを実行するように構成され得る。一例として、システムは、(1)初期アセンブリの位置504Aに関する入力の生成を開始し、(2)位置504Aにおける位置に対する更新を完了する前に、初期アセンブリの位置506Aに関する入力の生成を開始し得る。アセンブリの更新を並列化することにより、システムは、(例えば、必要な時間が短縮されることによって)アセンブリを生成するプロセスをより効率的にする。システムは、複数のプロセッサを使用し、かつ/または複数のアプリケーションスレッドを使用することにより、プロセスを並列化し得る。

10

#### 【0163】

図5の実施形態は、ゲノムアセンブリの一部を更新することを示しているが、いくつかの実施形態は、タンパク質配列またはその一部を更新するために、図示されたプロセスを実施し得る。例えば、初期アセンブリはタンパク質配列であり得る。次に、システムは、タンパク質配列内の位置に関する入力を生成して、機械学習モデル502に提供し得る。システムは、複数のアミノ酸の各々が位置において存在する尤度(例えば、確率)を示す出力を取得し得る。次に、システムは、初期タンパク質配列を更新して、更新されたタンパク質配列を取得し得る。

20

#### 【0164】

図6は、本明細書に記載の技術のいくつかの実施形態による、アセンブリを生成するための例示的な畳み込みニューラルネットワークモデル600を示す。いくつかの実施形態では、畳み込みニューラルネットワークモデル600は、図3Aを参照して上記のプロセス300を実行することによってトレーニングされ得る。いくつかの実施形態では、プロセス300から取得されたトレーニングされた畳み込みニューラルネットワークモデル600を使用して、図3Bを参照して上記のようにアセンブリを生成するためにプロセス310を実行し得る。

30

#### 【0165】

いくつかの実施形態では、モデル600は、シークエンシングデータから生成された入力、およびシークエンシングデータから生成されたアセンブリを受信するように構成される。一例として、モデル600は、図1A~図1Cを参照して上記のアセンブリシステム104によって使用される機械学習モデルであり得る。シークエンシングデータは、生物学的ポリマー配列(例えば、ヌクレオチド配列またはアミノ酸配列)を含み得る。いくつかの実施形態では、システムは、1つまたは複数の特徴の値を決定し、決定された値をモデル600への入力として提供するように構成され得る。一例として、システムは、アセンブリ内の位置の近傍における特徴の値を決定し、位置の近傍において決定された値をモデル600への入力として提供し得る。入力の例および入力を生成するための技術が本明細書で説明されている。

40

#### 【0166】

図6の例示的な実施形態では、モデル600は、モデル600に提供された入力を受信する第1の畳み込み層602を含む。第1の層602において、システムは、モデル600に提供された入力を、 $3 \times 5 \times 64$ の行列として表される64個の $3 \times 5$ フィルタにより畳み込む。例えば、システムは、出力を得るために、 $3 \times 5 \times 64$ の行列の各チャネルにより(例えば、図4Cに示されるような) $10 \times 25$ の入力マトリクスを畳み込み得る。第1の層602は、システムが畳み込みからの出力に適用する活性化関数としてReLU関数を含む。いくつかの実施形態では、第1の層602はまた、畳み込みの出力のサイ

50

ズを縮小するためのプーリング層を含み得る。

【0167】

図6の例示的な実施形態では、モデルは、第1の層602の出力を受信する第2の畳み込み層604を含む。第2の層604において、システムは、 $3 \times 5 \times 128$ の行列として表される一組の128個の $3 \times 5$ フィルタにより入力を畳み込む。システムは、第1の畳み込み層602からの出力を $3 \times 5 \times 128$ のフィルタセットにより畳み込み得る。第2の畳み込み層604は、システムが畳み込みからの出力に適用する活性化関数としてReLU関数を含む。いくつかの実施形態では、第2の層604はまた、畳み込みの出力のサイズを縮小するためのプーリング層を含み得る。次に、第2の畳み込み層604の出力は、第3の畳み込み層606に渡される。第3の層606において、システムは、 $3 \times 5 \times 256$ の行列として表される一組の256個の $3 \times 5$ フィルタにより入力を畳み込む。次に、システムは畳み込みからの出力にReLU活性化関数を適用する。いくつかの実施形態では、第3の層606はまた、畳み込みの出力のサイズを縮小するためのプーリング層を含み得る。

