

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7583153号  
(P7583153)

(45)発行日 令和6年11月13日(2024.11.13)

(24)登録日 令和6年11月5日(2024.11.5)

(51)国際特許分類	F I	
G 1 6 B 40/20 (2019.01)	G 1 6 B 40/20	
C 1 2 N 15/63 (2006.01)	C 1 2 N 15/63	Z Z N A
C 1 2 N 15/12 (2006.01)	C 1 2 N 15/12	
G 1 6 B 30/00 (2019.01)	G 1 6 B 30/00	

請求項の数 15 (全57頁)

(21)出願番号	特願2023-512747(P2023-512747)	(73)特許権者	597160510
(86)(22)出願日	令和3年8月20日(2021.8.20)		リジェネロン・ファーマシューティカルズ・インコーポレイテッド
(65)公表番号	特表2023-538139(P2023-538139 A)		REGENERON PHARMACEUTICALS, INC.
(43)公表日	令和5年9月6日(2023.9.6)		アメリカ合衆国10591-6707ニューヨーク州タリータウン、オールド・ソー・ミル・リバー・ロード777番
(86)国際出願番号	PCT/US2021/046975	(74)代理人	100105957
(87)国際公開番号	WO2022/040573		弁理士 恩田 誠
(87)国際公開日	令和4年2月24日(2022.2.24)	(74)代理人	100068755
審査請求日	令和5年6月9日(2023.6.9)		弁理士 恩田 博宣
(31)優先権主張番号	63/068,654	(74)代理人	100142907
(32)優先日	令和2年8月21日(2020.8.21)		弁理士 本田 淳
(33)優先権主張国・地域又は機関	米国(US)	(74)代理人	100152489

最終頁に続く

(54)【発明の名称】 配列生成および予測のための方法およびシステム

(57)【特許請求の範囲】

【請求項1】

遺伝的データをコンピュータにより受信することであって、前記遺伝的データは、第一の複数のヌクレオチド配列を含み、前記複数のヌクレオチド配列の各ヌクレオチド配列は、関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む、受信することと、

前記第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列をコンピュータにより決定することと、

閾値を満たす前記関連発現スコアに基づいて、前記第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列をコンピュータにより決定することと、

前記第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列をコンピュータにより生成することと、

コアプロモーターとして標識された前記第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された前記第四の複数のヌクレオチド配列に基づいて、訓練データセットをコンピュータにより生成することと、

前記訓練データセットに基づいて、生成モデルをコンピュータにより訓練することと、前記生成モデルをコンピュータにより出力することと、を含む、コンピュータ実装方法。

【請求項2】

前記第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された前記第二の複数のヌクレオチド配列をコンピュータにより決定することが、

閾値を満たす前記関連発現スコアに基づいて、前記第一の複数のヌクレオチド配列から複数のTSSをコンピュータにより決定することと、

前記複数のTSSに基づいて、複数のサミットヌクレオチド塩基をコンピュータにより決定することと、

前記複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基をコンピュータにより決定することと、

各サミットヌクレオチド塩基および前記関連する複数の周辺塩基を、コアプロモーターとして標識された前記第二の複数のヌクレオチド配列としてコンピュータにより保存することと、を含む、請求項1に記載の方法。

【請求項3】

前記複数のTSSに基づいて、前記複数のサミットヌクレオチド塩基をコンピュータにより決定することは、前記複数のTSSのそれぞれについて、最強の遺伝子発現のキャップ解析(CAGE)シグナルを有するヌクレオチド塩基をコンピュータにより決定することを含む、請求項2に記載の方法。

【請求項4】

前記複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、前記関連する複数の周辺塩基をコンピュータにより決定することが、前記複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基をコンピュータにより決定することを含む、請求項2に記載の方法。

【請求項5】

5'方向の前記第一の複数のヌクレオチド塩基が49個のヌクレオチド塩基を含み、3'方向の前記第二の複数のヌクレオチド塩基が50個のヌクレオチド塩基を含む、請求項4に記載の方法。

【請求項6】

前記第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された前記第四の複数のヌクレオチド配列をコンピュータにより生成することが、

前記第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基をコンピュータにより決定することと、

コアプロモーターではないとして標識された第四の複数のヌクレオチド配列として、各関連する複数のシフト塩基をコンピュータにより保存することと、を含む、請求項1～5のいずれか一項に記載の方法。

【請求項7】

前記第三の複数のヌクレオチド配列の各ヌクレオチド配列について、前記関連する複数のシフト塩基をコンピュータにより決定することが、前記第三の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をコンピュータによりシフトさせることを含み、前記ヌクレオチド塩基の前記量は、前記第三の複数のヌクレオチド配列と類似のクロマチンランドスケープを持たせるのに十分近い距離を保つことと、その一方で、近隣にある調節要素を拾わない程度に十分に離れていることとのバランスを表す、請求項6に記載の方法。

【請求項8】

前記訓練データセットに基づいて、前記生成モデルをコンピュータにより訓練することが、

前記訓練データセットの各ヌクレオチド配列について、複数のシード配列および標的ヌクレオチド対をコンピュータにより生成することであって、各シード配列および標的ヌクレオチド対は、定義された長さを有するシード配列、および所定のヌクレオチド配列上の前記シード配列の直後に標的ヌクレオチドを含む、生成することと、

前記複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をコンピュータによりベクター化することと、

前記ベクター化されたシード配列および標的ヌクレオチド対に基づいて、前記生成モデ

10

20

30

40

50

ルをコンピュータにより訓練することと、を含む、請求項 1 ~ 7 のいずれか一項に記載の方法。

【請求項 9】

前記訓練データセットの各ヌクレオチド配列について、前記複数のシード配列および標的ヌクレオチド対をコンピュータにより生成することが、

前記関連発現スコアに基づいて、前記複数の T S S をコンピュータによりクラスタリングすることと、

T S S の各クラスターについて、四分位幅をコンピュータにより決定することと、

前記四分位幅に基づいて、各 T S S をシャープ T S S またはブロード T S S としてコンピュータにより標識することと、

シャープ T S S またはブロード T S S 標識に基づいて、前記訓練データセット中の前記ヌクレオチド配列をシャープ T S S 群またはブロード T S S 群にコンピュータにより分割することと、

前記定義された長さで、定義された工程サイズを有するスライディングウィンドウを、各ヌクレオチド配列にコンピュータにより適用することと、

前記スライディングウィンドウの各工程で、シード配列および標的ヌクレオチド対をコンピュータにより保存することと、を含む、請求項 8 に記載の方法。

【請求項 10】

前記複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をコンピュータによりベクター化することは、各ヌクレオチドをそれぞれの番号としてコンピュータによりコード化することを含む、請求項 8 に記載の方法。

【請求項 11】

前記生成モデルが、長短期メモリ ( L S T M ) リカレントニューラルネットワーク ( R N N ) を含む、請求項 1 ~ 10 のいずれか一項に記載の方法。

【請求項 12】

前記生成モデルに基づいて、ヌクレオチド配列をコンピュータにより生成することをさらに含み、前記生成モデルに基づいて、前記ヌクレオチド配列をコンピュータにより生成することが、

a ) シード配列をコンピュータにより受信することと、

b ) 前記シード配列に基づいて、次のヌクレオチドをコンピュータにより予測することと、

c ) 前記シード配列に前記次のヌクレオチドをコンピュータにより付加することと、

d ) 前記ヌクレオチド配列の所望の長さまで b ~ c をコンピュータにより繰り返すこととあって、前記ヌクレオチド配列がコアプロモーター配列である、繰り返すことと、を含む、請求項 1 ~ 11 のいずれか一項に記載の方法。

【請求項 13】

前記所望の長さは、50 ヌクレオチド ~ 100 ヌクレオチドである、請求項 12 に記載の方法。

【請求項 14】

前記コアプロモーター配列に基づいてプロモーターを操作すること、および前記プロモーターを核酸構築物に挿入することをさらに含む、請求項 12 に記載の方法。

【請求項 15】

予測モデルに、前記ヌクレオチド配列を提供することと、

前記予測モデルに基づいて、前記ヌクレオチド配列がコアプロモーターであることをコンピュータにより決定することと、を含む、請求項 12 に記載の方法。

【発明の詳細な説明】

【背景技術】

【0001】

関連出願の相互参照

本出願は、2020年8月21日に出願された米国仮特許出願第 63 / 068 , 654

10

20

30

40

50

号の出願日における優先権の利益を主張するものである。これ以前の出願の内容は、参照によりその全体が本明細書に組み込まれる。

【0002】

配列表への参照

2021年8月20日に提出された配列表は、2021年8月20日に「37595\_\_0033P1\_\_Sequence\_\_Listing.txt」というファイル名のテキストファイル（サイズ2,805バイト）として作成されたものであり、連邦規則法典第37巻特許法第1.52条(e)(5)に従い、本明細書において参照により援用されている。

【0003】

アデノ随伴ウイルス(AAV)は、遺伝子治療において導入遺伝子の送達のためのゴールドスタンダードである。免疫原性が低く感染性が強いなど、多くの利点を提供する一方で、一つの制限は、その厳密なDNAパッケージング能力である。多くの治療法がすでにこの限界に近づいている。組換えAAVベクターにコードされた他の特徴とともに、これは制御配列のためのスペースをほとんど残さない。一般的に使用されるウイルス性および内因性哺乳類プロモーターは両方ともこれらの制限を超えており、AAVを介した大きな導入遺伝子の送達には使用できない。したがって、短く効率的な制御配列に対する強いニーズがある。

【発明の概要】

【0004】

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む、遺伝的データを受信することと、閾値を満たす関連発現スコアに基づいて、第一の複数のヌクレオチド配列から複数のTSSを決定することと、複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基に対して、関連する複数の周辺塩基を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することと、第二の複数のヌクレオチド配列の各ヌクレオチド配列に対して、関連する複数のシフト塩基を決定することと、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、予測モデルに対する複数の特徴を決定することと、訓練データセットの第一の部分に基づいて、複数の特徴に従って予測モデルを訓練することと、訓練データセットの第二の部分に基づいて、予測モデルを試験することと、および試験に基づいて、予測モデルを出力することと、を含む方法が開示される。

【0005】

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、第二の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、予測モデルに対する複数の特徴を決定することと、訓練データセットの第一の部分に基づいて、複数の特徴に従って予測モデルを訓練することと、訓練データセットの第二の部分に基づいて、予測モデルを試験することと、および試験に基づいて、予測モデルを出力することと、を含む方法も開示される。

## 【 0 0 0 6 】

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点（ＴＳＳ）を含む、遺伝的データを受信することと、遺伝的データを正規化することと、関連発現スコアに基づいて、ＴＳＳをクラスター化することと、ＴＳＳの各クラスターについて、四分位幅を決定することと、四分位幅に基づいて、各ＴＳＳをシャープＴＳＳまたはブロードＴＳＳに標識することと、複数のＴＳＳに基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基に対して、関連する複数の周辺塩基を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することと、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列の各ヌクレオチド配列に対して、関連する複数のシフト塩基を決定することと、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットにおける各ヌクレオチド配列に対して、複数のシード配列および標的ヌクレオチド対を生成することと、複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクトル化することと、ベクトル化されたシード配列および標的ヌクレオチド対に基づいて、生成モデルを訓練することと、および生成モデルを出力することと、を含む方法も開示される。

10

20

## 【 0 0 0 7 】

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点（ＴＳＳ）を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、生成モデルを訓練することと、および生成モデルを出力することと、を含む方法も開示される。

30

## 【 0 0 0 8 】

また、ヌクレオチド配列を受信することと、訓練された予測モデルにヌクレオチド配列を提供することと、および予測モデルに基づいて、ヌクレオチド配列がコアプロモーターであることを決定することと、を含む方法が開示される。

## 【 0 0 0 9 】

また、（ a ）ヌクレオチド配列と配列長を受信すること、（ b ）訓練された生成モデルに、ヌクレオチド配列を提供すること、（ c ）生成モデルに基づいて、ヌクレオチド配列に関連付けられた次のヌクレオチドを決定すること、（ d ）ヌクレオチド配列に次のヌクレオチドを付与すること、（ e ）ヌクレオチド配列の長さが配列長に等しくなるまで b ~ d を繰り返すこと、および（ f ）ヌクレオチド配列をコアプロモーター配列として出力すること、を含む方法も開示される。

40

## 【 0 0 1 0 】

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点（ＴＳＳ）を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、閾値を

50

満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、生成モデルを訓練することと、を含む方法も開示される。

【 0 0 1 1 】

開示される方法のいずれかを実施するよう構成された装置が開示される。

装置が開示される方法のいずれかを実施するよう構成された、プロセッサが実行可能な指示実施形態を有する、コンピュータ可読媒体が開示される。

10

【 0 0 1 2 】

開示される方法および組成物のさらなる利点は、一部が、以下の記載において記載されるか、一部が、記載から理解されるか、または開示される方法および組成物の実施によって学んでもよい。開示される方法および組成物の利点は、添付の特許請求の範囲において特に指摘されている要素および組み合わせによって実現され、達成されるであろう。前述の一般的な説明および以下の詳細な説明は両方とも、請求される本発明の、あくまで例示的かつ説明的なものであって、限定的なものではないことを理解されたい。

【図面の簡単な説明】

【 0 0 1 3 】

20

本明細書において援用され、かつ本明細書の一部を成す添付の図面は、開示される方法および組成物の一部の実施形態を例証し、説明と共に、開示される方法および組成物の原理を説明する役割を果たすものである。

【図 1】図 1 は、例示的な操作環境を示す。

【図 2】図 2 は、例示的な方法を示す。

【図 3】図 3 は、例示的な方法を示す。

【図 4 A】図 4 A は、例示的な R N N ブロックのコンパクトな表記を示す。

【図 4 B】図 4 B は、R N N ブロックの拡張表記を示す。

【図 5】図 5 は、例示的な L S T M - R N N ブロックを示す。

【図 6】図 6 は、例示的な方法を示す。

30

【図 7】図 7 は、例示的な方法を示す。

【図 8】図 8 は、方法の一例を示す。8 0 2 は、配列番号 1 を示し、8 0 6 は、配列番号 2 を示し、8 1 0 は、配列番号 3 を示し、8 1 4 は、配列番号 4 を示し、8 1 8 は、配列番号 5 を示す。

【図 9】図 9 は、例示的な方法を示す。

【図 1 0】図 1 0 は、予測モデルの例示的な特徴を示す。

【図 1 1】図 1 1 は、例示的な方法を示す。

【図 1 2】図 1 2 は、例示的な方法を示す。

【図 1 3】図 1 3 は、プロモーターアッセイの一例を示す。

【図 1 4】図 1 4 は、コアプロモーターを制御する生成コアプロモーターの性能の比較を示す。

40

【図 1 5】図 1 5 は、例示的な操作環境を示す。

【図 1 6 - 1】図 1 6 は、例示的な方法を示す。

【図 1 6 - 2】同上。

【図 1 7】図 1 7 は、例示的な方法を示す。

【図 1 8 - 1】図 1 8 は、例示的な方法を示す。

【図 1 8 - 2】同上。

【図 1 9】図 1 9 は、例示的な方法を示す。

【図 2 0】図 2 0 は、例示的な方法を示す。

【図 2 1】図 2 1 は、例示的な方法を示す。

50

## 【発明を実施するための形態】

## 【0014】

下記の特定の実施形態およびそれに含まれる実施例についての発明を実施するための形態、ならびに図面およびその前後の説明を参照することによって、開示される方法および組成物についての理解を容易にすることができる。

## 【0015】

## A. 用語の定義

当然のことながら、本開示の方法および組成物は、記載されている特定の方法論、プロトコルおよび試薬に限定されるものではない。理由はこれらが、変更される可能性があるからである。本明細書中に使用されている用語は、あくまで特定の実施形態を説明すること

10

## 【0016】

本明細書および添付の特許請求の範囲で使用される場合、単数形「a」、「an」、および「the」は、文脈が明白に指示しない限り、複数の参照を含むことに留意されたい。したがって、例えば、「ある配列」への言及は、複数の配列を含み、「その配列」への言及は、一つまたは複数の配列および当業者に公知のその均等物などへの言及である。

## 【0017】

本明細書で使用される場合、用語「配列決定」または「シーケンサー」は、生体分子、例えば、DNAまたはRNAなどの核酸の配列を決定するために使用される多数の技術のいずれかを指す。例示的な配列決定方法としては、標的配列決定、単一分子のリアルタイム配列決定、エクソン配列決定、電子顕微鏡ベースの配列決定、パネル配列決定、トランジスタ介在性配列決定、直接配列決定、ランダムショットガン配列決定、サンガーシデオキシ末端配列決定、全ゲノム配列決定、ハイブリダイゼーションによる配列決定、パイロシークエンシング、二本鎖配列決定、サイクルシーケンシング、単一塩基伸長配列決定、固相配列決定、ハイスループット配列決定、超平行シグネチャシーケンシング、エマルションPCR、より低い変性温度PCR (COLD-PCR) での共増幅、マルチプレックスPCR、可逆的染料ターミネーターによる配列決定、対末端配列決定、短期配列決定、エキソヌクレアーゼ配列決定、ライゲーションによる配列決定、ショートリードシーケンシング、一分子配列決定、合成による配列決定、リアルタイムシーケンシング、逆ターミネーター配列決定、ナノポア配列決定、454配列決定、Solexa Genome Analyzer配列決定、SOLiD (商標) 配列決定、MS-PET配列決定、およびその組み合わせが挙げられるが、これらに限定されない。一部の実施形態では、配列決定は、例えば、IlluminaまたはApplied Biosystemsから市販されている遺伝子アナライザーなどの遺伝子アナライザーによって行うことができる。

20

30

## 【0018】

「ポリヌクレオチド」、「核酸」、「核酸分子」、または「オリゴヌクレオチド」は、ヌクレオシド間結合によって結合されたヌクレオシド (デオキシリボヌクレオシド、リボヌクレオシド、もしくはそのアナログを含む) の直鎖ポリマーを指す。典型的には、ポリヌクレオチドは、少なくとも三つのヌクレオシドを含む。オリゴヌクレオチドは、通常、数個の単量体単位、例えば、3~4個から数百個の単量体単位までのサイズ範囲に及ぶ。ポリヌクレオチドが、「ATGCCCTG」などの文字の配列で表される場合、ヌクレオチドは、左から右に5' 3'の順であり、別段示されない限り、「A」は、アデノシンを示し、「C」は、シトシンを示し、「G」は、グアノシンを示し、「T」は、チミジンを示すことは、理解されるだろう。文字A、C、G、およびTは、当該技術分野で標準的なように、塩基自体、ヌクレオシド、または塩基を含むヌクレオチドを指すように使用され得る。