10

【0168】

図6の例示的な実施形態では、モデル600は、5つの完全に接続された層を有する高密度層608を含み、各々が256の入力値を受信する。システムは、第3の畳み込み層606から取得された出力を凝縮して(*condense*)、高密度層608への入力として提供し得る。高密度層608は、複数の値を出力することができ、各値は、入力がモデル600に提供された位置において個々の生物学的ポリマー(例えば、ヌクレオチドまたはアミノ酸)が存在する尤度を示す。一例として、高密度層は5個の値を出力し得、各値は、ヌクレオチド(例えば、アデニン、シトシン、グアニン、チミン、および/またはヌクレオチド無し)がその位置に存在する尤度を示す。システムは、ソフトマックス(*softmax*)関数を高密度層608の出力に適用して、合計が1になる一組の確率値を取得し得る。図6の例示的な実施形態に示されるように、システムは、ソフトマックス関数を高密度層608の出力に適用して、個々のヌクレオチドがアセンブリ内のある位置に存在する確率を示す5個の確率の出力610を取得する。出力610は、(例えば、図5を参照して上で説明したように)アセンブリを更新するために使用し得る。

20

【0169】

図7は、本明細書に記載の技術のいくつかの実施形態による技術の性能結果を示している。各プロットは、従来の手法と比較して、技術によって提供される精度の向上を示す。図7では、カナ(Canu)およびミニアスム(Miniasm)は2つの従来のアセンブリ技術である。ミニアスム(Miniasm)+レコン(Racon)は、レコン・エラー訂正を適用したミニアスムを表す。カナ(Canu)+クォーラム(Quorum)は、カナから生成されたアセンブリを修正するために本明細書で説明する技術の実施である。ミニアスム+クォーラムは、ミニアスムから生成されたアセンブリを修正するために本明細書で説明する技術の実施である。

30

【0170】

図7に示すように、ミニアスム+クォーラムは、データの各サンプルに関して、ミニアスム+レコンよりもエラー率が大幅に低くなっている。例として、 $30 \times$  PacBio(*PacBio*)データからの大腸菌の場合、ミニアスム+クォーラム(連結点で表される)の各反復のエラー率は、100エラー/100キロベース(*kilo-bases*)満であるが、ミニアスム+レコンの最小エラー率は約200エラー/100キロベースである。別の例として、 $30 \times$  ONTデータからの大腸菌の場合、ミニアスム+クォーラムの各反復のエラー率は約400エラー/100キロベースであるが、ミニアスム+レコンのエラー率は約500エラー/100キロベースである。

40

【0171】

図7に示すように、カナ+クォーラムは、カナのみの結果よりも精度が向上している。カナには従来のエラー訂正技術が組み込まれているが、本明細書で説明する技術により、アセンブリ生成の精度が向上する。例として、 $30 \times$  ONTデータからの大腸菌の場合

50

、カヌのエラー率は500エラー/100キロベースを超えるが、カヌ+クォーラムの各反復のエラー率は350エラー/100キロベース未満である。

【0172】

図7に示されるように、本明細書に記載される技術は、エラー訂正を実行するために実質的に大量の計算時間を追加することなく、アセンブリの向上された精度を提供し得る。例として、ミニアスム+クォーラムは、実質的に同じCPU時間数で、ミニアスム+レコンよりも優れた精度を実現する。別の例として、カヌ+クォーラムは、アセンブリを修正するためのCPU時間数を大幅に増加させることなく、カヌ単独よりも高い精度を実現する。

【0173】

いくつかの実施形態では、本明細書で説明されるシステムおよび技術は、1つまたは複数のコンピューティングデバイスを使用して実施され得る。しかしながら、実施形態は、特定のタイプのコンピューティングデバイスによる動作に限定されない。さらなる例として、図8は、例示的なコンピューティングデバイス800のブロック図である。コンピューティングデバイス800は、1つまたは複数のプロセッサ802および1つまたは複数の有形の非一時的なコンピュータ可読記憶媒体（例えば、メモリ804）を含み得る。メモリ804は、有形の非一時的なコンピュータ記録可能媒体に、実行時に上記の機能のいずれかを実施するコンピュータプログラム命令を格納し得る。プロセッサ802は、メモリ804に接続され、そのようなコンピュータプログラム命令を実行して、機能を実現および実行させる。

【0174】

コンピューティングデバイス800はまた、コンピューティングデバイスが他のコンピューティングデバイスと（例えば、ネットワークを介して）通信することができるネットワーク入力/出力（I/O）インタフェース806を含み、かつ、1つまたは複数のユーザI/Oインタフェース808も含み、コンピューティングデバイスは、1つまたは複数のユーザI/Oインタフェース808を介してユーザに出力を提供し、かつユーザから入力を受信する。ユーザI/Oインタフェースは、キーボード、マウス、マイクロフォン、ディスプレイデバイス（例えば、モニタまたはタッチスクリーン）、スピーカ、カメラ、および/または他の様々なタイプのI/Oデバイスなどのデバイスを含み得る。

【0175】

上述した実施形態は、多くの方法で実施することができる。例として、実施形態は、ハードウェア、ソフトウェア、又はそれらの組み合わせを用いて実施し得る。ソフトウェアで実施する場合、ソフトウェアコードは、単一のコンピューティングデバイスで提供されるか、複数のコンピューティングデバイスに分散されるかに関係なく、任意の適切なプロセッサ（例えば、マイクロプロセッサ）またはプロセッサの集合上で実行することができる。上述した機能を実行する任意の構成要素又は構成要素の集合は、上述の機能を制御する1つまたは複数のコントローラとして一般的に考えられることを理解されたい。1つまたは複数のコントローラは、専用ハードウェア、またはマイクロコードまたはソフトウェアを使用して上記の機能を実行するようにプログラムされた汎用ハードウェア（例えば、1つまたは複数のプロセッサ）など、様々な方法で実施することができる。

【0176】

この点に関して、本明細書で説明される実施形態の1つの実施は、1つまたは複数のプロセッサ上での実行時に、1つまたは複数の実施形態の上記の機能を実行するコンピュータプログラム（即ち、複数の実行可能な命令）がエンコードされた少なくとも1つのコンピュータ可読記憶媒体（例えば、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、デジタル多用途ディスク（DVD）、または他の光ディスクストレージ、磁気カセット、磁気テープ、磁気ディスクストレージまたは他の磁気ストレージデバイス、または他の有形の非一時的なコンピュータ可読記憶媒体）を含むことを理解されたい。コンピュータ可読媒体は、本明細書で説明される技術の態様を実施するために、記憶されているプログラムが任意のコンピューティングデバイスにロードでき

10

20

30

40

50

るように移送可能である。加えて、実行時に、上述した機能の任意の1つを実行するコンピュータプログラムの参照は、ホストコンピュータ上で動作するアプリケーションプログラムに限定されないことを理解されたい。むしろ、コンピュータプログラムおよびソフトウェアという用語は、本明細書では一般的な意味で使用され、1つまたは複数のプロセッサをプログラムして本明細書で説明する技術の態様を実施するために使用することができる任意のタイプのコンピュータコード（例えば、アプリケーションソフトウェア、ファームウェア、マイクロコード、または他の形式のコンピュータ命令）を指す。

【0177】

本開示の様々な特徴および態様は、単独で、2以上の任意の組み合わせにおいて、または前述の実施形態において具体的に開示されていない様々な構成で使用することができ、従って、その用途において、上述の説明または図面に示されている構成要素の詳細および構成に限定されない。例として、一実施形態で説明された態様は、別の実施形態で説明された態様と任意の方法で組み合わせることができる。