40

## 【0019】

用語「DNA (デオキシリボ核酸)」は、それぞれが、四つの核酸塩基、すなわち、アデニン (A)、チミン (T)、シトシン (C)、およびグアニン (G) のうちの一つを含

50

む、デオキシリボヌクレオシドを含むヌクレオチドの鎖を指す。用語「RNA（リボ核酸）」は、それぞれが、四つの核酸塩基、すなわち、A、ウラシル（U）、G、およびCのうちの一つを含む、四つのタイプのリボヌクレオシドを含むヌクレオチドの鎖を指す。ヌクレオチドの特定の対は、相補的な様式で互いに特異的に結合する（相補的塩基対と呼ばれる）。DNAでは、アデニン（A）は、チミン（T）と対形成し、シトシン（C）は、グアニン（G）と対形成する。RNAでは、アデニン（A）は、ウラシル（U）と対形成し、シトシン（C）は、グアニン（G）と対形成する。第一の核酸鎖が、第一の鎖のヌクレオチドに相補的であるヌクレオチドからなる第二の核酸鎖に結合するとき、この二つの鎖は、結合して、二本鎖を形成する。本明細書で使用される場合、「核酸配列決定データ」、「核酸配列決定情報」、「核酸配列」、「ヌクレオチド配列」、「ゲノム配列」、「遺伝子配列」または「フラグメント配列」もしくは「核酸配列決定読み取り」は、DNAまたはRNAなどの核酸の分子（例えば、全ゲノム、全トランスクリプトーム、エキソーム、オリゴヌクレオチド、ポリヌクレオチド、またはフラグメント）におけるヌクレオチド塩基の順序（例えば、アデニン、グアニン、シトシン、およびチミンまたはウラシル）を示す任意の情報またはデータを示す。本教示は、キャピラリー電気泳動、マイクロアレイ、ライゲーションベースのシステム、ポリメラーゼベースのシステム、ハイブリダイゼーションベースのシステム、直接的または間接的ヌクレオチド識別システム、パイロシーケンシング、イオンベースもしくはpHベースの検出システム、および電子署名ベースのシステムを含むが、これらに限定されない、すべての利用可能な様々な技術、プラットフォームまたはテクノロジーを使用して得られる配列情報を企図すると、理解されるべきである。

10

20

#### 【0020】

「ベクター」は、プラスミド、ファージ、ウイルス構築物、またはコスミドなどのレプリコンであり、これに別のDNAセグメントが付着する場合がある。ベクターを使用して、DNAセグメントを細胞内に形質導入し、発現させる。

#### 【0021】

「プロモーター」または「プロモーター配列」は、細胞内でRNAポリメラーゼを結合し、メッセンジャーRNA、リボソームRNA、核小体RNAの核内低分子、または任意のRNAポリメラーゼI、II、またはIIIの任意のクラスによって転写された任意の種類のRNAなどのポリヌクレオチドまたはポリペプチドコード配列の転写を開始することができるDNA制御領域である。

30

#### 【0022】

「任意選択的な」または「任意選択的に」は、後述されている事象、状況または材料が起こる場合もあれば起こらない場合もあるか、存在する場合もあれば存在しない場合もあることを意味すると共に、この記載には、前述の事象、状況または材料が起こる場合の例および起こらない場合の例、または存在する場合の例および存在しない場合が包含されることを意味する。

#### 【0023】

この明細書の記載および特許請求の範囲を通じて、語「含む（comprise）」およびこの語の変形、例えば「含む（comprising）」および「含む（comprises）」などは、「～を含むがこれに限定されない」を意味し、例えば、他の追加のもの、コンポーネント、整数、または工程を除外することを意図するものではない。特に、一つまたは複数の工程または動作を含むものとして記載される方法では、それぞれの工程が、列挙されているものを含むこと（その工程が、「からなる」などの限定する用語を含まない限り）が具体的に企図されており、それは、それぞれの工程が、例えば、工程に挙げられていない他の追加のもの、コンポーネントまたは工程を排除することが意図されていないことを意味している。

40

#### 【0024】

「例示的な」は、「の一例」を意味し、好ましい構成または理想的な構成の表示を伝達することを意図するものではない。「など」は、限定的な意味で使用されるものではなく

50



、説明を目的に使用される。

【0025】

本明細書では、範囲は、「約」一つの特定の値から、かつ／または「約」別の特定の値までとして表現される場合がある。こうした範囲が表されるとき、具体的に企図され、開示されることが考慮される範囲は、文脈が別途具体的に示さない限り、一つの特定の値からおよび／または他の特定の値の範囲である。同様に、値が近似値として表現されている場合には、先行する「約」を使用することにより、特定の値が別の実施形態を形成することが理解されるであろうし、具体的には、文脈が別途具体的に示さない限り、開示されることが考慮されるべき実施形態が企図される。これらの範囲の各々の終点は、文脈が別途具体的に示さない限り、他の終点と関連して、かつ他の終点とは独立して有意であることがさらに理解されるであろう。最後に、明示的に開示された範囲内に含まれる個々の値および値のサブレンジの全ても、具体的に企図されており、文脈が別段示さない限り、開示されているとみなされるべきであることが理解されるべきである。前述は、特定の事例において、これらの実施形態の一部またはすべてが明示的に開示されているか否かにかかわらず、適用される。

10

【0026】

B．制御配列を設計するためのアプローチ

図1は、AAVベクター内にパッケージングされたAAVベクターおよびDNAの概略図を示す。AAV遺伝子治療ベクターのパッケージ化されたDNAには、AAVゲノムの3'末端の一つおよび5'末端の一つの、二つの逆位末端反復(ITS)配列が必要であることが判明した。AAVの二つのITSは合計で約0.2~0.3kbであるため、これら二つのITSの間に導入できる外来性DNA(対象遺伝子を含む)は4.4kbよりも小さくしなければならない。一例として、ITSは、2x145bpの長さであってもよい(左ITSと右ITSは同一である)。2つのITS間の外来性DNAの長さが許容される最大値(4~4.4kb)に近い場合、パッケージング効率は大幅に低下する。外来性DNAは、図1に示すように、プロモーター(エンハンサー要素を含むまたは含まない)、目的の導入遺伝子/遺伝子、ポリAを含むことができるが、これらに限定されない。含まれる外来性DNA要素が多いほど、目的の遺伝子は小さくなり得る。したがって、目的の遺伝子のサイズを大きくすることができる一つの選択肢は、プロモーターなどの制御配列を設計するための機械学習アプローチにおいて本明細書に記載されるプロモーターサイズを小さくすることである。

20

30

【0027】

本明細書では、制御配列を設計するための機械学習アプローチを記載する。記載された方法は、前処理のためのデータ方法(訓練データの生成)、および予測モデルおよび生成モデルを生成するための方法に分離され得る。したがって、図2に示すように、210でプロモーター配列データセットを決定することを含む、方法200が本明細書に記載される。プロモーター配列データセットの一部、全部、またはバリエーションを使用して、220で生成モデルのための訓練データセットを生成してもよい。生成モデルは、プロモーター配列データセットに従って訓練されることに基づいて、プロモーター配列を生成するように構成されてもよい。プロモーター配列データセットの一部、全部、またはバリエーションを使用して、220で予測モデルのための訓練データセットを生成してもよい。生成モデルのための訓練データセットを使用して、240で生成モデルを訓練してもよい。予測モデルのための訓練データセットを使用して、240で予測モデルを訓練してもよい。予測モデルは、生成モデルの品質管理機構としての役目を果たし得る。生成モデルを使用して、260でコアプロモーター配列を生成してもよい。予測モデルを使用して、270でコアプロモーター配列を、コアプロモーターとして、またはコアプロモーターではないとして分類してもよい。したがって、生成されたコアプロモーター配列を実験設定で試験する前に、予測モデルを使用して、生成モデルの設定をベンチマークし、生成された配列が、内因性配列で訓練されたモデルに基づいて、コアプロモーター活性に対して陽性であると予測されるかどうかを試験してもよい。予測モデルと生成モデルとの間のデータ漏洩を避け

40

50

るために、すべての配列は、いかなる重複も含まないように相互参照されてもよい。

#### 【0028】

##### C. 訓練データの生成方法

機械学習では、訓練データセットは、さらなる適用および利用のためのベースラインとしての役割を果たす、初期データセットであってもよい。一部の実施形態では、訓練データセットは、標識される。一部の実施形態では、訓練データセットは、標識されない。一つまたは複数の機械学習技術を使用して、訓練データセットを分析し、特徴（説明変数または独立変数と呼んでもよい）と結果（必要に応じて、目的変数または従属変数と呼んでもよい）との間の関係を一般化するモデルを作成してもよい。

#### 【0029】

したがって、プロモーター配列データから訓練データセットを作成するための方法が説明される。訓練データセットは、訓練されるモデル（例えば、生成モデル対予測モデル）に基づいて、異なる方法に従って作成されてもよい。ヒトゲノムの候補コアプロモーターのリストを含む候補コアプロモーターデータ（プロモーター配列データ）を決定してもよい。一実施形態では、候補コアプロモーターデータは、公開されているソースからダウンロードしてもよい。候補コアプロモーターデータは、転写開始点（TSS）プロファイリングデータとして、例えば、FANTOM5からダウンロードしてもよい。一実施形態では、プロモーター配列データは、遺伝子発現のキャップ解析（CAGE）データを含んでもよい。CAGEは、転写開始点およびそのプロモーターをマッピングするための技術である。CAGEは、ゲノム全体の転写開始検出を可能にし、プロモーター使用解析を含む、組織/細胞/条件特異的転写開始点（TSS）の同時同定によるハイスループット遺伝子発現プロファイリングをもたらす。CAGEは、mRNAの5'末端から最初の20ヌクレオチドに由来するDNAタグのコンカテマーの調製および配列決定に基づいており、これは分析されたサンプル中のmRNAの元の濃度（RNA頻度）を反映する。CAGEデータでは、生物学的状態（サンプル）のパネル全体のTSSピークは、TSSピークの各々が、隣接し、関連するTSSを含む、DPI（分解ベースのピーク識別）によって識別されてもよい。TSSピークは、プロモーターを定義するためのアンカーとして、およびプロモーターレベルの発現分析の単位として使用され得る。TSSサミットは、所与のコアプロモーター内で最も強いシグナルを有するヌクレオチド塩基の座標を指し得る。

#### 【0030】

プロモーター配列データは、生成モデル用の訓練データセットおよび予測モデル用の訓練データセットを生成するために使用され得る。

##### 1. 生成モデルのための訓練データの生成

一実施形態では、プロモーター配列データは、必要に応じてフィルタリングされてもよい。例えば、プロモーター配列データは、種、成人/小児の状態、器官/組織の関連性などによって配列のみを保持するようにフィルタリングされてもよい。例えば、プロモーター配列データは、ヒト、成人、および/または肝臓データに関連付けられた配列のみを保持するようにフィルタリングされてもよい。プロモーター配列データは、正規化されてもよい。例えば、べき乗分布は正規化に使用され得る。一実施形態では、複数のライブラリ（例えば、成人肝臓から一つと成人腎臓から一つ）を使用する場合、正規化が用いられてもよい。ライブラリ内で、配列リードまたはタグの数が、TSSの相対的な強度を示している場合がある。しかしながら、シーケンスされたリードの総数は各実験で異なる可能性があるため、リードまたはタグの絶対数は、異なるシーケンシング実験間で比較できない。それらを比較できるようにするために、シーケンスのリードまたはタグのカウントが正規化される。この正規化は、例えば、リード数を倍率で割る（例えば、リードの総数を100万で割る）ことによって行われてもよい。別の実施例では、正規化は、べき乗分布に基づいてもよい。CAGEタグの数以下の数（ $\leq$ ）によってサポートされるTSSの数は、べき乗分布によって近似され得る逆累積分布に従う。例えば、典型的なCAGEライブラリでは、10個のCAGEタグでサポートされる1,000個のTSSと、1,000個のCAGEタグでサポートされる10個のTSSがある。各ライブラリは、そ

10

20

30

40

50

の実験的に決定された分布を、分析中のライブラリとほぼ同様の仮想参照分布に当てはめることによって正規化することができる。例えば、肝臓と腎臓の両方のサンプルの上位 1,000 個の TSS (CAGE タグ数による) のみが分析される場合、正規化が有用であり得る。肝臓サンプルがより深く配列決定されている (総リード数が多い) 場合、ランキングは腎臓で活性な TSS によって支配されるか、またはその逆であり得る。適切に正規化した後、腎臓と肝臓の両方で上位約 500 のピークが、組み合わせたサンプルの上位 1,000 の TSS を構成するはずである。

#### 【0031】

次いで、TSS をコアプロモーターと区別してもよい。コアプロモーターは、近くの TSS のクラスターとみなすことができる。これらの TSS は、わずかに異なる 5' 末端を有する機能的に等価な mRNA を生じさせる。個別のコアプロモーターは、例えば 100 bp のウィンドウ内に複数の小さな TSS の「ピーク」があるようなブロード転写開始パターン、または例えば、10 bp のウィンドウ内に単一の高い TSS ピークといくつかの非常に小さな TSS ピークがあるようなシャープ転写開始パターンを有することができる (これは一例である)。コアプロモーター、その後クラスを決定するために、どの TSS が共通クラスターに属するかを決定してもよい。したがって、プロモーター配列データ中の TSS は、クラスター化されてもよく、結果として生じるコアプロモーター (または TSS クラスター) の四分位幅が決定されてもよい (例えば、下分位点 = 0.1、上分位点 = 0.9)。一実施形態では、二つの独立した TSS が互いに 20 bp 以下 ( $\leq$ ) 離れている場合、二つの独立した TSS が同一のクラスター (したがってコアプロモーター) に属する、距離ベースのクラスタリングを使用してもよい。さらに、TSS は、クラスタリングに含まれる最小限のリード数によってサポートされるように要求されてもよい。

#### 【0032】

クラスタリング後、TSS クラスター (= コアプロモーター) 幅を決定することができる。そのためには、クラスター全体に沿って移動し、TSS タグの総数を累積和として数えることができる。その和が下分位点 (例えば、合計の 0.1 または 10%) に当たる位置は、コアプロモーター開始として定義されてもよく、その和が、上分位点 (例えば、合計の 0.9 または 90%) に当たる位置は、コアプロモーター終了として定義されてもよい。これらの二つの位置の間の幅は、四分位幅と呼んでもよい。

#### 【0033】

TSS は、四分位幅に基づいてプロモータークラスにピン化してもよい。プロモータークラスは、例えば、狭いゲノム領域内に転写が起こるシャープ型プロモーター、およびより大きなゲノム領域に TSS が分散されるブロード型プロモーターがあり得る。コアプロモーターは、四分位幅によってランク付けされてもよく、下半分はシャープ (幅が小さい)、上半分はブロード (幅が広い) として標識されてもよい。シャープ型プロモーターおよびブロード型プロモーターは、それぞれ TATA ボックスおよび CpG アイランドと関連している可能性が高い。次いで、候補コアプロモーター配列は、各 TSS について、TSS サミットを 5' 方向に数塩基、および 3' 方向に数塩基拡張することによって決定され得る。例えば、長さ 100 bp の候補コアプロモーター配列を作成するために、5' 方向に 49 bp、および 3' 方向に 50 bp のヌクレオチドを決定することができる。候補コアプロモーター配列は、CAGE シグナルに従ってフィルタリングされてもよい。閾値未満の CAGE シグナルを有する候補コアプロモーター配列は除外されてもよく、得られたコアプロモーター配列はコアプロモーターとして標識されてもよい。閾値は、例えば、10 を超える正規化されたカウントであってもよい。別の実施例では、閾値は、約 5 ~ 約 15 であり得る。カウント分布は、非常に長い「尾」を持つ分布であってもよく、ほとんどのコアプロモーターは、少数の CAGE タグのみによってサポートされているのに対し、多くのタグによってサポートされるコアプロモーターはごくわずかであることを意味する。10 を超えるカットオフを選択すると、最強のコアプロモーターのみが検討されるようになる。あるいは、例えば、約 1,000 ピーク ~ 約 3,000 ピークまでの最高数のピークを使用することができる。閾値として選択されたピークの数、総数と強度との間のトレ

ードオフを表す。コアプロモーターを多く選択するほど、弱いコアプロモーターを含むことによってシグナルが「希釈」される。このカットオフは、これら二つの相反する力の間のバランスである。他の閾値の例としては、限定されるものではないが、5（約上位5000）、10（約上位3000）、25（約上位1500）、50（約上位1000）などの絶対数が挙げられる。

#### 【0034】

一実施形態では、制御配列のセットは、コアプロモーター配列のセットを、5'方向または3'方向に塩基数だけシフトさせることによって生成されてもよい。例えば、候補コアプロモーター配列を5'方向に50,000bpシフトさせることによって。シフトした塩基数は、類似のクロマチンランドスケープを有するのに十分近い距離を保ちつつ、隣接する調節要素を選ばないほど十分に離れていることのバランスを表す。5'方向へのシフトは、遺伝子本体（3'方向に延び、哺乳類ゲノムでは50kb超の長さであることが多い）へのシフトを防止する。別の実施形態では、制御配列のセットは、ゲノム全体からランダム配列を選択することによって生成されてもよい。制御配列は、任意のCAGEピークと重複する任意の制御配列を除去するためにフィルタリングされてもよく、制御配列は、コアプロモーターではないとして標識されてもよい。