10

【0178】

「ほぼ」、「実質的に」および「約」という用語は、いくつかの実施形態では目標値の $\pm 20\%$ 以内、いくつかの実施形態では目標値の $\pm 10\%$ 以内、いくつかの実施形態では目標値の $\pm 5\%$ 以内、およびいくつかの実施形態では目標値の $\pm 2\%$ 以内を意味するために使用され得る。「ほぼ」および「約」という用語は、目標値を含むことができる。

【0179】

また、本明細書で開示されるコンセプトは、方法として具現化されてもよく、その一例が提供されている。方法の一部として実行される処理は、任意の適切な方法で順序が付けられてもよい。従って、実施形態は、例示的な実施形態では逐次的な工程として示されているが、図示されている順序とは異なる順序で工程を実施すること、及びいくつかの工程を同時に実施することも可能である。

20

【0180】

請求項の要素を修飾するために、請求項に「第1」、「第2」、「第3」等の順序を示す用語が使用されているが、これは、請求項のある1つの要素の優先度や、先行性や、順序を示すか、又はある方法を実施する時間的な順序を示すものではなく、単なる標識として同じ名称を有する（但し、通常用語を使用する）他の要素からある名前を有する別の請求項の要素を区別するために使用されている。

30

【0181】

また、本明細書で使用されている言い回しや用語は、説明を目的としたものであり、限定的なものとは見なすべきではない。本明細書における「含む」、「備える」、「有する」、「含有する」、「含む」、およびそれらの変形の使用は、その後列挙される項目およびその均等物ならびに追加の項目を包含することを意味する。