#### 【0035】

一実施形態では、遺伝的データを受信することを含む生成モデル用の訓練データセットを生成するための方法が記述されている。遺伝的データは、第一の複数のヌクレオチド配列を含むことができる。第一の複数のヌクレオチド配列は、プロモーター配列を含み得る。複数のヌクレオチド配列の各ヌクレオチド配列は、関連発現スコアを有する少なくとも一つの転写開始点（TSS）を含み得る。関連発現スコアは、CAGEピークを含み得る。遺伝的データは、プロモーター配列データを含んでもよい。遺伝的データは正規化されてもよい。べき乗法の適用を含む、当技術分野で公知の任意の正規化技術を使用してもよい。TSSは、関連発現スコアに基づいてクラスター化されてもよく、TSSの各クラスターについて、四分位幅が決定されてもよい。四分位幅は、各TSSをシャープTSSまたはブロードTSSとして標識するために使用され得る。複数のサミットヌクレオチド塩基が、TSSで決定されてもよい。サミットヌクレオチド塩基を決定することは、最も強いCAGEシグナルを有するヌクレオチド塩基を決定することを含み得る。各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定してもよい。関連する複数の周辺塩基を決定することは、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定し、それによって候補コアプロモーター配列を形成することを含み得る。5'方向の第一の複数のヌクレオチド塩基は、49個のヌクレオチド塩基を含むことができ、3'方向の第二の複数のヌクレオチド塩基は、50個のヌクレオチド塩基を含むことができる。各サミットヌクレオチド塩基およびその関連する複数の周辺塩基は、コアプロモーターとして標識された第二の複数のヌクレオチド配列（候補コアプロモーター配列）として保存されてもよい。第三の複数のヌクレオチド配列は、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から決定され得る。閾値は、例えば、10を超える正規化されたカウントであってもよい。別の実施例では、閾値は、約5～約15であり得る。カウント分布は、非常に長い「尾」を持つ分布であってもよく、ほとんどのコアプロモーターは、少数のCAGEタグのみによってサポートされているのに対し、多くのタグによってサポートされるコアプロモーターはごくわずかであることを意味する。10を超えるカットオフを選択すると、最強のコアプロモーターのみが検討されるようになる。あるいは、例えば、約1,000ピーク～約3,000ピークまでの最高数のピークを使用することができる。閾値として選択されたピークの数、総数と強度との間のトレードオフを表す。コアプロモーターを多く選択するほど、弱いコアプロモーターを含むことによってシグナルが「希釈」される。このカットオフは、これら二つの相反する力の間のバランスである。他の閾値の例としては、限定されるものではないが、5（約上位5000）、10（約上位3000）、25（約上位1500）、50（約上位1000）

10

20

30

40

50

などの絶対数が挙げられる。第三の複数のヌクレオチド配列は、コアプロモーター配列のセットを含んでもよい。コアプロモーター配列のセットは、ヒトゲノムアセンブリ (hg19) 中にNsを含有する任意の配列に対してさらにフィルタリングされてもよい。

【0036】

制御配列のセットは、第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定し、関連する複数のシフト塩基のそれぞれを、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列として保存することによって生成され得る。

【0037】

コアプロモーター配列のセットおよび制御配列のセットは、生成モデルのための訓練データセットとして保存されてもよい。

2. 予測モデルのための訓練データの生成

プロモーター配列データは、CAGEピーク閾値を適用することによってフィルタリングされてもよい。例えば、プロモーター配列データの上位CAGEピークのみを使用してもよい。一実施形態では、CAGEピーク閾値は、予測モデル用に選択されたコアプロモーターの強度が、生成モデル用に選択されたコアプロモーターの強度と合致するように設定されてもよい。コアプロモーターの強度が合致する場合、新規に生成されたコアプロモーターの分類は、より信頼性が高い場合がある。しかしながら、予測モデルに採用される機械学習モデル (例えば、ロジスティック回帰モデル) が、生成モデルに採用される機械学習モデル (例えば、ニューラルネットワーク) よりもバイアスが大きい可能性があるため、コアプロモーターの数は、予測モデルの方が少なくてもよい。そのため、予測モデルは過剰適合する可能性が低く、より少ない例で訓練することができる。生成モデルにおける過剰適合を避けるために、予測モデルよりも多くの実施例を訓練に使用してもよい。予測モデルと生成モデルの両方に対して同様の強度のコアプロモーターを確保するために、閾値はそれに応じて選択されるべきであるが、予測モデルに対してそれほど重要ではない可能性がある生成モデルに対して十分な数のコアプロモーターを確保するために、入力データを選択すべきである。

【0038】

フィルタリングされたプロモーター配列データをさらにフィルタリングして、生成モデルのために生成された訓練データセットのピークのいずれかと重複する任意の配列データを除去してもよい。次いで、コアプロモーター配列のセットは、各TSSについて、TSSサミットを5'方向に数塩基、および3'方向に数塩基拡張することによって決定され得る。例えば、長さ100bpのコアプロモーター配列を作成するために、5'方向に49bp、および3'方向に50bpのヌクレオチドを決定することができる。コアプロモーター配列のセットは、ヒトゲノムアセンブリ (hg19) 中にNsを含有する任意の配列に対してさらにフィルタリングされてもよい。

【0039】

制御配列のセットは、コアプロモーター配列のセットを、5'方向または3'方向に塩基数だけシフトさせることによって生成されてもよい。例えば、候補コアプロモーター配列を5'方向に50,000bpシフトさせることによって。制御配列をフィルタリングして、任意のCAGEピークと重複する任意の制御配列、および生成モデルのための制御配列のセットと重複する任意の制御配列を除去することができ、制御配列は、コアプロモーターではないとして標識されてもよい。

【0040】

一実施形態では、遺伝的データを受信することを含む予測モデル用の訓練データセットを生成するための方法が説明されている。遺伝的データは、第一の複数のヌクレオチド配列を含むことができる。第一の複数のヌクレオチド配列は、プロモーター配列を含み得る。複数のヌクレオチド配列の各ヌクレオチド配列は、関連発現スコアを有する少なくとも一つのTSSを含み得る。関連発現スコアは、CAGEピークを含み得る。遺伝的データは、プロモーター配列データを含んでもよい。第一の複数のヌクレオチド配列をフィルタ

10

20

30

40

50

リングして、生成モデルのために生成された訓練データセットのピークのいずれかと重複する任意の配列データを除去することができる。複数のサミットヌクレオチド塩基が、TSSで決定されてもよい。サミットヌクレオチド塩基を決定することは、最も強いCAGEシグナルを有するヌクレオチド塩基を決定することを含み得る。各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定してもよい。関連する複数の周辺塩基を決定することは、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定し、それによって候補コアプロモーター配列を形成することを含み得る。5'方向の第一の複数のヌクレオチド塩基は、49個のヌクレオチド塩基を含むことができ、3'方向の第二の複数のヌクレオチド塩基は、50個のヌクレオチド塩基を含むことができる。各サミットヌクレオチド塩基およびその関連する複数の周辺塩基は、コアプロモーターとして標識された第二の複数のヌクレオチド配列(コアプロモーター配列のセット)として保存されてもよい。コアプロモーター配列のセットは、例えば、ヒトゲノムアセンブリ(hg19、hg38など)中にNsを含有する任意の配列に対してさらにフィルタリングされてもよい。

10

#### 【0041】

制御配列のセットは、第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定し、関連する複数のシフト塩基のそれぞれを、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列(制御配列のセット)として保存することによって生成され得る。制御配列のセットをフィルタリングして、任意のCAGEピークと重複する任意の制御配列、および生成モデルのための制御配列のセットと重複する任意の制御配列を除去することができる。

20

#### 【0042】

コアプロモーター配列のセットおよび制御配列のセットは、予測モデルの訓練データセットとして保存されてもよい。

#### D. 生成モデル

##### 1. 生成モデルを生成する方法

本明細書において、一つまたは複数のリカレントニューラルネットワーク(RNN)、例えば、長短期メモリ(LSTM)リカレントニューラルネットワーク(LSTM-RNN)などを使用して、長さの異なる新規のコアプロモーター配列を生成する技術が開示される。LSTM-RNNモデルをシード配列に適用して、シード配列から可能性の高い次のヌクレオチドを予測することができる。可能性の高い次のヌクレオチドを、シード配列に連結し、得られた配列を、LSTM-RNNモデルに戻して、シード配列と以前に決定した可能性の高い次のヌクレオチドから別の可能性の高い次のヌクレオチドを予測することができる。LSTM-RNNモデルを使用して、任意の長さのコアプロモーター配列を生成してもよい。

30

#### 【0043】

RNNは、ユニット間の接続が有向サイクルを形成する人工ニューラルネットワークの一種である。RNNは、ネットワークが動的な時間的挙動を示すことができる内部状態を有する。例えば、フィードフォワードニューラルネットワークとは異なり、RNNは、その内部メモリを使用して、任意の入力シーケンスを処理できる。LSTM-RNNは、標準的なニューラルネットワークユニットの代わりに、またはそれに加えて、LSTMユニットをさらに含む。LSTMユニット、またはブロックは、任意の長さの時間にわたって値を記憶または保存できる、「スマート」ユニットである。LSTMブロックには、その入力に記憶するのに十分重要である時、値を記憶し続けるべきまたは忘れるべき時、および値を出力するべき時を決定するゲートを含む。

40

#### 【0044】

議論の明確化のために、本開示全体を通して議論される実施形態は、LSTM-RNNに関して議論される。しかしながら、様々なタイプのRNNが、例えば、メモリ細胞シーケンシング操作に関するバリエーション(例えば、双方向、単方向、後方視単方向、または

50

後方視ウィンドウを有する前方視単方向)またはメモリ細胞タイプに関するバリエーション(例えば、LSTMバリエーションまたはゲーティングされた反復単位(GRU))を含む、記載した実施形態で使用され得る。BLSTM-RNNでは、出力は過去と未来両方の状態情報に依存する。さらに、乗算ユニットのゲートにより、メモリ細胞は、過去と未来両方のイベントの長いシーケンスにわたって情報を保存し、アクセスすることができる。さらに、他のタイプの位置認識ニューラルネットワーク、または連続予測モデルを、LSTM-RNNまたはBLSTM-RNNの代わりに使用することができる。本開示全体を通して示されるように、LSTM-RNNは、入力層、出力層、および一つまたは複数の非表示層を含む。

#### 【0045】

10

図3、図4A、図4B、および図5はそれぞれ、ニューラルネットワーク300、RNNブロック400、およびLSTM RNNブロック500の概要を提供するために示される。図3は、ニューラルネットワーク300の一例を示す。ニューラルネットワーク300は、入力ノード、ブロック、またはユニット302、出力ノード、ブロック、またはユニット304、および非表示ノード、ブロック、またはユニット306を含む。入力ノード302は、接続308を介して非表示ノード306に接続され、非表示ノード306は、接続310を介して出力ノード304に接続される。

#### 【0046】

入力ノード302は、入力データに対応する一方、出力ノード304は、入力データの関数として出力データに対応する。例えば、入力ノード302は、入力配列に対応してもよく、出力ノード304は、出力配列またはヌクレオチドに対応してもよい。ノード306は、ニューラルネットワークモデル自体がノードを生成する非表示ノードである。ノード306の一つの層のみが図示されるが、実際には、ノード306の二つ以上の層が通常存在する。

20

#### 【0047】

したがって、ニューラルネットワーク300を構築するために、手動でまたは別の方法で既に出力データにマッピングされた入力データ形式で訓練データが、ネットワーク300を生成するニューラルネットワークモデルに提供される。したがって、モデルは、非表示ノード306、入力ノード302と非表示ノード306との間の接続310の重み、非表示ノード306と出力ノードとの間の接続310の重み、および非表示ノード306自体の層間の接続の重みを生成する。その後、ニューラルネットワーク300は、出力データが未知である入力データに対して使用して、所望の出力データを生成することができる。

30

#### 【0048】

RNNは、ニューラルネットワークの一種である。一般的なニューラルネットワークは、入力データを処理して出力データを生成する間、いかなる中間データも保存しない。比較すると、RNNはデータを持続し、そうでない一般的なニューラルネットワークよりもその分類能力を向上させることができる。

#### 【0049】

図4Aは、RNNであるニューラルネットワーク300の非表示ノード306を代表する、RNNブロック400の一例をコンパクトに表記したものを示す。RNNブロック400は、入力ノード302のうちの一つからつながる図3の接続308であり得る、または別の非表示ノード306からつながる接続であり得る、入力接続402を有する。RNNブロック400は同様に、出力ノード304のうちの一つにつながる図3の接続310であり得る、または別の非表示ノード306につながる接続であり得る、出力接続404を有する。

40

#### 【0050】

RNNブロック400は概して、出力接続404上に提供される情報を生成するために、入力接続402上に提供される情報に対して(少なくとも)実行される処理406を含むと言われる。処理406は、典型的には関数の形態である。例えば、関数は、出力接続404を入力接続402にマッピングする、同一性活性化関数であってもよい。関数は、

50

入力接続 402 に基づいて、範囲 (0、1) 内の値を出力できる、ロジスティックシグモイド関数などのシグモイド活性化関数であってもよい。関数は、入力接続 402 に基づいて、範囲 (-1、1) 内の値を出力できる、ハイパボリックロジスティックタンジェント関数などのハイパボリックタンジェント関数であってもよい。

#### 【0051】

RNNブロック 400 はまた、それ自体の時間的後継回路に戻る時間的ループ接続 408 を有する。接続 408 は、RNNブロック 400 を再帰的にするものであり、複数のノード内のかかるループの存在は、ニューラルネットワーク 300 を再帰的にするものである。したがって、RNNブロック 400 が接続 404 上に出力する情報 (または他の情報) は、接続 408 上に持続することができ、これに基づいて、接続 402 上で受信された新しい情報を処理することができる。すなわち、RNNブロック 400 が接続 404 上に出力する情報は、RNNブロック 400 が次に入力接続 402 上で受信する情報とマージ、または連結され、処理 406 を介して処理される。

10

#### 【0052】

図 4B は、RNNブロック 400 の拡張表記を示す。RNNブロック 400' および接続 402'、404'、406'、408' は、同じ RNNブロック 400 および接続 402、404、406、408 であるが、時間的には遅い時間のものである。したがって、図 4B は、RNNブロック 400' が、早い時間で (同じ) RNNブロック 400 によって提供される接続 406 に提供される情報を、遅い時間で受信することを示す。遅い時間での RNNブロック 400' は、それ自体が接続 406' 上でさらに遅い時間にそれ自体に情報を提供することができる。

20

#### 【0053】

LSTM-RNN は、RNN の一種である。理論上の一般的な RNN は、短期的および長期的に情報を維持することができる。しかしながら、実際には、こうした RNN は、長期的に情報を維持する能力が証明されていない。より厳密には、一般的な RNN は、実質的に長期的な依存関係を学習することができない。つまり、RNN は、以前に比較的長い期間処理した情報に基づいて情報を処理することができない。比較すると、LSTM-RNN は、長期的な依存関係を学習できる特別なタイプの RNN であり、それゆえ、長期的に情報を維持することのできるタイプの RNN である。

#### 【0054】

30

図 5 は、例示的な LSTM-RNNブロック 500' を示す。LSTM-RNNブロック 500' は、図 4A および図 4B の RNNブロック 400 / 400' の接続 402 / 402' および 404 / 404' ならびに処理 406 / 406' と同等の入力接続 502'、出力接続 504' および処理 506' を有する。しかしながら、RNNブロック 400 / 400' の時間的インスタンスを接続する単一の時間的ループ接続 408 / 408' を有するのではなく、LSTM-RNNブロック 500' は、LSTM-RNNブロック 500 の時間的インスタンス間で情報が持続する、二つの時間的ループ接続 508' および 510' を有する。

#### 【0055】

入力接続 502' 上の情報は、LSTM-RNNブロックの以前の時間的インスタンスから接続 508 上に提供された持続的情報にマージされ、処理 506' を受ける。処理 506' の結果は、もしあるとしても、LSTM-RNNブロックの以前の時間的インスタンスから接続 510 上に提供された持続的情報とどのように組み合わせられるかは、ゲート 512' および 514' を介して制御される。接続 502' および 508' のマージ情報に基づいて動作するゲート 512' は、接続 510 上の持続情報の通過 (または非通過) を許可する、要素別積演算子 516' を制御する。同じ基準で動作するゲート 514' は、処理 506' の出力の通過 (または非通過) を許可する要素別演算子 518' を制御する。

40

#### 【0056】

演算子 516' および 518' の出力は、追加演算子 520' を介して合計され、LSTM-RNNブロック 500' の現在のインスタンスの接続 510' 上の持続的情報として渡される。したがって、接続 510' の持続情報に接続 510 の持続情報を反映する程度、およ

50



び接続 5 1 0 ' のこの情報が処理 5 0 6 ' の出力を反映する程度は、ゲート 5 1 2 ' および 5 1 4 ' によって制御される。そのため、情報は、必要に応じて、L S T M - R N N ブロックの複数の時間的インスタンス全体、またはそれらにわたって持続し得る。

【 0 0 5 7 】

L S T M - R N N ブロック 5 0 0 ' の現在のインスタンスの出力は、それ自体が、R N N の次の層への接続 5 0 4 ' 上で提供され、また接続 5 0 8 ' 上で L S T M - R N N ブロックの次の時間的インスタンスにも持続する。この出力は、別の要素別積演算子 5 2 2 ' によって提供され、ゲート 5 2 4 ' および 5 2 6 ' によってそれぞれ制御されるように、接続 5 1 0 ' 上にも提供される情報と、接続 5 0 2 ' および 5 0 8 上のマージされた情報との組み合わせを通過させる。このようにして、次に、図 5 の L S T M - R N N ブロック 5 0 0 ' は、長期情報および短期情報の両方を持続させることができる。

10

【 0 0 5 8 】

以下でさらに詳細に説明するように、例えば L S T M - R N N などの一つまたは複数の ( R N N ) を使用して、プロモーター配列を決定してもよい。L S T M - R N N は、訓練プロセスを介して、一連の入力特徴に基づいて一連のパラメータを予測するモデルを提供する。図 6 は、L S T M - R N N を訓練するための方法 6 0 0 の一例を示す。一部の変形では、訓練プロセスは、専用ソフトウェア (例えば、k e r a s、T e n s f o r F l o w スクリプト、C U R R E N N T ツールキット、または市販のツールキットの修正版など) および / または専用ハードウェアを使用して実施され得る。

【 0 0 5 9 】

20

工程 6 1 0 で、コンピューティングデバイスは、本明細書に記述されたように、訓練データセットを決定し得る。訓練データセットは、コアプロモーターとして標識されたコアプロモーター配列のセットと、コアプロモーターではないとして標識された制御配列のセットとを含んでもよい。訓練データセットは、L S T M - R N N の訓練および検証に使用され得る。一部の変形では、訓練データの第一の部分を訓練に使用してもよく、訓練データの第二の部分を試験 / 検証に使用してもよい。コアプロモーター配列のセットは、異なるクラス (シャープ / ブロード) の配列を、シード配列と予測標的の対に分割してもよい。シード配列は、例えば、1 0 個のヌクレオチドなど、任意の長さのヌクレオチド配列を含んでもよい。予測標的は、シード配列のすぐ後にあるヌクレオチドを含んでもよい。シード配列 / 予測標的対は、工程サイズ 1 のスライディングウィンドウアプローチを使用してコアプロモーター配列を分割することによって生成され得る。次いで、配列対は、数値符号化 (例えば、「A」: 0、「C」: 1、「G」: 2、「T」: 3) を使用してベクトル化されてもよい。コアプロモーターの異なるクラスは、その生物学において根本的に異なる。典型的には、プラスミドベースまたはウイルスベクターベースの導入遺伝子では、「シャープ」プロモータークラスのみが使用される。二つのクラスは、その配列内容が異なるため、二つのクラスは、訓練または配列生成のために混合することができない。言い換えれば、両方のクラスのコアプロモーターが訓練に使用された場合、「シャープ」コアプロモーターに通常見られるコアプロモーターモチーフからのシグナルは希釈され、非機能性配列の生成につながる可能性がある。代わりに、二つのタイプが分離され、その後、二つの別々に訓練されたモデルに基づいて新しい配列が生成され、それによって、「シャープ」または「ブロード」コアプロモーターの新規のインスタンスを生成することができる。どのタイプのコアプロモーターが生成されるかは、用途により異なるが、最も一般的には、「シャープ」タイプである。結論として、分離は用途にはあまり関係しないが (通常、シャープコアプロモーターのみが使用されるため)、シャープとブロードの混合が正しいシグナルを希釈することを考えると、適切な訓練にはより関係がある。しかしながら、別個のクラスを使用して、「ブロード」タイプのコアプロモーターを生成することもできる。