【 図 1 A 】

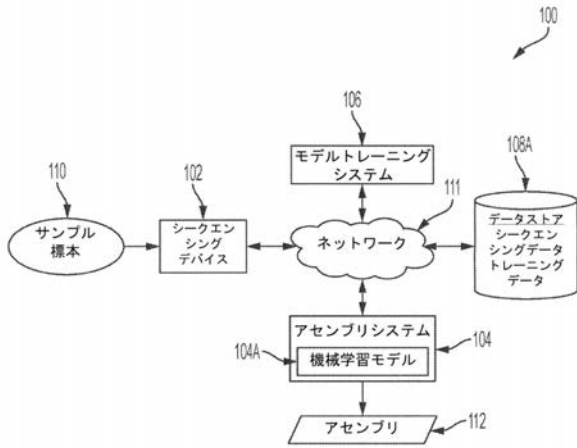


FIG. 1A

【 図 1 B 】

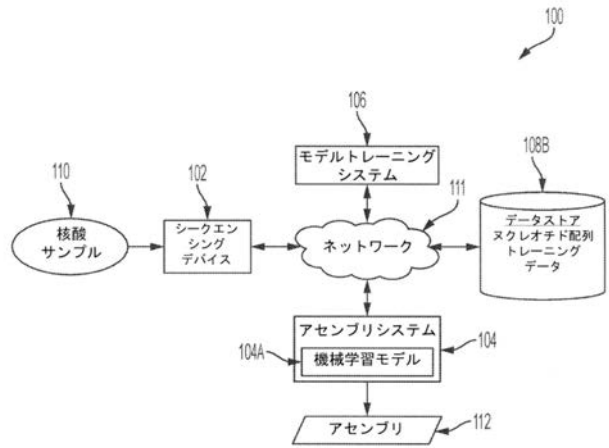


FIG. 1B

【 図 1 C 】

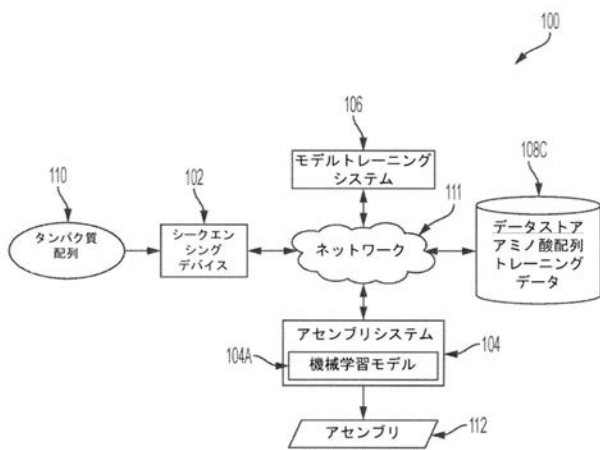


FIG. 1C

【 図 2 A 】

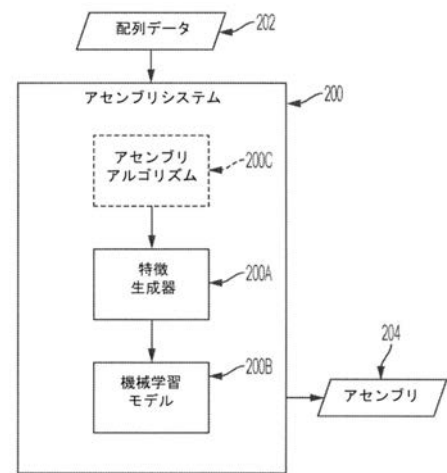


FIG. 2A

【 図 2 B 】

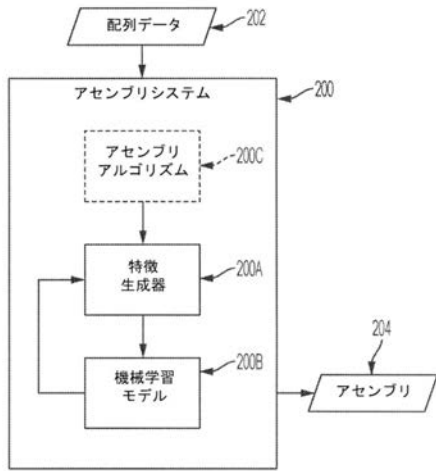


FIG. 2B

【 図 2 C 】

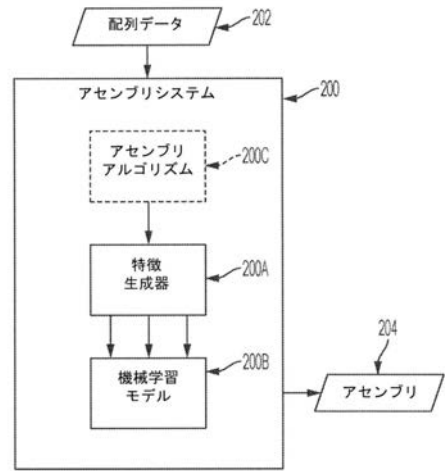


FIG. 2C

【 図 2 D 】

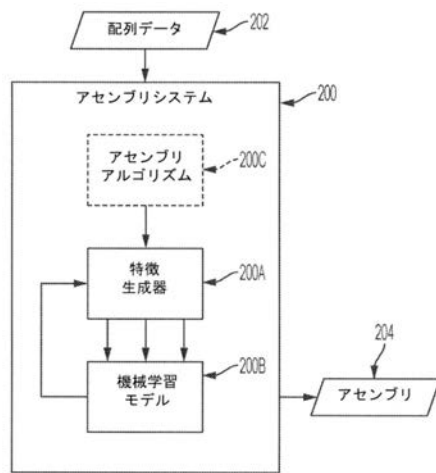


FIG. 2D

【 図 3 A 】

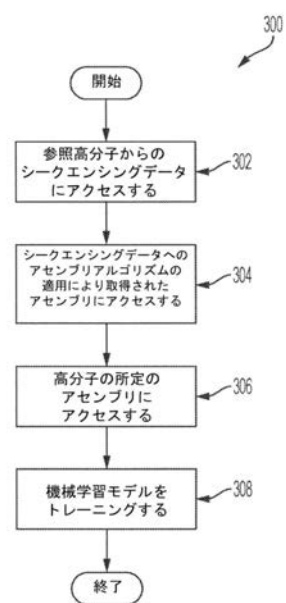


FIG. 3A



【 図 5 】

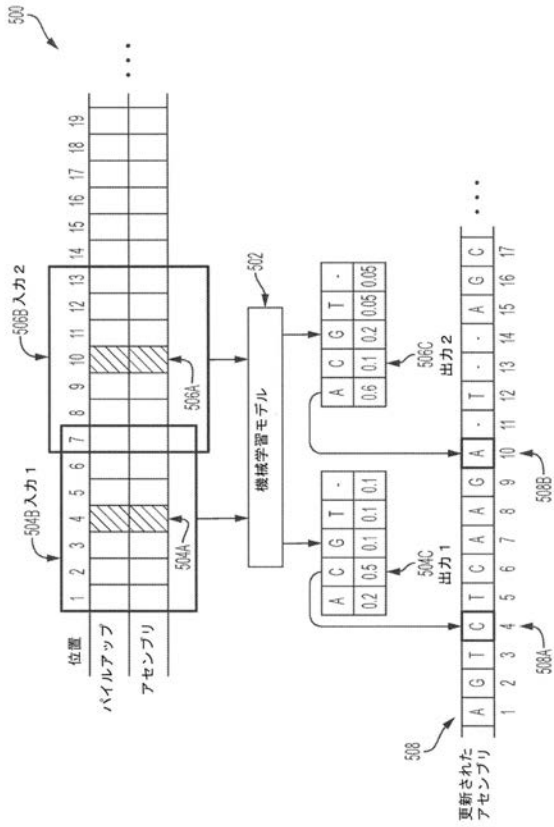


FIG. 5

【 図 6 】

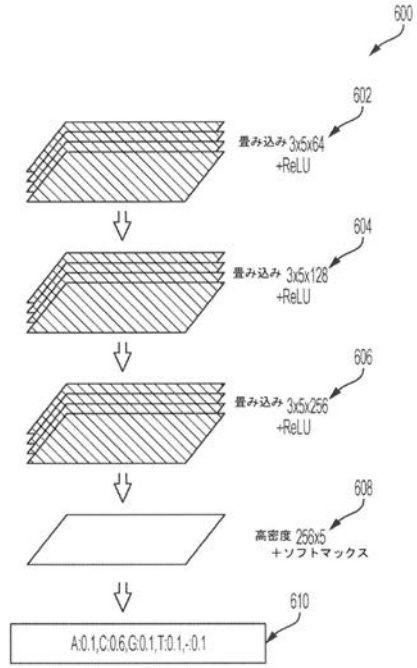


FIG. 6

【 図 7 】

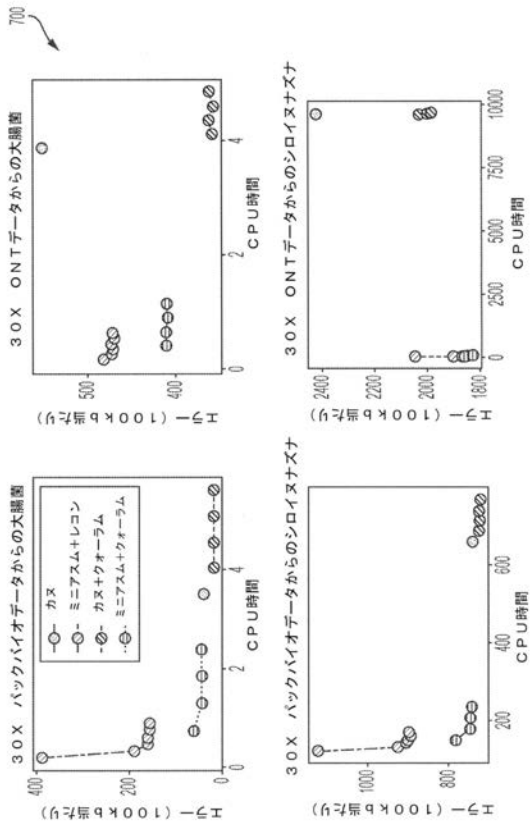


FIG. 7

【 図 8 】

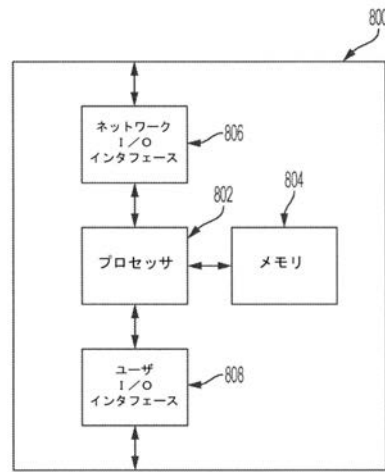


FIG. 8

## 【 国際調査報告 】

## INTERNATIONAL SEARCH REPORT

International application No PCT/US2019/032065
---

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> INV. G16B30/20 ADD. G16B40/20      G16B40/30		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) G16B		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	NICHOLAS J LOMAN ET AL: "A complete bacterial genome assembled de novo using only nanopore sequencing data", NATURE METHODS, vol. 12, no. 8, 15 June 2015 (2015-06-15), pages 733-735, XP055609123, New York ISSN: 1548-7091, DOI: 10.1038/nmeth.3444 abstract p. 736, section "Online Methods", sub-section "Assembly pipeline and software" p. 736, section "Online Methods", sub-section "Computing the consensus sequence using signal data" -/--	1-66