30

40

【 0 0 6 0 】

工程 6 2 0 で、コンピューティングデバイスは、L S T M - R N N のモデルセットアップを実行し得る。モデルセットアップには、L S T M - R N N の接続の重みを初期化する

50

ことが含まれ得る。一部の配置では、事前訓練は、重みを初期化するために実施され得る。その他の場合、重みは、分布スキームに従って初期化されてもよい。一つの分布スキームは、平均が0、標準偏差が0.1の正規分布に従って、重みの値を無作為化することを含む。一実施形態では、各ヌクレオチドは重みと関連付けられてもよい。例えば、重みは、A: 0.2、C: 0.05、G: 0.6、T: 0.15のように割り当てることができる。

#### 【0061】

工程630で、コンピューティングデバイスは、LSTM-RNNのモデル訓練を実行し得る。一部の变形では、モデル訓練は、例えば、最急降下、確率的勾配降下、または弾性逆伝播(RPROP)を含む、一つまたは複数の訓練技術を実施することを含み得る。訓練技術は、訓練セットをモデルに適用し、モデルを調整する。モデル訓練をスピードアップするために、並列訓練を実施することができる。例えば、訓練セットは、他のバッチと並列に処理される100個の配列のバッチに分割されてもよい。

10

#### 【0062】

工程640で、コンピューティングデバイスは、LSTM-RNNのモデル検証を実行し得る。一部の变形では、検証セットを、訓練されたモデルに適用し、ヒューリスティックを追跡してもよい。ヒューリスティックは、一つまたは複数の停止条件と比較されてもよい。停止条件が満たされる場合、訓練プロセスは終了し得る。一つまたは複数の停止条件のいずれも満たされない場合、工程630のモデル訓練および工程640の検証が繰り返される場合がある。工程630および640の各反復は、訓練エポックと呼ばれることがある。追跡され得るヒューリスティックの一部には、二乗誤差の合計(SSE)、加重二乗誤差の合計(W SSE)、回帰ヒューリスティック、または訓練エポックの数が含まれる。

20

#### 【0063】

##### 2. 生成モデルを使用する方法

図7は、プロモーター配列を決定するための訓練されたLSTM-RNN710を使用する例示的な流れを示す。LSTM-RNN710は、入力配列720(例えば、「シード」)を受けように構成されてもよい。入力配列720は、ヌクレオチド配列を含んでもよい。ヌクレオチド配列は、プロモーター配列(例えば、コアプロモーター配列)を含んでもよい。入力配列720は、長さを有してもよい。長さは、例えば、約5ヌクレオチド~約100ヌクレオチドであってもよい。入力配列720は、10ヌクレオチド長であってもよい。他の入力配列長、例えば、1、2、3、4、5、6、7、8、9、11、12、13、14、15、16、17、18、19、または20ヌクレオチド長が企図される。一実施形態では、入力配列720のソースは、訓練データセットからのランダムに選択されたコアプロモーターのランダム配列または第一の10nt(または任意の他の入力配列長)であってもよい。一実施形態では、実世界の例からの特定のコアプロモーターモチーフにつながる10nt(または任意の他の入力配列長)は、そのモチーフを有するコアプロモーターの生成を強制するために、入力配列720として選択されてもよい。別の実施形態では、ランダム配列を入力配列として使用してもよく、出力配列を特定のモチーフの存在についてスクリーニングしてもよい。

30

40

#### 【0064】

入力配列720は、入力層、一つまたは複数の非表示層、および出力層を介して、入力配列720を処理するLSTM-RNN710に入力される。LSTM-RNN710は、出力層を介して、可能性の高い次のヌクレオチドを出力する。したがって、LSTM-RNN710は、出力配列730を生成するために入力配列720に付加され得る、可能性の高い次のヌクレオチドを予測するように構成されてもよい。LSTM-RNN710は、ヌクレオチド確率に基づいて、可能性の高い次のヌクレオチドを予測するように構成されてもよい。ヌクレオチド確率は、例えば、A: 0.2、C: 0.05、G: 0.6、T: 0.15であり得る。LSTM-RNN710は、生成された出力配列730を入力として取り、出力配列730を新しい入力配列720として効果的に処理するように構成

50

されてもよい。LSTM-RNN710は、出力配列730に対して所望の長さが達成されるまで、可能性の高い次のヌクレオチドを繰り返し予測するように構成されてもよい。所望の出力長は、例えば、約20ヌクレオチド～約100ヌクレオチドであってもよい。所望の出力長は、50ヌクレオチドであってもよい。LSTM-RNN710は、前のヌクレオチドのプロモーター配列が与えられた新しいプロモーター配列において、次のヌクレオチドの確率分布を生成する。これにより、LSTM-RNN710は、一度に一つの新しいプロモーター配列を生成することができる。一実施形態では、任意の数のコアプロモーターを生成することができ、その後、GC含量またはコアプロモーターモチーフ含量（予測モデルで使用されるものと類似した特徴）などの特定の態様についてスクリーニングすることができる。

10

#### 【0065】

図8は、コアプロモーター配列の生成の視覚的描写である。シード配列802は、LSTM-RNN710に入力されてもよく、所望の最終コアプロモーター配列長が指定されてもよい。図8の実施例では、シード配列802は、10ヌクレオチド長であり、所望の最終コアプロモーター配列長は、14ヌクレオチドである。LSTM-RNN710は、シード配列802が与えられると、次に可能性の高いヌクレオチド804を予測し得る。次に可能性の高いヌクレオチド804をシード配列802に添加（連結）して、配列806を作成してもよい。LSTM-RNN710は、配列806の少なくとも一部が与えられると、次に可能性の高いヌクレオチド808を予測し得る。例えば、n個のヌクレオチドのスライディングウィンドウを使用して、次に可能性の高いヌクレオチド808を予測してもよい。スライディングウィンドウは、例えば、5、6、7、8、9、10、11、12、13、14、または15ヌクレオチド長であってもよい。次に可能性の高いヌクレオチド808を配列806に添加（連結）して、配列810を作成してもよい。LSTM-RNN710は、配列810の少なくとも一部が与えられると、次に可能性の高いヌクレオチド812を予測し得る。例えば、n個のヌクレオチドのスライディングウィンドウを使用して、次に可能性の高いヌクレオチド812を予測してもよい。次に可能性の高いヌクレオチド812を配列810に添加（連結）して、配列814を作成してもよい。LSTM-RNN710は、配列814の少なくとも一部が与えられると、次に可能性の高いヌクレオチド816を予測し得る。例えば、n個のヌクレオチドのスライディングウィンドウを使用して、次に可能性の高いヌクレオチド816を予測してもよい。次に可能性の高いヌクレオチド816を配列814に添加（連結）して、最終コアプロモーター配列818を作成してもよい。最終コアプロモーター配列818の長さは、所望の最終コアプロモーター配列長と等しく、最終コアプロモーター配列818は、コアプロモーターとして出力されてもよい。

20

30

#### 【0066】

一実施形態では、a)ヌクレオチド配列と配列長を受信すること、b)訓練された生成モデルに、ヌクレオチド配列を提供すること、c)生成モデルに基づいて、ヌクレオチド配列に関連付けられた次のヌクレオチドを決定すること、d)ヌクレオチド配列に次のヌクレオチドを付与すること、e)ヌクレオチド配列の長さが配列長に等しくなるまでb～dを繰り返すこと、およびf)ヌクレオチド配列をコアプロモーター配列として出力すること、を含む方法が記述される。配列長は、約50ヌクレオチド～約100ヌクレオチドとすることができる。

40

#### 【0067】

本明細書に記載されるように、定義された長さおよび/または定義された配列含量を有するコアプロモーター配列を生成するための方法およびシステムが提供される。生成されたコアプロモーター配列は、例えば、前臨床研究および遺伝子治療の両方におけるAAV、アデノ、またはレンチウイルスなどのウイルスベクター；Cas9発現、Cre-lox、TALENまたはジンクフィンガーヌクレアーゼを駆動するものなどのゲノム編集ベクター；発光または蛍光レポータープラスミド；高収率の抗体発現プラスミド；クローニングおよびエンジニアリングプラスミド；化学遺伝学（例えばDREADD）、および光

50

遺伝学などの任意の種類の導入遺伝子に使用することができる。一部の態様では、生成されたコアプロモーターのうちの一つまたは複数を含む構築物またはベクターは、核酸構築物と呼んでもよい。したがって、開示された核酸構築物は、生成されたコアプロモーターおよび一つまたは複数の導入遺伝子を含み得る。一部の態様では、開示される核酸構築物は、任意の発現ベクター、ウイルスまたは非ウイルスであってもよい。一部の態様では、核酸構築物は、直鎖状または環状であってもよい。環状核酸構築物は、プラスミドまたはベクターと呼んでもよい。

#### 【0068】

したがって、生成されたコアプロモーター配列は、遺伝子治療のためのウイルスベクター上の導入遺伝子に使用され得る。現在、AAVは遺伝子治療のゴールドスタンダードだが、ゲノムサイズが非常に限られているため、AAVベクターで使用されるいかなる要素も、サイズを最適化する必要がある。例えばCas9などの一部の導入遺伝子のコード配列は、サイズを最適化することができないが、他の調節要素（コアプロモーターなど）は可能である。したがって、100nt長であり得る内因性コアプロモーターを使用する代わりに、開示された方法およびシステムは、50nt長であり、定義されたモチーフ含量を有し、したがって、任意の内因性コアプロモーターよりも遺伝子発現を駆動するのにさらに効率的なコアプロモーターを生成することができる。

10

#### 【0069】

一実施形態では、生成されたコアプロモーター配列の配列含量は、遺伝子治療設定における自然免疫反応を回避するために定義されてもよい。コアプロモーターは、多くの場合、高い含量のCpGジヌクレオチドを有し、これはTLRベースの自然免疫反応をトリガすることができる。説明された生成モデルを使用して、CpGジヌクレオチドを欠くコアプロモーター配列を生成することができる。

20

#### 【0070】

方法は、コアプロモーター配列に基づいてプロモーターを操作することをさらに含んでもよい。方法は、プロモーターを核酸構築物に挿入することをさらに含んでもよい。プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含む。方法は、核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することをさらに含んでもよい。一部の態様では、本開示の核酸構築物を含む、任意の公知のウイルスベクターを作製することができる。一部の態様では、方法は、生成されたコアプロモーターを含む任意の公知の非ウイルスベクター（例えば、DNAベースのベクター）を作製することを含んでもよい。

30

#### 【0071】

用語「発現ベクター」は、（例えば、転写制御要素に連結された）細胞による発現に適した形態の導入遺伝子を含む任意のベクター（例えば、プラスミド、コスミドまたはファージ染色体）を含む。一部の態様では、プラスミドは一般的に使用されるDNAベクターの形態であるため、「プラスミド」と「ベクター」は互換的に使用される。さらに、本発明は、同等の機能を果たす他のベクターを含むことを意図している。

#### 【0072】

##### E. 予測モデル

40

##### 1. 予測モデルを生成する方法

ここで図9を参照すると、予測モデルを生成するための方法が説明されている。説明された方法は、所定の配列に対するプロモーター状態（例えば、プロモーター/非プロモーター）を予測するよう構成されている少なくとも一つのMLモジュール930である、訓練モジュール920による、一つまたは複数の訓練データセット910の分析に基づき、訓練するための機械学習（「ML」）技術を使用してもよい。

#### 【0073】

訓練データセット910は、コアプロモーター配列（YES）として標識されたコアプロモーター配列のセット、およびコアプロモーター配列（NO）ではないとして標識された制御配列のセットを含んでもよい。このようなデータは、全体的または部分的に、本明

50

細書に記載のプロモーター配列データから導出されてもよい。

【0074】

コアプロモーター配列のセットおよび制御配列のセットのサブセットは、訓練データセット910または試験データセットにランダムに割り当てられてもよい。一部の実施では、訓練データセットまたは試験データセットへのデータの割り当ては完全に無作為ではない場合がある。この場合、一つまたは複数の基準が、割り当て中に使用され得る。一般に、任意の好適な方法を使用して、データを訓練データセットまたは試験データセットに割り当ててもよい一方で、はいおよびいいえの標識分布が、訓練データセットおよび試験データセットにおいていくらか類似していることを保証し得る。

【0075】

訓練モジュール920は、一つまたは複数の特徴選択技術により、訓練データセット910における複数のコアプロモーター配列（例えば、はいとして標識された）および/または複数の制御配列（例えば、いいえとして標識された）から特徴セットを抽出することによって、MLモジュール930を訓練してもよい。訓練モジュール920は、正の例（例えば、はいであると標識された）の統計上有意な特徴および負の例（例えば、いいえであると標識された）の統計上有意な特徴を含む訓練データセット910から、特徴セットを抽出することによって、MLモジュール930を訓練してもよい。

【0076】

訓練モジュール920は、様々な方法で、訓練データセット910から特徴セットを抽出してもよい。訓練モジュール920は、異なる特徴抽出技術を使用して、各回に特徴抽出を複数回実施し得る。一例では、異なる技術を使用して生成される特徴セットは各々が、異なる機械学習ベースの分類モデル940を生成するために使用され得る。例えば、最も高い品質の測定基準を伴う特徴セットが、訓練における使用のために選択され得る。訓練モジュール920は、新規の配列（例えば、未知のプロモーター状態を有する）が、プロモーターである可能性が高いか、または低いかを示すよう構成されている、一つまたは複数の機械学習ベースの分類モデル940A~940Nを構築するための特徴セットを使用してもよい。

【0077】

訓練データセット910を分析して、訓練データセット910における特徴とはい/いいえの標識の間の任意の依存性、関連性、および/または相関を決定してもよい。識別された相関は、異なるはい/いいえの標識と関連する特徴のリストの形態を有してもよい。本明細書で使用される場合、用語「特徴」は、データのある項目が、一つまたは複数の特定のカテゴリ内にあるか否かを決定するために使用され得るデータの項目の任意の特徴を指し得る。例として、本明細書に記載される特徴は、一つまたは複数の配列パターン、GC含量、CpG含量、公知のコアプロモーター配列モチーフ、ATG頻度、および/または相対エントロピーを含んでもよい。公知のコアプロモーター配列モチーフの発生については、TSSに対する相対位置決め方式も考慮され得る。相対エントロピーは、固定ヌクレオチド分布に基づき、ある配列がランダム配列とどれほど類似しているかの尺度である（コアプロモーター配列はランダムではない）。ランダムDNAの確率の例は、ランダムDNAの確率：A：0.3、C：0.2、T：0.2、G：0.3である。図10は、プロモーター状態を予測する様々な特徴の相対的有意性を示す。例えば、CpGジヌクレオチド含量は、コアプロモーターの同一性と最も強い正の相関を有する。ランダムDNAは、比較的少数のCpGを有する一方で、コアプロモーターは、多数を有することができる（コアプロモーターのサブセットは、CpGアイランドと呼ばれる）。TATAボックスの有無も重要である。

【0078】

図9に戻ると、特徴選択技術は、一つまたは複数の特徴選択ルールを含み得る。一つまたは複数の特徴選択ルールは、特徴発生ルールを含み得る。特徴発生ルールは、訓練データセット910においていずれの特徴が閾値の回数にわたって生じるかを決定すること、および閾値を満たすそれらの特徴を、特徴として特定することを含み得る。

10

20

30

40

50

## 【 0 0 7 9 】

単一の特徴選択ルールを、特徴を選択するために適用してもよく、または複数の特徴選択ルールを、特徴を選択するために適用してもよい。特徴選択ルールは、カスケード方式で適用されてもよく、特徴選択ルールは、特定の順序で適用され、以前のルールの結果に適用される。例えば、特徴発生ルールは、訓練データセット 9 1 0 に適用されて、特徴の第一のリストを生成し得る。特徴の最終リストは、一つまたは複数の特徴群（例えば、プロモーター状態を予測するために使用され得る特徴の群）を決定するためのさらなる特徴選択技術により分析されてもよい。任意の好適な計算技術を使用して、フィルター方法、ラッパー方法、および／または埋め込み方法などの任意の特徴選択技術を使用して、特徴群を特定し得る。一つまたは複数の特徴群は、フィルター方法に従い選択されてもよい。フィルター方法には、例えば、ピアソンの相関、線形判別分析、分散分析（ANOVA）、カイ二乗、それらの組み合わせなどが含まれる。フィルター方法に従った特徴の選択は、任意の機械学習アルゴリズムから独立している。代わりに、特徴は、転帰変数（例えば、はい／いいえ）との相関について、様々な統計検定におけるスコアに基づいて選択され得る。

10

## 【 0 0 8 0 】

別の例として、一つまたは複数の特徴群は、ラッパー方法により選択されてもよい。ラッパー方法は、特徴のサブセットを使用し、特徴のサブセットを使用して機械学習モデルを訓練するように構成され得る。以前のモデルから引き出された推論に基づいて、特徴は、サブセットから追加および／または削除され得る。ラッパー方法は、例えば、前方特徴量選択、後方特徴量削減、再帰的特徴量削減、それらの組み合わせなどを含む。一例として、前方特徴選択を使用して、一つまたは複数の特徴群を識別してもよい。前方特徴量選択は、機械学習モデルにおける特徴なしに始まる反復方法である。各反復において、モデルを最良に改善する特徴が、新たな変数の追加によって機械学習モデルの性能が改善されなくなるまで加えられる。一例として、後方排除を使用して、一つまたは複数の特徴群を識別してもよい。後方削減は、機械学習モデルにおける全ての特徴で始まる反復方法である。各反復では、最下位の特徴が、特徴の除去時に改善が観察されなくなるまで除去される。再帰的特徴除去を使用して、一つまたは複数の特徴群を識別してもよい。再帰的特徴量削減は、性能が最良である特徴サブセットを見出すことを目指す貪欲最適化アルゴリズムである。再帰的特徴量削減によって、モデルが反復的に作成され、各反復で最良または最悪の性能の特徴を別にしておく。再帰的特徴量削減によって、全ての特徴が消耗するまで、特徴が残っている次のモデルが構築される。再帰的特徴量削減によって、次に、それらの削減の順序に基づいて特徴がランク付けされる。

20

30

## 【 0 0 8 1 】

さらなる例として、一つまたは複数の特徴群は、埋め込み方法により選択されてもよい。埋め込み方法によって、フィルター方法とラッパー方法の質が組み合わされる。埋め込み方法には、例えば、過学習を低下させるためのペナルティ機能を実施する、最小絶対収縮および選択演算子（LASSO）およびリッジ回帰が含まれる。例えば、LASSO回帰によって、係数の大きさの絶対値に相当するペナルティを加えるL1正則化が実施され、リッジ回帰によって、係数の大きさの二乗に相当するペナルティを加えるL2正則化が実施される。

40

## 【 0 0 8 2 】

訓練モジュール 9 2 0 によって特徴セットが生成された後、訓練モジュール 9 2 0 によって、特徴セットに基づいて、機械学習ベースの分類モデル 9 4 0 が生成され得る。機械学習ベースの分類モデルは、機械学習技術を使用して生成される、データ分類のための複雑な数学的モデルを指し得る。一例では、機械学習ベースの分類モデル 9 4 0 は、境界特徴を表すサポートベクトルのマップを含み得る。この例では、境界特徴は、ある特徴セット内の最高ランクの特徴から選択されても、かつ／またはそれらを表してもよい。

## 【 0 0 8 3 】

訓練モジュール 9 2 0 は、それぞれの分類カテゴリー（例えば、はい、いいえ）につい

50

ての機械学習ベースの分類モデル 9 4 0 A ~ 9 4 0 N を構築するための訓練データセット 9 1 0 から決定または抽出された特徴セットを使用してもよい。いくつかの例では、機械学習ベースの分類モデル 9 4 0 A ~ 9 4 0 N を、単一の機械学習ベースの分類モデル 9 4 0 に組み合わせてもよい。同様に、ML モジュール 9 3 0 は、単一もしくは複数の機械学習ベースの分類モデル 9 4 0 を含有する単一の分類指標、および / または単一もしくは複数の機械学習ベースの分類モデル 9 4 0 を含有する複数の分類指標を表し得る。

#### 【 0 0 8 4 】

特徴を、機械学習アプローチ、例えば判別分析；決定木；最近傍（NN）アルゴリズム（例えば、k - NN モデル、レプリケータ - NN モデルなど）；統計アルゴリズム（例えば、ベイジアンネットワークなど）；クラスタリングアルゴリズム（例えば、k 平均値、平均値シフトなど）；ニューラルネットワーク（例えば、リザーバネットワーク、人工ニューラルネットワークなど）；サポートベクター機械（SVM）；ロジスティック回帰アルゴリズム；線形回帰アルゴリズム；マルコフモデルまたはチェーン；主成分分析（PCA）（例えば、線形モデルについて）；多層パーセプトロン（MLP）ANN（例えば、非線形モデルについて）；リザーバネットワークの複製（例えば、非線形モデルについて、通常は時系列について）；ランダムフォレスト分類；それらの組み合わせおよび / または同様のものを使用して訓練された分類モデルにおいて組み合わせてもよい。得られた ML モジュール 9 3 0 は、プロモーター状態を新規の配列に割り当てるための、それぞれの特徴についての決定ルールまたはマッピングを含んでもよい。

#### 【 0 0 8 5 】

一実施形態では、訓練モジュール 9 2 0 は、畳み込みニューラルネットワーク（CNN）として機械学習ベースの分類モデル 9 4 0 を訓練してもよい。CNN は、少なくとも一つの畳み込み特徴層および最終の分類層（softmax）につながる三つの完全に連結した層を含む。最終の分類層を最終的に適用して、当該技術分野で公知の softmax 関数を使用して、完全に結び付けられた層の出力を組み合わせてもよい。

#### 【 0 0 8 6 】

特徴および ML モジュール 9 3 0 は、試験データセット内の配列のプロモーター状態を予測するために使用され得る。一例では、それぞれの配列の予測結果は、配列がプロモーターである可能性または確率に対応する信頼レベルを含む。信頼レベルは、ゼロから一の間の値であってもよく、それは、配列が、はい / いいえのプロモーター状態に属する可能性を表してもよい。一例では、二つの状態（例えば、はいおよびいいえ）があるとき、信頼レベルは、値 p に対応してもよく、それは、特定の配列が、第一の状態（例えば、はい）に属する可能性を指す。この場合では、値  $1 - p$  は、特定の配列が、第二の状態（例えば、いいえ）に属する可能性を指し得る。一般に、複数の信頼レベルは、試験データセットの各配列について、および三つ以上の状態がある場合、各特徴について提供され得る。最も高性能の特徴は、各試験配列について得られた結果を、各試験配列についての公知のはい / いいえのプロモーター状態と比較することによって決定されてもよい。一般に、最も高性能の特徴は、公知のはい / いいえプロモーター状態と密接に一致する結果を有するであろう。最も高性能の特徴を使用して、配列のはい / いいえプロモーター状態を予測してもよい。例えば、新規の配列が、決定 / 受信されてもよい。新規の配列は、最も高性能の特徴に基づき、新規の配列を、プロモーター（はい）またはプロモーターではない（いいえ）のいずれかとして分類し得る ML モジュール 9 3 0 に適用されてもよい。

#### 【 0 0 8 7 】

図 1 1 は、訓練モジュール 9 2 0 を使用して、ML モジュール 9 3 0 を生成するための例となる訓練方法 1 1 0 0 を説明するフローチャートである。訓練モジュール 9 2 0 によって、教師あり、教師なし、および / または半教師あり（例えば、補強ベース）の機械学習ベースの分類モデル 9 4 0 を実施することができる。図 1 1 に例証する方法 1 1 0 0 は、教師あり学習方法の例であり；訓練方法のこの例の変形を以下で考察するが、しかし、他の訓練方法は、教師なしおよび / または半教師ありの機械学習モデルを訓練するために類似的に実施することができる。

10

20

30

40

50

## 【 0 0 8 8 】

訓練方法 1 1 0 0 は、工程 1 1 1 0 において第一の配列データを決定（例えば、アクセス、受信、検索など）してもよい。配列データは、コアプロモーター配列の標識されたセットおよび制御配列の標識されたセットを含んでもよい。標識は、プロモーター状態（例えば、はいまたはいいえ）に対応してもよい。

## 【 0 0 8 9 】

訓練方法 1 1 0 0 は、工程 1 1 2 0 において、訓練データセットおよび試験データセットを生成してもよい。訓練データセットおよび試験データセットは、標識された配列を訓練データセットまたは試験データセットのいずれかに無作為に割り当てることによって、生成されてもよい。一部の実施では、訓練または試験データとしての標識された配列の割り当ては、完全に無作為でなくてもよい。一例として、標識された配列の大部分を使用して、訓練データセットを生成してもよい。例えば、標識された配列の 7 5 % を使用して、訓練データセットを生成してもよく、2 5 % を使用して、試験データセットを生成してもよい。別の例では、標識された受容体配列の 8 0 % を使用して、訓練データセットを生成してもよく、2 0 % を使用して、試験データセットを生成してもよい。

## 【 0 0 9 0 】

訓練方法 1 1 0 0 は、工程 1 1 3 0 において、例えば、プロモーター状態（例えば、はい対いいえ）の異なる分類の中で区別するための分類指標によって使用することができる一つまたは複数の特徴を決定（例えば、抽出、選択など）してもよい。一例として、訓練方法 1 1 0 0 は、標識された配列からセットの特徴を決定してもよい。さらなる例では、特徴のセットは、訓練データセットまたは試験データセットのいずれかにおいて標識された配列以外の標識された配列から決定されてもよい。言い換えると、標識された配列は、機械学習モデルの訓練のためよりむしろ、特徴の決定のため使用され得る。このような標識された配列を使用して、特徴の初期のセットを決定してもよく、それは、訓練データセットを使用してさらに低減されてもよい。例として、本明細書に記載される特徴は、一つまたは複数の配列パターン、G C 含量、C p G 含量、公知のコアプロモーター配列モチーフ、A T G 頻度、および/または相対エントロピーを含んでもよい。公知のコアプロモーター配列モチーフの発生については、T S S に対する相対位置決め方式も考慮され得る。相対エントロピーは、固定ヌクレオチド分布に基づき、ある配列がランダム配列とどれほど類似しているかの尺度である（コアプロモーター配列はランダムではない）。ランダム D N A の確率の例は、ランダム D N A の確率：A : 0 . 3、C : 0 . 2、T : 0 . 2、G : 0 . 3 である。図 1 0 は、プロモーター状態を予測する様々な特徴の相対的有意性を示す。例えば、C p G ジヌクレオチド含量は、コアプロモーターの同一性と最も強い正の相関を有する。ランダム D N A は、比較的少数の C p G を有する一方で、コアプロモーターは、多数を有することができる（コアプロモーターのサブセットは、C p G アイランドと呼ばれる）。T A T A ボックスの有無も重要である。

## 【 0 0 9 1 】

図 1 1 に戻ると、訓練方法 1 1 0 0 は、工程 1 1 4 0 で、一つまたは複数の特徴を使用して、一つまたは複数の機械学習モデルを訓練し得る。一例では、機械学習モデルは、教師あり学習を使用して訓練され得る。別の例では、教師なし学習および半教師ありを含む、他の機械学習技術が用いられてもよい。1 1 4 0 で訓練された機械学習モデルは、解決される問題および/または訓練データセットで利用可能なデータに応じて、異なる基準に基づいて選択され得る。例えば、機械学習分類器は、異なる程度のバイアスを受け得る。したがって、二つ以上の機械学習モデルを、1 1 4 0 で訓練し、工程 1 1 5 0 で最適化し、改善し、相互検証することができる。

## 【 0 0 9 2 】

訓練方法 1 1 0 0 は、1 1 6 0 で予測モデルを構築するために、一つまたは複数の機械学習モデルを選択し得る。予測モデルは、試験データセットを使用して評価してもよい。予測モデルは、試験データセットを分析し、工程 1 1 7 0 において予測されるプロモーター状態を生成してもよい。予測されるプロモーター状態を、工程 1 1 8 0 において評価し



て、こうした値が、所望の精度レベルを達成したかどうかを決定することができる。予測モデルの性能は、予測モデルによって示される複数のデータ点の多数の真の陽性、偽陽性、真の陰性、および/または偽陰性の分類に基づいて、多数の方法で評価され得る。

【0093】

例えば、予測モデルの偽陽性は、予測モデルが、実際にはプロモーターではない配列を、誤ってプロモーターとして分類した回数を指し得る。逆に、予測モデルの偽陰性は、機械学習モデルが、実際には配列がプロモーターであるのに、プロモーター配列をプロモーターではないと分類した回数を指し得る。真陰性および真陽性は、予測モデルによって一つまたは複数の配列が、プロモーター、またはプロモーターではないとして正しく分類された回数を指し得る。これらの測定に関連するのは、想起および精度の概念である。一般に、想起とは、真陽性および偽陰性の合計に対する真陽性の比率を指し、それによって予測モデルの感度が定量化される。同様に、精度は、真の陽性と偽陽性との合計の正陽性の比を指す。このような所望の精度レベルに達すると、訓練期が終了し、予測モデル（例えば、MLモジュール930）が、工程1190において出力されてもよく、しかしながら、所望の精度レベルに達していないとき、訓練方法1100のその後の反復は、例えば、配列データのより大きな収集を考慮するなどの変動を伴って、工程1110において開始して行われてもよい。

10

【0094】

## 2. 予測モデルを使用する方法

図12は、機械学習ベースの分類器を使用して、ヌクレオチド配列がプロモーターであるかどうかを決定するための例示的なプロセスフローの図である。図12に図示するように、非分類配列1210は、MLモジュール930への入力として提供されてもよい。MLモジュール930は、分類結果1220に到達するために、機械学習ベースの分類器を使用して、未分類の配列1210を処理してもよい。

20

【0095】

分類結果1220は、非分類配列1210の一つまたは複数の特徴を識別し得る。例えば、分類結果1220は、非分類配列1210のプロモーター状態を識別してもよい（例えば、非分類配列1210がプロモーター機能を実行する可能性が高いかどうか）。

【0096】

MLモジュール930を使用して、生成モデル（例えば、LSTM-RNN710）によって生成された配列を分類してもよい。予測モデル（例えば、MLモジュール930）は、生成モデル（例えば、LSTM-RNN710）の品質管理機構として機能し得る。生成モデルによって生成された配列を実験設定で試験する前に、予測モデルを使用して、生成された配列がコアプロモーター活性に対して陽性であると予測されるかどうかを試験してもよい。

30

【0097】

## F. 使用方法

一部の態様では、特定のプロモーター配列は、本明細書に記載される方法のうちの一つまたは複数によって生成および/または識別される。プロモーター配列が識別されると、プロモーターを産生または操作することができる。一部の態様では、産生された、操作された、および合成されたという用語は、互換的に使用され得る。一部の態様では、プロモーターは、共通使用の技術に従って化学的に合成することができる。例えば、Beaucage et al. (1981) Tet. Lett. 22: 1859、および米国特許第4,668,777号を参照のこと。その全体が参照により本明細書に取り込まれる。かかる化学オリゴヌクレオチド合成は、Perkin-Elmer Corp., Foster City, Calif., USAの一部門であるApplied BiosystemsによるBiosearch 4600または8600 DNA合成装置、およびPerceptive Biosystems, Framingham, Mass., USAによるExpediteなどの市販の装置を使用して実施することができる。プロモーターはまた、Yu et al. Recent Pat DNA Gene Seq, 2012,

40

50

A p r ; 6 ( 1 ) : 1 0 - 2 1 に開示された特許に記載された技術のいずれかを使用して合成することができる。各文献は参照によりその全体が組み込まれる。

【 0 0 9 8 】

プロモーターが産生されると、プロモーターは核酸構築物に挿入され得る。一部の態様では、核酸構築物は、ウイルスベクターを産生するために使用されるプラスミドを含むがこれに限定されないプラスミドであってもよい。一部の態様では、プロモーターは、ウイルスの作製に使用されるプラスミドに既に存在する導入遺伝子（すなわち、対象遺伝子）の上流に挿入され得る。一部の態様では、プロモーター配列は、核酸配列を形成する導入遺伝子の上流に挿入されてもよく、次いで、核酸配列は、ウイルスの作製に使用されるプラスミドに挿入されてもよい。一部の態様では、プロモーターは、導入遺伝子を挿入する前にプラスミド内に挿入することができる。任意の公知のクローニング方法を使用して、プロモーターを含むプラスミドまたは核酸配列を作製することができる。例えば、一部の態様では、A A V などのウイルスベクターの作製に使用されるプラスミドを、一つまたは複数の制限エンドヌクレアーゼで切断することができる。5'末端および3'末端に制限エンドヌクレアーゼ部位を含み、間に特定のプロモーターを有する核酸配列を、制限エンドヌクレアーゼ部位に特異的な制限エンドヌクレアーゼを用いて切断することができる。一部の態様では、5'末端および3'末端に制限エンドヌクレアーゼ部位を含み、間に特定のプロモーターを有する核酸配列は、プロモーターの下流の導入遺伝子をさらに含み、制限エンドヌクレアーゼ部位の間にもある。一部の態様では、プラスミドの切断に使用される制限エンドヌクレアーゼは、5'末端および3'末端に制限エンドヌクレアーゼ部位を含み、間に特定のプロモーターを有する核酸配列の切断に使用される制限エンドヌクレアーゼと同一である。同じ制限エンドヌクレアーゼで切断することで、プラスミドに付着末端、および5'末端および3'末端に制限エンドヌクレアーゼ部位を含み、間に特定のプロモーターを有する核酸配列が産生される。一部の態様では、特定のプロモーターを含む核酸配列は、制限エンドヌクレアーゼで切断されたプラスミドに既に付着している特定の制限エンドヌクレアーゼ部位を末端に既に含有するように化学的に合成することができるため、制限エンドヌクレアーゼで核酸配列を切断する工程を省くことができる。最終的に、核酸配列および類似の付着末端を有するプラスミドを互いに接触させ、導入遺伝子の上流に特定のプロモーターを含む環状プラスミドの形成を可能にする。次いで、環状プラスミドを使用して、公知のウイルス産生方法を使用して、A A V などのウイルスを産生することができる。

【 0 0 9 9 】

G . 実施例

以下の実施例は、本方法およびシステムを例証する。以下の実施例は、その限定を意図するものではない。

【 0 1 0 0 】

本明細書に開示される方法を使用してコアプロモーター配列を識別した後、コアプロモーター活性を、生物学的システムを使用して試験することができる。図 1 3 は、プロモーターアッセイがどのように設計され得るかの一例の概要を示す概略図である。個々のコアプロモーター（C P）候補（対照または生成されたコアプロモーター）を試験するために、二つのレポーター構築物を設計した。レポーター構築物は、N a n o L u c ルシフェラーゼ（N l u c、P r o m e g a）のコード配列、S V 4 0 後期ポリ阿德ニル化シグナル、および肝特異的エンハンサー（K h e r a d p o u r e t a l . , d o i : h t t p : / / d x . d o i . o r g / 1 0 . 1 1 0 1 / g r . 1 4 4 8 9 9 . 1 1 2 から）、またはエンハンサー省略のいずれかを含有し、ユニバーサルリバースプライマー結合部位が続く。さらに、P C R 中にコアプロモーター候補を追加するために、N a n o L u c コード配列の上流にプライマー結合部位を導入した。活性化エンハンサーをレポーター遺伝子の下流に配置して、エンハンサー内で転写が開始されないようにすることができる。