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents : "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 30 July 2019		Date of mailing of the international search report 09/08/2019
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer Heidrich, Alexander

2

## INTERNATIONAL SEARCH REPORT

International application No PCT/US2019/032065
---

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	-& NICHOLAS J LOMAN ET AL: "A complete bacterial genome assembled de novo using only nanopore sequencing data - Supplementary Material", NATURE METHODS, vol. 12, no. 8, 15 June 2015 (2015-06-15), pages 733-735, XP055609143, New York ISSN: 1548-7091, DOI: 10.1038/nmeth.3444 p. 7, section "1.2 Computing the Consensus Sequence" p. 8, "Consensus Algorithm": p. 9	1-66
X	----- ROBERT VASER ET AL: "Fast and accurate de novo genome assembly from long uncorrected reads", GENOME RESEARCH, vol. 27, no. 5, 18 January 2017 (2017-01-18), pages 737-746, XP055608901, US ISSN: 1088-9051, DOI: 10.1101/gr.214270.116 p. 737, section "Results" pp. 741 to 745, section "Methods"	1-66
X	----- Ryan Poplin ET AL: "Creating a universal SNP and small indel variant caller with deep neural networks", bioRxiv, 20 March 2018 (2018-03-20), XP055585250, DOI: 10.1101/092890 Retrieved from the Internet: URL:https://www.biorxiv.org/content/biorxiv/early/2018/03/20/092890.full.pdf [retrieved on 2019-05-03] abstract pages 4,6	1-66
X	----- Ruibang Luo ET AL: "Clairvoyante: a multi-task convolutional deep neural network for variant calling in Single Molecule Sequencing", bioRxiv, 26 April 2018 (2018-04-26), XP055608907, DOI: 10.1101/310458 Retrieved from the Internet: URL:https://www.biorxiv.org/content/biorxiv/early/2018/04/28/310458.full-text.pdf [retrieved on 2019-07-25] abstract page 4 - page 8 ----- -/--	1-66

**INTERNATIONAL SEARCH REPORT**

International application No  
PCT/US2019/032065

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	ALLEX C F ET AL: "Neural network input representations that produce accurate consensus sequences from DNA fragment assemblies", BIOINFORMATICS, OXFORD UNIVERSITY PRESS, SURREY, GB, vol. 15, no. 9, 1 September 1999 (1999-09-01), pages 723-728, XP002267211, ISSN: 1367-4803, DOI: 10.1093/BIOINFORMATICS/15.9.723 figures 1-8 -----	1-66

## フロントページの続き

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

(74)代理人 100152489

弁理士 中村 美樹

(72)発明者 ツァオ、ミン ドウック

アメリカ合衆国 0 2 4 7 8 マサチューセッツ州 ベルモント フレデリック ストリート 8

Fターム(参考) 4B029 AA07 BB15 FA15

4B063 QA11 QQ42 QQ52 QQ79 QS39