【 0 1 0 1 】

レポーター構築物を、二本鎖 D N A サンプル（g B l o c k、I D T）として注文し、

水中で  $10 \text{ ng} / \mu \text{I}$  に希釈した。ルシフェラーゼ活性をアッセイするために細胞内に導入できるレポーター構築物を生成するために、これら二つのテンプレートをPCRのユニバーサルテンプレートとして使用し、その間、それぞれのコアプロモーター候補を5'オリゴヌクレオチドを介して導入した。そのために、コアプロモーター配列、続いてレポーター構築物の5'プライマー結合部位に対応する配列を含有するオリゴヌクレオチドを注文した。得られたdsDNA PCR産物は、コアプロモーター候補、Nluc CDS、SV40後期ポリA、およびエンハンサー（図13下）、またはバックグラウンド対照としてのエンハンサーなし（図13上）のいずれかからなる。この直鎖状dsDNA産物は、ルシフェラーゼ活性を評価するために、選択した細胞株へのトランスフェクションに直接使用することができる。PCRについては、 $25 \mu \text{L}$ の $2 \times \text{Q5 hot start master mix}$ （NEB）、 $2.5 \mu \text{L}$ のフォワードオリゴ、 $2.5 \mu \text{L}$ のリバースオリゴ、 $1 \text{ ng}$ のgFM15テンプレート $20 \mu \text{L}$ の $\text{H}_2\text{O}$ を使用した。この反応混合物を、以下のPCRプログラムで増幅した： $98^\circ \text{C}$ で30秒、 $98^\circ \text{C}$ で10秒、 $68^\circ \text{C}$ で30秒、 $72^\circ \text{C}$ で15秒、工程2へさらに24回、 $72^\circ \text{C}$ で2分間、 $4^\circ \text{C}$ で保持。最後に、PCR反応を、Ampure XPビーズ（Beckman Coulter）を使用して、0.55ビーズ対PCR反応比で精製した。コアプロモーター候補をプラスミド上ではなく、PCR産物のコンテキスト内で試験することにより、他の交絡配列が存在しないことを保証する。

#### 【0102】

図14は、生成されたコアプロモーター（表1に示す配列）を対照コアプロモーターと比較したインビトロプロモーターアッセイを示す。図14に示すように、「no-CP」は対照、「Serpina1」は内因性コアプロモーター、「SCP1」は合成コアプロモーター、「GCP7」、「GCP10」、「GCP\_MTE」、および「GCP\_MTE\_V2」は、開示された方法に従って生成されるシャープコアプロモーター、「GCP18」は、開示された方法に従って生成されるランダムコアプロモーターである。生成されたコアプロモーター（CP）および対照の倍率変化値（平均 $\pm 95\%$ 信頼区間、CI）を示す。すべてのデータポイントは、エンハンサーなしのノーコアプロモーターコントロールの平均値に対して正規化される。コアプロモーターはタイプ別にグループ化される：対照はコアプロモーターを含まない陰性対照である；内因性は肝臓発現Serpina1遺伝子のコアプロモーターである；合成はSCP1（Kadonaga lab）などの操作されたコアプロモーターである；generated\_sharpは、内因性シャープコアプロモーターで訓練されたモデルに基づいて生成されたコアプロモーターである；generated\_randomは、内因性ランダム配列で訓練されたモデルに基づくコアプロモーターである（追加の陰性対照として機能する）。左の二つのパネルは、エンハンサーを含まないレポーター構築物に由来する倍率変化（ベースライン）を示し、右の二つのパネルは、エンハンサー（Huh-7の内因性肝エンハンサー、HEK293-HZの最小SFFVエンハンサー）を含むレポーター構築物に由来する倍率変化を示す。上の二つのパネルは、Huh-7細胞から得られたデータであり、下の二つのパネルは、HEK293細胞で行われた実験から得られたデータである。

#### 【0103】

【表 1】

表 1：生成されたコアプロモーターの配列。

ID	配列	配列番号
SERP1 NA1	CGTTGCCCCCTCTGGATCCACTGCTTAAA TACGGACGAGGACAGGGCCCTGTCTCCT CAGCTTCAGGCACCACTGACCTGGG ACAGTGAA	6
SCP1	CCCTAGGGTACTTATATAAGGGGGTGGG GGCGCGTTTCGTCCCTCAGTCGCGATCGAA CACTCGAGCCGAGCAGACGTGCCTACGG ACCGG	7
GCP7	AGCTGAAACCAACTCTTGAGCAATATAA AAGCTGCTGCCCCGGGACCCAGCGCAGAC GGCGGGCGGCGGCGGCGGCGGCGGCGGCG GCGGCGGCGCAGCGCT	8
GCP10	CGCCGTGGCTATAAAAGCACTGCACACC CCGCCAACCCAAACCCCGGCAA	9
GCP_M TE	GGGAACTGGTATAAAAGGGCCGGCGCTG GTTACCCAGTCCTTGGCGCCCCCTCGAG CCGAGCAGACGTGTCTAGTAGATCTCAC	10
GCP_M TE_V2	GGGAACTGGTATAAAAGGGCCGGCGCTC GTGGCGTTACCCAGTCCTTGGCGCCCCC TCGAGCCGAGCAGACGTGTCTAGTAGAT CTCAC	11
GCP18	GTCTCTCCAGTTGGATCAGGTAGATAAC TTTTTGAACATTTTCTTATTGGGAAGA TCTGGGTTCCATTCTGCTCTCTGGGATT GCAGGTGTGAGCCACA	12

10

20

## 【0104】

図 14 に示すアッセイは、Hu h - 7 細胞および HEK 293 - HZ 細胞のトランスフ  
ェクション、続いてルシフェラーゼアッセイを伴う。1 ウェル当たり  $1 \times 10^4$  個の Hu  
h - 7 細胞または HEK 293 - HZ 細胞を、DMEM + 10% FBS 中、96 ウェルプ  
レートに播種し、24 時間後に、Mirus TransIT - LT1 トランスフェクシ  
ョン試薬 (Mirus Bio、#MIR 2304) を使用して、0.1  $\mu$ g のレポーター  
構築物でトランスフェクトした。トランスフェクション対照として、ホタルルシフェ  
ラーゼプラスミドを 1 : 9 の比率で共トランスフェクトした (ホタルプラスミド : Nano  
Luc PCR 産物)。ルシフェラーゼ活性をアッセイするために、細胞を溶解し、Na  
no - Glo デュアルルシフェラーゼアッセイシステム (Promega 社、#N161  
0) を用いてトランスフェクションの 24 時間後に処理した。SpectraMax i  
3 プレートリーダー (Molecular Devices) を使用して、ルシフェラー  
ゼ活性を測定した。

30

40

## 【0105】

計算方法

訓練データ

ヒトゲノムの候補コアプロモーターのリストは、R パッケージ CAGER (doi : 1  
0.18129/B9.BIOC.CAGER) を使用して FANTOM5 (doi : 1  
0.1038/sdata.2017.112) から転写開始点 (TSS) プロファイリ  
ングデータとしてダウンロードするか、または FANTOM5 から直接ダウンロードした  
。このデータから、予測モデルおよび生成モデルに対して、コアプロモーターの二つの別  
個のリストを作成した (以下のモデルの説明を参照)。生成モデルについて、FANTO

50

M5データセットをフィルタリングして、ヒトの成人肝データ(サンプル: liver\_\_adult\_\_pool1)のみを保持した。データを正規化し(方法=「power Law」、fitInRange=c(5、1000)、アルファ=1.05、T=1\*10<sup>6</sup>)、TSSをクラスタリングし、四分位幅を計算した(qLow=0.1、qUp=0.9)。その後、その四分位幅に基づいて、シャープTSSおよびブロードTSSにピン化された。次に、コアプロモーター候補を、TSSサミットを5'方向に49bp、および3'方向に50bp延長することによって作成した。CAGEシグナルに基づいて、2,950の最強のコアプロモーターのみが維持された。対照セットとして、これらのコアプロモーターを5'方向に50,000bpシフトさせ、シフト後に任意のCAGEピークと重複するものをフィルタリングして除去した(n\_control=2915)。

10

#### 【0106】

予測モデルでは、FANTOM5データセット全体の上位CAGEピーク(hg19.cage\_\_peak\_\_phase1and2combined\_\_coord.bedの列5でカットオフ>50,000)を取り、生成モデルの任意のピークと重複するものをフィルタリングして除外した。生成モデルのピークと同様に、TSSサミットを、5'方向に49bp、3'方向に50bp延長して、コアプロモーターの最終リストを作成した。これらのコアプロモーターを5'方向に50,000bpシフトさせて、陰性対照領域のリストを作成し、これを生成モデルについて任意のCAGEピークまたは陰性対照領域に対してさらにフィルタリングした。上記のコアプロモーター配列の全てを、ヒトゲノムアセンブリ(hg19)中のNsを含有する任意の配列に対してフィルタリングした。予測モデルでは、コアプロモーターの各カテゴリーを、標識に関連付けられた配列として保存した(1=コアプロモーター、0=陰性対照)。

20

#### 【0107】

##### 予測モデル

予測モデルを訓練するために、標識された配列を取り(「訓練データ」のセクションを参照)、コアプロモーター生物学に関連する特徴を抽出した。GC含量、ATおよびCGジヌクレオチド頻度、ATG頻度、コアプロモーターモチーフの発生、および相対エントロピーが計算された(それぞれ、A、C、G、Tの頻度が、0.3、0.2、0.2、0.3のランダム配列と比較して)。モチーフの発生については、TSSに対する相対位置決め方式も考慮された。次に、訓練データセットを、scikit-learnのtrain\_\_test\_\_splitを使用して、訓練セットおよび検証セットに分割した(80%訓練、20%検証、ラベル階層化あり、Scikit-learn:Machine Learning in Python、Pedregosa et al., JMLR 12, pp.2825-2830, 2011, バージョン0.21.2)。上記は、19の特徴を有する11,339例の訓練セットおよび2,835の検証例をもたらした。イエローブリック(doi:10.5281/zenodo.3687330、バージョン0.9.1)に実装されたf1加重スコアを使用して、L1正則化を伴うロジスティック回帰モデルのハイパーパラメータを選択した:penalty='l1'、solver='liblinear'、multi\_\_class='auto'、C=0.5。ROCや特徴の重要性などのすべての指標は、イエローブリックを使用して可視化された。

30

40

#### 【0108】

##### 生成モデル

生成モデルを訓練するために、異なるクラスのコアプロモーター配列を、予測標的として以下のヌクレオチドを有する10ntシード配列の対に分割した。シード/標的対を生成するために、コアプロモーター配列を、工程サイズ1のスライディングウィンドウアプローチを使用して分割した。これにより、シャープCPについては94,869対、ブロードCPについては170,640対、ランダムな対照については261,720対となった。これらの配列対を、単純な数値符号化('A':0、'C':1、'G':2、'T':3)を使用してベクター化した。

#### 【0109】

50

TensorFlowバックエンド(1.13.0-dev20190126)を備えたkeras(2.2.4)を使用して、長短期メモリ(LSTM)リカレントニューラルネットワーク(RNN)を実装した。LSTMは128ユニットで、その後、ソフトマックスアクティベーション関数を有する単一の高密度出力層が続いた。RMSprop(lr=0.001)オプティマイザおよびカテゴリカルな交差エントロピーを、損失関数として使用した。このモデルは、それぞれの対(上記参照)を入力として使用し、バッチサイズ128で25エポック訓練させたが、早期停止を採用した(損失監視、忍耐=1、分\_デルタ=0.001)。新しい配列を生成するために、学習した確率分布からのサンプリングを使用して、「シード」配列の新しいヌクレオチドを予測することができる。次いで、新たに生成された配列を、配列生成の反復プロセスにおける配列生成の別のサイクルの別のシードとして使用することができる。このプロセスに確率を加えるために、最初に学習した確率を、ソフトマックス温度によって再計量してから、この新しく導出された分布からサンプリングした(温度=0.8)。最後に、このアプローチを使用して、シャープでブロードかつランダムなコアプロモーター配列からの入力を使用して訓練されたモデルに基づいて、新規のコアプロモーター配列を生成した。これらの配列の長さは、50~100ntの範囲であった。

#### 【0110】

図15は、ネットワーク1504を通じて接続された計算デバイス1501およびサーバ1502の非限定的な例を含む環境1500を描写するブロック図である。一態様では、いずれの記載の方法のいくつかまたは全ての工程も、本明細書に記載の計算デバイスで実行することができる。コンピューティングデバイス1501は、配列データ1520(例えば、CAGEデータなどのプロモーター配列データ)、訓練データ1522(例えば、標識された配列データ:コアプロモーター配列および制御配列)、生成モジュール1524(例えば、任意の補助的な訓練モジュールを含む、LSTM-RNN710)、予測モジュール1526(例えば、任意の補助的な訓練モジュールを含む、MLモジュール930)などの一つまたは複数を保存するように構成された一つまたは複数のコンピュータを備え得る。サーバ1502は、配列データ1520を保存するように構成した一つまたは複数のコンピュータを含むことができる。複数のサーバ1502は、ネットワーク1504を通じて計算デバイス1501と通信することができる。一実施形態では、サーバ1502は、CAGE実験によって生成されたデータのためのリポジトリを備えてもよい。

#### 【0111】

計算デバイス1501およびサーバ1502は、ハードウェアアーキテクチャに関して、一般にプロセッサ1508、メモリシステム1510、入力/出力(I/O)インターフェース1512、およびネットワークインターフェース1514を含む、デジタルコンピュータであってもよい。これらの構成要素(1508、1510、1512、および1514)は、ローカルインターフェース1516を介して通信的に連結される。ローカルインターフェース1516は、例えば、当該技術分野で公知の一つまたは複数のバスまたは他の有線もしくは無線接続であってもよいが、これに限定されない。ローカルインターフェース1516は、コントローラ、バッファ(キャッシュ)、ドライバ、リピータ、およびレシーバなどの、通信を可能にするための追加の要素(簡略化のために省略される)を有してもよい。さらに、ローカルインターフェースは、前述の構成要素間の適切な通信を可能にするためのアドレス、制御、および/またはデータ接続を含んでもよい。

#### 【0112】

プロセッサ1508は、特にメモリシステム1510に保存される、ソフトウェアを実行するためのハードウェアデバイスであってもよい。プロセッサ1508は、任意のカスタム作製または市販のプロセッサ、中央処理ユニット(CPU)、計算デバイス1501およびサーバ1502に関連付けられたいくつかのプロセッサの中の補助プロセッサ、半導体ベースのマイクロプロセッサ(マイクロチップもしくはチップセットの形態)、またはソフトウェア命令を実行するための一般に任意のデバイスとすることができる。計算デバイス1501および/またはサーバ1502が動作中である時、プロセッサ1508は

、メモリシステム 1510 内に保存されているソフトウェアを実行して、メモリシステム 1510 へのおよびそこからのデータを通信し、ソフトウェアに従って、計算デバイス 1501 およびサーバ 1502 の動作を一般に制御するように構成されてもよい。

【0113】

I/O インターフェース 1512 を使用して、一つまたは複数のデバイスまたは構成要素からユーザ入力を受信する、かつ/またはそれらへとシステム出力を提供することができる。ユーザ入力は、例えば、キーボードおよび/またはマウスを介して提供されてもよい。システム出力は、表示デバイスおよびプリンタ（図示せず）を介して提供されてもよい。I/O インターフェース 1512 は、例えば、シリアルポート、パラレルポート、小型コンピュータシステムインターフェース（SCSI）、赤外（IR）インターフェース、無線周波数（RF）インターフェース、および/またはユニバーサルシリアルバス（USB）インターフェースを含んでもよい。

10

【0114】

ネットワークインターフェース 1514 は、計算デバイス 1501 および/またはネットワーク 1504 上のサーバ 1502 から送信および受信するために使用することができる。ネットワークインターフェース 1514 は、例えば、10BaseT Ethernet アダプタ、100BaseT Ethernet アダプタ、LAN PHY Ethernet アダプタ、Token Ring アダプタ、ワイヤレスネットワークアダプタ（例えば、WiFi、セルラー、サテライト）、または任意の他の好適なネットワークインターフェースデバイスを含んでもよい。ネットワークインターフェース 1514 は、ネットワーク 1504 上での適切な通信を可能にするためのアドレス、制御、および/またはデータ接続を含んでもよい。

20

【0115】

メモリシステム 1510 は、揮発性メモリ素子（例えば、ランダムアクセスメモリ（DRAM、SRAM、SDRAM などの RAM））および不揮発性メモリ素子（例えば、ROM、ハードドライブ、テープ、CDROM、DVDROM など）のいずれか一つまたはその組み合わせを含んでもよい。さらに、メモリシステム 1510 は、電子、磁気、光学、および/または他の型の保存媒体を組み込んでもよい。メモリシステム 1510 は、様々な構成要素が互いに離れて位置するが、プロセッサ 1508 によってアクセスすることができる、分散型アーキテクチャを有し得ることに留意されたい。

30

【0116】

メモリシステム 1510 内のソフトウェアは、一つまたは複数のソフトウェアプログラムを含んでもよく、これらの各々は、論理機能を実施するための実行可能な命令の順序付けされたリストを含む。図 15 の例では、コンピューティングデバイス 1501 のメモリシステム 1510 におけるソフトウェアは、配列データ 1520、訓練データ 1522、生成モジュール 1524、予測モジュール 1526、および適当な操作システム（OS）1518 を含むことができる。図 15 の例では、サーバ 1502 のメモリシステム 1510 内のソフトウェアは、配列データ 1520、および好適なオペレーティングシステム（OS）1518 を含むことができる。オペレーティングシステム 1518 は、他のコンピュータプログラムの実行を本質的に制御し、スケジューリング、入力 - 出力制御、ファイルおよびデータ管理、メモリ管理、および通信制御、ならびに関連するサービスを提供する。

40

【0117】

例証の目的で、アプリケーションプログラムおよびオペレーティングシステム 1518 などの他の実行可能なプログラム構成要素は、本明細書では別々のブロックとして例証されているが、そのようなプログラムおよび構成要素は、計算デバイス 1501 および/またはサーバ 1502 の異なる保存構成要素内で、様々な時間に存在し得ることが認識される。生成モジュール 1524 および/または予測モジュール 1526 の実施は、何らかの形態のコンピュータ可読媒体に保存されても、それを通過して伝達されてもよい。本開示の方法のいずれも、コンピュータ可読媒体上に具現化されたコンピュータ可読命令によっ

50

て実行することができる。コンピュータ可読媒体は、コンピュータによってアクセス可能な任意の利用可能媒体とすることができる。例として、かつ限定を意図するものではないが、コンピュータ可読媒体は、「コンピュータストレージ媒体」および「通信媒体」を含み得る。「コンピュータ記憶媒体」は、コンピュータ可読命令、データ構造、プログラムモジュール、または他のデータなどの、情報を記憶するための任意の方法または技術で実施される、揮発性および不揮発性の取り外し可能な媒体および取り外し不能な媒体を含み得る。例示的なコンピュータ記憶媒体は、RAM、ROM、EEPROM、フラッシュメモリもしくは他の記憶技術、CD-ROM、デジタル多用途ディスク(DVD)もしくは他の光学記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶デバイスもしくは他の磁気記憶デバイス、または所望の情報の記憶に使用することができ、かつコンピュータによってアクセスすることができる任意の他の媒体を含み得る。

10

**【0118】**

一実施形態では、予測モジュール1526は、図16に示す方法1600を実行するように構成され得る。方法1600は、単一の計算デバイス、複数の電子デバイス、および同様のものによって、全体的または部分的に実施されてもよい。方法1600は、遺伝的データを受信することを含み、遺伝的データは、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列は、1601で関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む。関連発現スコアは、遺伝子発現のキャップ解析(CAGE)ピークを含み得る。

**【0119】**

20

方法1600は、閾値を満たす関連発現スコアに基づいて、1602で第一の複数のヌクレオチド配列から複数のTSSを決定することを含み得る。

方法1600は、複数のTSSに基づいて、1603で複数のサミットヌクレオチド塩基を決定することを含み得る。複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することは、複数のTSSの各々について、最強のCAGEシグナルを有するヌクレオチド塩基を決定することを含み得る。

**【0120】**

方法1600は、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、1604で関連する複数の周辺塩基を決定することを含み得る。複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することは、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定することを含み得る。5'方向の第一の複数のヌクレオチド塩基が49個のヌクレオチド塩基を含み得、3'方向の第二の複数のヌクレオチド塩基が50個のヌクレオチド塩基を含む。

30

**【0121】**

方法1600は、1605でコアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することを含み得る。

**【0122】**

40

方法1600は、第二の複数のヌクレオチド配列の各ヌクレオチド配列について、1606で関連する複数のシフト塩基を決定することを含み得る。第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することは、第二の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含み得る。

**【0123】**

方法1600は、1607でコアプロモーターではないとして標識された第三の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することを含み得る。

方法1600は、コアプロモーターとして標識された第二の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づい

50



て、1608で訓練データセットを生成することを含み得る。

【0124】

方法1600は、訓練データセットに基づいて、1609で予測モデルについての複数の特徴を決定することを含み得る。予測モデルについての複数の特徴は、GC含量、ATおよびCGジヌクレオチド頻度、ATG頻度、コアプロモーターモチーフの発生、相対エントロピー、および関連するTSSに対する相対的位置決め方式のうちの一つまたは複数を含んでもよい。

【0125】

1610では、方法1600は、訓練データセットの第一の部分に基づいて、1610で複数の特徴に従って予測モデルを訓練することを含んでもよい。

10

方法1600は、訓練データセットの第二の部分に基づいて、1611で予測モデルを試験することを含み得る。

【0126】

方法1600は、試験に基づいて、1612で予測モデルを出力することを含み得る。

一実施形態では、方法1600は、生成モデルで使用されるTSSの発現スコアと重複する発現スコアを有する第一の複数のヌクレオチド配列から、複数のTSSの任意のTSSをフィルタリングすることをさらに含むことができる。一実施形態では、方法1600は、ヒトゲノムアセンブリ(hg19)中にNsを含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることをさらに含むことができる。

【0127】

20

一実施形態では、予測モジュール1526は、図17に示す方法1700を実行するように構成され得る。方法1700は、単一の計算デバイス、複数の電子デバイス、および同様のものによって、全体的または部分的に実施されてもよい。方法1700は、遺伝的データを受信することを含み、遺伝的データは、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列は、1701で関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む。関連発現スコアは、遺伝子発現のキャップ解析(CAGE)ピークを含んでもよい。

【0128】

方法1700は、第一の複数のヌクレオチド配列に基づいて、1702でコアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することを含み得る。第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することは、閾値を満たす関連発現スコアに基づいて、第一の複数のヌクレオチド配列から複数のTSSを決定することと、複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することと、各サミットヌクレオチド塩基および関連する複数の周辺塩基を、コアプロモーターとして標識された第二の複数のヌクレオチド配列として保存することと、を含み得る。

30

【0129】

複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することは、複数のTSSのそれぞれについて、最強のCAGEシグナルを有するヌクレオチド塩基を決定することを含む場合がある。複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することは、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定することを含み得る。5'方向の第一の複数のヌクレオチド塩基は、49個のヌクレオチド塩基を含み得、3'方向の第二の複数のヌクレオチド塩基は、50個のヌクレオチド塩基を含む。

40

【0130】

方法1700は、第二の複数のヌクレオチド配列に基づいて、1703でコアプロモーターではないとして標識された第三の複数のヌクレオチド配列を決定することを含み得る。第二の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識され

50

た第三の複数のヌクレオチド配列を決定することは、第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することと、各関連する複数のシフト塩基を、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列として保存することと、を含み得る。

【0131】

第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することは、第二の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含み得る。

【0132】

方法1700は、コアプロモーターとして標識された第二の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、1704で訓練データセットを生成することを含み得る。

10

【0133】

方法1700は、訓練データセットに基づいて、1705で予測モデルについての複数の特徴を決定することを含み得る。予測モデルについての複数の特徴は、GC含量、ATおよびCGジヌクレオチド頻度、ATG頻度、コアプロモーターモチーフの発生、相対エントロピー、および関連するTSSに対する相対的位置決め方式のうちの一つまたは複数を含んでもよい。

【0134】

方法1700は、訓練データセットの第一の部分に基づいて、1706で複数の特徴に従って予測モデルを訓練することを含み得る。

20

方法1700は、訓練データセットの第二の部分に基づいて、1707で予測モデルを試験することを含み得る。

【0135】

方法1700は、試験に基づいて、1708で予測モデルを出力することを含み得る。

一実施形態では、コラム1700の方法はまた、生成モデルで使用されるTSSの発現スコアと重複する発現スコアを有する第一の複数のヌクレオチド配列から、複数のTSSの任意のTSSをフィルタリングすることを含んでもよい。一実施形態では、方法1700は、ヒトゲノムアセンブリ(hg19)中にNsを含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることをさらに含むことができる。

30

【0136】

一実施形態では、生成モジュール1524は、図18に示す方法1800を実施するように構成され得る。方法1800は、単一の計算デバイス、複数の電子デバイス、および同様のものによって、全体的または部分的に実施されてもよい。方法1800は、遺伝的データを受信することを含み、遺伝的データは、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列は、1801で関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む。関連発現スコアは、遺伝子発現のキャップ解析(CAGE)ピークを含んでもよい。

【0137】

方法1800は、1802で遺伝的データを正規化することを含み得る。

40

方法1800は、関連発現スコアに基づいて、1803でTSSをクラスタリングすることを含み得る。

【0138】

方法1800は、TSSの各クラスターについて、1804で四分位幅を決定することを含み得る。

方法1800は、四分位幅に基づいて、1805で各TSSをシャープTSSまたはブロードTSSとして標識することを含み得る。

【0139】

方法1800は、複数のTSSに基づいて、1806で複数のサミットヌクレオチド塩基を決定することを含み得る。複数のTSSに基づいて、複数のサミットヌクレオチド塩

50

基を決定することは、複数の T S S の各々について、最強の C A G E シグナルを有するヌクレオチド塩基を決定することを含み得る。

【 0 1 4 0 】

方法 1 8 0 0 は、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、1 8 0 7 で関連する複数の周辺塩基を決定することを含み得る。複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することは、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5 ' 方向の第一の複数のヌクレオチド塩基および 3 ' 方向の第二の複数のヌクレオチド塩基を決定することを含み得る。5 ' 方向の第一の複数のヌクレオチド塩基は、4 9 個のヌクレオチド塩基を含み、3 ' 方向の第二の複数のヌクレオチド塩基は、5 0 個のヌクレオチド塩基を含む。

10

【 0 1 4 1 】

方法 1 8 0 0 は、1 8 0 8 でコアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することを含み得る。

【 0 1 4 2 】

方法 1 8 0 0 は、閾値を満たす関連発現スコアに基づいて、1 8 0 9 で第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することを含み得る。

方法 1 8 0 0 は、第三の複数のヌクレオチド配列の各ヌクレオチド配列について、1 8 1 0 で関連する複数のシフト塩基を決定することを含み得る。第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することは、第三の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含み得る。

20

【 0 1 4 3 】

方法 1 8 0 0 は、1 8 1 1 でコアプロモーターではないとして標識された第四の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することを含み得る。

方法 1 8 0 0 は、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、1 8 1 2 で訓練データセットを生成することを含み得る。

【 0 1 4 4 】

30

方法 1 8 0 0 は、訓練データセットの各ヌクレオチド配列について、1 8 1 3 で複数のシード配列および標的ヌクレオチド対を生成することを含み得る。訓練データセット中の各ヌクレオチド配列について、複数のシード配列および標的ヌクレオチド対を生成することは、シャープ T S S またはブロード T S S 標識に基づいて、訓練データセット中のヌクレオチド配列をシャープ T S S 群またはブロード T S S 群に分割することと、定義された長さのスライディングウィンドウを適用し、定義された工程サイズを各ヌクレオチド配列に有ることと、スライディングウィンドウの各工程でシード配列および標的ヌクレオチド対を保存することと、を含み得る。

【 0 1 4 5 】

方法 1 8 0 0 は、1 8 1 4 で、複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクター化することを含み得る。各シード配列および標的ヌクレオチド対は、定義された長さを有するシード配列、および所定のヌクレオチド配列上のシード配列の直後に標的ヌクレオチドを含み得る。定義された長さは、例えば、1 0 塩基であってもよい。複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクター化することは、各ヌクレオチドをそれぞれの番号としてコード化することを含み得る。

40

【 0 1 4 6 】

方法 1 8 0 0 は、ベクター化されたシード配列および標的ヌクレオチド対に基づいて、1 8 1 5 で生成モデルを訓練することを含み得る。生成モデルは、長短期メモリ ( L S T M ) リカレントニューラルネットワーク ( R N N ) を含んでもよい。

50

## 【0147】

方法1800は、1816で生成モデルを出力することを含んでもよい。

一実施形態では、方法1800は、ヒトゲノムアセンブリ(hg19)中にNsを含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることを含み得る。

## 【0148】

一実施形態では、方法1800は、生成モデルに基づいて、ヌクレオチド配列を生成することを含み得る。ヌクレオチド配列は、例えば、コアプロモーター配列であってもよい。方法1800は、コアプロモーター配列に基づいてプロモーターを操作することを含み得る。

10

## 【0149】

生成モデルに基づいて、ヌクレオチド配列を生成することは、(a)シード配列を受信することと、(b)シード配列に基づいて、次のヌクレオチドを予測することと、(c)シード配列に次のヌクレオチドを付加することと、(d)ヌクレオチド配列の所望の長さに達するまでb~cを繰り返すことと、を含み得る。所望の長さは、例えば、約50ヌクレオチド~約100ヌクレオチドであってもよい。

## 【0150】

一実施形態では、方法1800は、プロモーターを核酸構築物に挿入することを含み得る。プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含み得る。

20

## 【0151】

一実施形態では、方法1800は、核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することを含んでもよい。

一実施形態では、生成モジュール1524は、図19に示す方法1900を実施するように構成され得る。方法1900は、単一の計算デバイス、複数の電子デバイス、および同様のものによって、全体的または部分的に実施されてもよい。方法1900は、遺伝的データを受信することを含み、遺伝的データは、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列は、1901で関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む。関連発現スコアは、遺伝子発現のキャップ解析(CAGE)ピークを含んでもよい。

30

## 【0152】

方法1900は、第一の複数のヌクレオチド配列に基づいて、1902でコアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することを含み得る。第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することは、複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することと、各サミットヌクレオチド塩基および関連する複数の周辺塩基を、コアプロモーターとして標識された第二の複数のヌクレオチド配列として保存することと、を含み得る。

## 【0153】

40

複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することは、複数のTSSのそれぞれについて、最強のCAGEシグナルを有するヌクレオチド塩基を決定することを含む。複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することは、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定することを含み得る。5'方向の第一の複数のヌクレオチド塩基が49個のヌクレオチド塩基を含み得、3'方向の第二の複数のヌクレオチド塩基が50個のヌクレオチド塩基を含む。

## 【0154】

方法1900は、閾値を満たす関連発現スコアに基づいて、1903で第二の複数のヌ

50

クレオチド配列から第三の複数のヌクレオチド配列を決定することを含み得る。

方法 1900 は、第三の複数のヌクレオチド配列に基づいて、1904 でコアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することを含み得る。第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することは、第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することと、各関連する複数のシフト塩基を、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列として保存することと、を含み得る。

【0155】

第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することは、第三の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含み得る。

【0156】

方法 1900 は、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、1905 で訓練データセットを生成することを含み得る。

【0157】

方法 1900 は、訓練データセットに基づいて、1906 で生成モデルを訓練することを含んでもよい。生成モデルは、例えば、長短期メモリ (LSTM) リカレントニューラルネットワーク (RNN) を含んでもよい。訓練データセットに基づいて、生成モデルを訓練することは、訓練データセット中の各ヌクレオチド配列について、複数のシード配列と標的ヌクレオチド対を生成することと、複数のシード配列と標的ヌクレオチド対の各シード配列と標的ヌクレオチド対をベクター化することと、ベクター化されたシード配列と標的ヌクレオチド対に基づいて、生成モデルを訓練することと、を含み得る。

【0158】

各シード配列および標的ヌクレオチド対は、定義された長さを有するシード配列、および所定のヌクレオチド配列上のシード配列の直後に標的ヌクレオチドを含み得る。定義された長さは、例えば、10 塩基であってもよい。複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクター化することは、各ヌクレオチドをそれぞれの番号としてコード化することを含み得る。

【0159】

方法 1900 は、1907 で生成モデルを出力することを含んでもよい。

一実施形態では、方法 1900 は、遺伝的データを正規化することを含み得る。一実施形態では、方法 1900 は、関連発現スコアに基づいて、TSS をクラスタリングすることと、TSS の各クラスターについて、四分位幅を決定することと、四分位幅に基づいて、各 TSS をシャープ TSS またはブロード TSS として標識することと、を含み得る。

【0160】

一実施形態では、方法 1900 は、生成モデルに基づいて、ヌクレオチド配列を生成することを含み得る。ヌクレオチド配列は、例えば、コアプロモーター配列であってもよい。

【0161】

一実施形態では、方法 1900 は、コアプロモーター配列に基づいてプロモーターを操作することを含み得る。方法 1900 は、プロモーターを核酸構築物に挿入することを含み得る。プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含み得る。

【0162】

一実施形態では、方法 1900 は、核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することを含んでもよい。

一実施形態では、方法 1900 は、ヒトゲノムアセンブリ (hg19) 中に Ns を含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることを含み得る。

10

20

30

40

50

## 【 0 1 6 3 】

一実施形態では、方法 1 9 0 0 は、生成モデルに基づいて、ヌクレオチド配列を生成することを含み得る。ヌクレオチド配列は、例えば、コアプロモーター配列であってもよい。方法 1 9 0 0 は、コアプロモーター配列に基づいてプロモーターを操作することを含み得る。

## 【 0 1 6 4 】

生成モデルに基づいて、ヌクレオチド配列を生成することは、( a ) シード配列を受信することと、( b ) シード配列に基づいて、次のヌクレオチドを予測することと、( c ) シード配列に次のヌクレオチドを付加することと、( d ) ヌクレオチド配列の所望の長さに達するまで b ~ c を繰り返すことと、を含み得る。所望の長さは、例えば、約 5 0 ヌクレオチド ~ 約 1 0 0 ヌクレオチドであってもよい。

10

## 【 0 1 6 5 】

一実施形態では、方法 1 9 0 0 は、プロモーターを核酸構築物に挿入することを含み得る。プロモーターを核酸構築物に挿入することは、導入遺伝子の上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含み得る。

## 【 0 1 6 6 】

一実施形態では、方法 1 9 0 0 は、核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することを含んでもよい。

一実施形態では、予測モジュール 1 5 2 6 は、図 2 0 に示す方法 2 0 0 0 を実行するように構成され得る。方法 2 0 0 0 は、単一の計算デバイス、複数の電子デバイス、および同様のものによって、全体的または部分的に実施されてもよい。方法 2 0 0 0 は、2 0 1 0 でヌクレオチド配列を受信することを含んでもよい。ヌクレオチド配列を受信することは、複数のヌクレオチド配列を受信することを含み得、複数のヌクレオチド配列は、生成モデルによって生成された。

20

## 【 0 1 6 7 】

方法 2 0 0 0 は、訓練された予測モデルに、2 0 2 0 でヌクレオチド配列を提供することを含み得る。

方法 2 0 0 0 は、予測モデルに基づいて、2 0 3 0 でヌクレオチド配列がコアプロモーターであることを決定することを含み得る。

## 【 0 1 6 8 】

一実施形態では、方法 2 0 0 0 は、ヌクレオチド配列がコアプロモーターであるという決定に基づいて、一つまたは複数の基準に従ってヌクレオチド配列をフィルタリングすることを含み得る。一つまたは複数の基準は、例えば、G C 含量またはモチーフのうちの一つまたは複数を含んでもよい。

30

## 【 0 1 6 9 】

一実施形態では、方法 2 0 0 0 は、遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点 ( T S S ) を含む、遺伝的データを受信することと、閾値を満たす関連発現スコアに基づいて、第一の複数のヌクレオチド配列から複数の T S S を決定することと、複数の T S S に基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基に対して、関連する複数の周辺塩基を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することと、第二の複数のヌクレオチド配列の各ヌクレオチド配列に対して、関連する複数のシフト塩基を決定することと、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、予測モデルに対する複数の特徴を決定することと、訓練データセットの第一の部分に基づいて、複数の特徴に従って予測モデルを訓練す

40

50

ることと、訓練データセットの第二の部分に基づいて、予測モデルを試験することと、および試験に基づいて、予測モデルを出力することと、を含み得る。

【0170】

一実施形態では、生成モジュール1524は、図21に示す方法2100を実施するように構成され得る。方法2100は、単一の計算デバイス、複数の電子デバイス、および同様のものによって、全体的または部分的に実施されてもよい。方法2100は、2110でヌクレオチド配列および配列長を受信することを含んでもよい。

【0171】

方法2100は、訓練された生成モデルに、2120でヌクレオチド配列を提供することを含み得る。

方法2100は、生成モデルに基づいて、2130でヌクレオチド配列と関連付けられた次のヌクレオチドを決定することを含み得る。

【0172】

方法2100は、2140で、次のヌクレオチドをヌクレオチド配列に付加することを含んでもよい。

方法2100は、2150でヌクレオチド配列の長さが配列長と等しくなるまで、工程2120~2140を繰り返すことを含み得る。配列長は、例えば、約50ヌクレオチド~約100ヌクレオチドであってもよい。

【0173】

方法2100は、2160でコアプロモーター配列としてヌクレオチド配列を出力することを含み得る。

一実施形態では、方法2100は、コアプロモーター配列に基づいてプロモーターを操作することを含み得る。

【0174】

一実施形態では、方法2100は、プロモーターを核酸構築物に挿入することを含み得る。プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含み得る。

【0175】

一実施形態では、方法2100は、核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することを含んでもよい。

一実施形態では、方法2100は、遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、および訓練データセットに基づいて、生成モデルを訓練することと、を含み得る。

【0176】

一実施形態では、方法2100は、一つまたは複数の基準に従ってヌクレオチド配列をフィルタリングすることを含み得る。一つまたは複数の基準は、GC含量またはモチーフのうちの一つまたは複数を含んでもよい。

【0177】

記載された方法、システム、および装置、ならびにそれらの変形に照らして、本明細書では、以下に本発明のより具体的に記述された特定の実施形態を説明する。しかし、これらの特に列挙された実施形態は、本明細書に記載される異なるまたはより一般的な教示を含む任意の異なる特許請求の範囲に対して何らかの限定効果を有すると解釈されるべきで

10

20

30

40

50

はなく、または「特定の」実施形態が、その中に文字通り使用される言語の固有の意味以外の何らかの方法で、何らかの形で限定されると解釈されるべきでもない。

【0178】

実施形態1：

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点(TSS)を含む、遺伝的データを受信することと、閾値を満たす関連発現スコアに基づいて、第一の複数のヌクレオチド配列から複数のTSSを決定することと、複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基に対して、関連する複数の周辺塩基を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することと、第二の複数のヌクレオチド配列の各ヌクレオチド配列に対して、関連する複数のシフト塩基を決定することと、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、予測モデルに対する複数の特徴を決定することと、訓練データセットの第一の部分に基づいて、複数の特徴に従って予測モデルを訓練することと、訓練データセットの第二の部分に基づいて、予測モデルを試験することと、および試験に基づいて、予測モデルを出力することと、を含む方法。

10

20

【0179】

実施形態2：

関連発現スコアが、遺伝子発現のキャップ解析(CAGE)ピークを含む、先行する実施形態のいずれか一つに記載の実施形態。

【0180】

実施形態3：

複数のTSSに基づいて、複数のサミットヌクレオチド塩基を決定することが、複数のTSSのそれぞれについて、最強のCAGEシグナルを有するヌクレオチド塩基を決定することを含む、先行する実施形態のいずれか一つに記載の実施形態。

30

【0181】

実施形態4：

複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することが、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定することを含む、先行する実施形態のいずれか一つに記載の実施形態。

【0182】

実施形態5：

5'方向の第一の複数のヌクレオチド塩基が49個のヌクレオチド塩基を含み、3'方向の第二の複数のヌクレオチド塩基が50個のヌクレオチド塩基を含む、実施形態4に記載の実施形態。

40

【0183】

実施形態6：

第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することが、第二の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含む、先行する実施形態のいずれか一つに記載の実施形態。

【0184】

実施形態7：

50



予測モデルについての複数の特徴が、G C 含量、A T および C G ジヌクレオチド頻度、A T G 頻度、コアプロモーターモチーフの発生、相対エントロピー、および関連する T S S に対する相対的位置決め方式のうちの一つまたは複数の含む、先行する実施形態のいずれか一つに記載の実施形態。

【0185】

実施形態 8 :

生成モデルで使用される T S S の発現スコアと重複する発現スコアを有する第一の複数のヌクレオチド配列から、複数の T S S の任意の T S S をフィルタリングすることをさらに含む、先行する実施形態のいずれか一つに記載の実施形態。

【0186】

実施形態 9 :

ヒトゲノムアセンブリ (hg19) 中に N s を含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることをさらに含む、先行する実施形態のいずれか一つに記載の実施形態。

【0187】

実施形態 10 :

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点 (T S S) を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、第二の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、予測モデルに対する複数の特徴を決定することと、訓練データセットの第一の部分に基づいて、複数の特徴に従って予測モデルを訓練することと、訓練データセットの第二の部分に基づいて、予測モデルを試験することと、および試験に基づいて、予測モデルを出力することと、を含む方法。

【0188】

実施形態 11 :

第二の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列を決定することが、第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することと、各関連する複数のシフト塩基を、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列として保存することと、を含む、実施形態 10 に記載の実施形態。

【0189】

実施形態 12 :

第二の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することが、第二の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含む、実施形態 11 に記載の実施形態。

【0190】

実施形態 13 :

関連発現スコアが、遺伝子発現のキャップ解析 (CAGE) ピークを含む、実施形態 10 ~ 12 のいずれかに記載の実施形態。

【0191】

実施形態 14 :

第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することが、閾値を満たす関連発現スコアに基づいて、第一の複数のヌクレオチド配列から複数の T S S を決定することと、複数の T S S に基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩

10

20

30

40

50

基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することと、各サミットヌクレオチド塩基および関連する複数の周辺塩基を、コアプロモーターとして標識された第二の複数のヌクレオチド配列として保存することと、を含む、実施形態 10 ~ 13 のいずれかに記載の実施形態。

【0192】

実施形態 15 :

複数の TSS に基づいて、複数のサミットヌクレオチド塩基を決定することが、複数の TSS のそれぞれについて、最強の CAGE シグナルを有するヌクレオチド塩基を決定することを含む、実施形態 14 に記載の実施形態。

【0193】

実施形態 16 :

複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することが、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5' 方向の第一の複数のヌクレオチド塩基および 3' 方向の第二の複数のヌクレオチド塩基を決定することを含む、実施形態 14 に記載の実施形態。

【0194】

実施形態 17 :

5' 方向の第一の複数のヌクレオチド塩基が 49 個のヌクレオチド塩基を含み、3' 方向の第二の複数のヌクレオチド塩基が 50 個のヌクレオチド塩基を含む、実施形態 16 に記載の実施形態。

【0195】

実施形態 18 :

予測モデルについての複数の特徴が、GC 含量、AT および CG ジヌクレオチド頻度、ATG 頻度、コアプロモーターモチーフの発生、相対エントロピー、および関連する TSS に対する相対的位置決め方式のうちの一つまたは複数を含む、実施形態 10 ~ 17 のいずれかに記載の実施形態。

【0196】

実施形態 19 :

生成モデルで使用される TSS の発現スコアと重複する発現スコアを有する第一の複数のヌクレオチド配列から、複数の TSS の任意の TSS をフィルタリングすることをさらに含む、実施形態 14 ~ 18 のいずれかに記載の実施形態。

【0197】

実施形態 20 :

ヒトゲノムアセンブリ (hg19) 中に Ns を含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることをさらに含む、実施形態 10 ~ 19 のいずれかに記載の実施形態。

【0198】

実施形態 21 :

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点 (TSS) を含む、遺伝的データを受信することと、遺伝的データを正規化することと、関連発現スコアに基づいて、TSS をクラスター化することと、TSS の各クラスターについて、四分位幅を決定することと、四分位幅に基づいて、各 TSS をシャープ TSS またはブロード TSS に標識することと、複数の TSS に基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基に対して、関連する複数の周辺塩基を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することと、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列の各ヌクレオチド配列に対して、関連する複数のシフト塩基を決定することと、

10

20

30

40

50

コアプロモーターではないとして標識された第四の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットにおける各ヌクレオチド配列に対して、複数のシード配列および標的ヌクレオチド対を生成することと、複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクトル化することと、ベクトル化されたシード配列および標的ヌクレオチド対に基づいて、生成モデルを訓練することと、および生成モデルを出力することと、を含む方法。

【 0 1 9 9 】

10

実施形態 2 2 :

関連発現スコアが、遺伝子発現のキャップ解析 ( C A G E ) ピークを含む、実施形態 2 1 に記載の実施形態。

【 0 2 0 0 】

実施形態 2 3 :

複数の T S S に基づいて、複数のサミットヌクレオチド塩基を決定することが、複数の T S S のそれぞれについて、最強の C A G E シグナルを有するヌクレオチド塩基を決定することを含む、実施形態 2 1 ~ 2 2 のいずれかに記載の実施形態。

【 0 2 0 1 】

実施形態 2 4 :

20

複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することが、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5'方向の第一の複数のヌクレオチド塩基および3'方向の第二の複数のヌクレオチド塩基を決定することを含む、実施形態 2 1 ~ 2 3 のいずれかに記載の実施形態。

【 0 2 0 2 】

実施形態 2 5 :

5'方向の第一の複数のヌクレオチド塩基が49個のヌクレオチド塩基を含み、3'方向の第二の複数のヌクレオチド塩基が50個のヌクレオチド塩基を含む、実施形態 2 4 に記載の実施形態。

30

【 0 2 0 3 】

実施形態 2 6 :

第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することが、第三の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含む、実施形態 2 1 ~ 2 5 のいずれかに記載の実施形態。

【 0 2 0 4 】

実施形態 2 7 :

ヒトゲノムアセンブリ ( h g 1 9 ) 中に N s を含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることをさらに含む、実施形態 2 1 ~ 2 6 のいずれかに記載の実施形態。

40

【 0 2 0 5 】

実施形態 2 8 :

各シード配列および標的ヌクレオチド対は、定義された長さを有するシード配列、および所定のヌクレオチド配列上のシード配列の直後に標的ヌクレオチドを含む、実施形態 2 1 ~ 2 7 のいずれかに記載の実施形態。

【 0 2 0 6 】

実施形態 2 9 :

定義された長さが10塩基である、実施形態 2 8 に記載の実施形態。

実施形態 3 0 :

50

訓練データセット中の各ヌクレオチド配列について、複数のシード配列および標的ヌクレオチド対を生成することは、シャープTSSまたはブロードTSS標識に基づいて、訓練データセット中のヌクレオチド配列をシャープTSS群またはブロードTSS群に分割することと、定義された長さのスライディングウィンドウを適用し、定義された工程サイズを各ヌクレオチド配列に有ることと、スライディングウィンドウの各工程でシード配列および標的ヌクレオチド対を保存することと、を含む、実施形態21~29のいずれかに記載の実施形態。

【0207】

実施形態31：

複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクター化することは、各ヌクレオチドをそれぞれの番号としてコード化することを含む、実施形態21~30のいずれかに記載の実施形態。

10

【0208】

実施形態32：

生成モデルが、長短期メモリ(LSTM)リカレントニューラルネットワーク(RNN)を含む、実施形態21~31のいずれかに記載の実施形態。

【0209】

実施形態33：

生成モデルに基づいて、ヌクレオチド配列を生成することをさらに含む、実施形態21~32のいずれかに記載の実施形態。

20

【0210】

実施形態34：

生成モデルに基づいて、ヌクレオチド配列を生成することは、(a)シード配列を受信することと、(b)シード配列に基づいて、次のヌクレオチドを予測することと、(c)シード配列に次のヌクレオチドを付加することと、(d)ヌクレオチド配列の所望の長さに達するまでb~cを繰り返すことと、を含む、実施形態33に記載の実施形態。

【0211】

実施形態35：

所望の長さは、約50ヌクレオチド~約100ヌクレオチドである、実施形態34に記載の実施形態。

30

【0212】

実施形態36：

ヌクレオチド配列がコアプロモーター配列である、実施形態33~35のいずれかに記載の実施形態。

【0213】

実施形態37：

コアプロモーター配列に基づいてプロモーターを操作することをさらに含む、実施形態36に記載の実施形態。

【0214】

実施形態38：

プロモーターを核酸構築物に挿入することをさらに含む、実施形態37に記載の実施形態。

40

【0215】

実施形態39：

プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含む、実施形態38に記載の実施形態。

【0216】

実施形態40：

核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することをさらに

50

含む、実施形態 38 ~ 39 のいずれかに記載の実施形態。

【0217】

実施形態 41 :

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点 (TSS) を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、生成モデルを訓練することと、および生成モデルを出力することと、を含む、方法。

10

【0218】

実施形態 42 :

遺伝的データを正規化することをさらに含む、実施形態 41 に記載の実施形態。

実施形態 43 :

関連発現スコアに基づいて、TSS をクラスタリングすることと、TSS の各クラスターについて、四分位幅を決定することと、四分位幅に基づいて、各 TSS をシャープ TSS またはブロード TSS として標識することと、をさらに含む、実施形態 41 ~ 42 のいずれかに記載の実施形態。

20

【0219】

実施形態 44 :

第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することが、第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することと、各関連する複数のシフト塩基を、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列として保存することと、を含む、実施形態 41 ~ 43 のいずれかに記載の実施形態。

【0220】

実施形態 45 :

訓練データセットに基づいて、生成モデルを訓練することは、訓練データセット中の各ヌクレオチド配列について、複数のシード配列と標的ヌクレオチド対を生成することと、複数のシード配列と標的ヌクレオチド対の各シード配列と標的ヌクレオチド対をベクター化することと、ベクター化されたシード配列と標的ヌクレオチド対に基づいて、生成モデルを訓練することと、を含む、実施形態 41 ~ 44 のいずれかに記載の実施形態。

30

【0221】

実施形態 46 :

関連発現スコアが、遺伝子発現のキャップ解析 (CAGE) ピークを含む、実施形態 41 ~ 45 のいずれかに記載の実施形態。

40

【0222】

実施形態 47 :

第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することは、複数の TSS に基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することと、各サミットヌクレオチド塩基および関連する複数の周辺塩基を、コアプロモーターとして標識された第二の複数のヌクレオチド配列として保存することと、を含む、実施形態 41 ~ 46 のいずれかに記載の実施形態。

【0223】

50

**実施形態 48 :**

複数の T S S に基づいて、複数のサミットヌクレオチド塩基を決定することが、複数の T S S のそれぞれについて、最強の C A G E シグナルを有するヌクレオチド塩基を決定することを含む、実施形態 47 に記載の実施形態。

**【0224】****実施形態 49 :**

複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、関連する複数の周辺塩基を決定することが、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基について、5' 方向の第一の複数のヌクレオチド塩基および 3' 方向の第二の複数のヌクレオチド塩基を決定することを含む、実施形態 47 ~ 48 のいずれかに記載の実施形態。

10

**【0225】****実施形態 50 :**

5' 方向の第一の複数のヌクレオチド塩基が 49 個のヌクレオチド塩基を含み、3' 方向の第二の複数のヌクレオチド塩基が 50 個のヌクレオチド塩基を含む、実施形態 49 に記載の実施形態。

**【0226】****実施形態 51 :**

第三の複数のヌクレオチド配列の各ヌクレオチド配列について、関連する複数のシフト塩基を決定することが、第三の複数のヌクレオチド配列の各ヌクレオチド配列から、ある量のヌクレオチド塩基をシフトさせることを含む、実施形態 41 ~ 50 のいずれかに記載の実施形態。

20

**【0227】****実施形態 52 :**

ヒトゲノムアセンブリ ( h g 19 ) 中に N s を含有する第二の複数のヌクレオチド配列の任意のヌクレオチド配列をフィルタリングすることをさらに含む、実施形態 41 ~ 51 のいずれかに記載の実施形態。

**【0228】****実施形態 53 :**

各シード配列および標的ヌクレオチド対は、定義された長さを有するシード配列、および所定のヌクレオチド配列上のシード配列の直後に標的ヌクレオチドを含む、実施形態 45 ~ 52 のいずれかに記載の実施形態。

30

**【0229】****実施形態 54 :**

定義された長さが 10 塩基である、実施形態 53 に記載の実施形態。

**実施形態 55 :**

訓練データセット中の各ヌクレオチド配列について、複数のシード配列および標的ヌクレオチド対を生成することは、シャープ T S S またはブロード T S S 標識に基づいて、訓練データセット中のヌクレオチド配列をシャープ T S S 群またはブロード T S S 群に分割することと、定義された長さのスライディングウィンドウを適用し、定義された工程サイズを各ヌクレオチド配列に有ることと、スライディングウィンドウの各工程でシード配列および標的ヌクレオチド対を保存することと、を含む、実施形態 45 ~ 54 のいずれかに記載の実施形態。

40

**【0230】****実施形態 56 :**

複数のシード配列および標的ヌクレオチド対の各シード配列および標的ヌクレオチド対をベクター化することは、各ヌクレオチドをそれぞれの番号としてコード化することを含む、実施形態 45 ~ 55 のいずれかに記載の実施形態。

**【0231】****実施形態 57 :**

50

生成モデルが、長短期メモリ ( L S T M ) リカレントニューラルネットワーク ( R N N ) を含む、実施形態 4 1 ~ 5 6 のいずれかに記載の実施形態。

【 0 2 3 2 】

実施形態 5 8 :

生成モデルに基づいて、ヌクレオチド配列を生成することをさらに含む、実施形態 4 1 ~ 5 7 のいずれかに記載の実施形態。

【 0 2 3 3 】

実施形態 5 9 :

生成モデルに基づいて、ヌクレオチド配列を生成することは、( a ) シード配列を受信することと、( b ) シード配列に基づいて、次のヌクレオチドを予測することと、( c ) シード配列に次のヌクレオチドを付加することと、d ) ヌクレオチド配列の所望の長さに達するまで b ~ c を繰り返すことと、を含む、実施形態 5 8 に記載の実施形態。

【 0 2 3 4 】

実施形態 6 0 :

所望の長さは、約 5 0 ヌクレオチド ~ 約 1 0 0 ヌクレオチドである、実施形態 5 9 に記載の実施形態。

【 0 2 3 5 】

実施形態 6 1 :

ヌクレオチド配列がコアプロモーター配列である、実施形態 4 1 ~ 6 1 のいずれかに記載の実施形態。

【 0 2 3 6 】

実施形態 6 2 :

コアプロモーター配列に基づいてプロモーターを操作することをさらに含む、実施形態 6 1 に記載の実施形態。

【 0 2 3 7 】

実施形態 6 3 :

プロモーターを核酸構築物に挿入することをさらに含む、実施形態 6 2 に記載の実施形態。

【 0 2 3 8 】

実施形態 6 4 :

プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含む、実施形態 6 3 に記載の実施形態。

【 0 2 3 9 】

実施形態 6 5 :

核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することをさらに含む、実施形態 6 4 に記載の実施形態。

【 0 2 4 0 】

実施形態 6 6 :

ヌクレオチド配列を受信することと、訓練された予測モデルにヌクレオチド配列を提供することと、および予測モデルに基づいて、ヌクレオチド配列がコアプロモーターであることを決定することと、を含む方法。

【 0 2 4 1 】

実施形態 6 7 :

ヌクレオチド配列を受信することは、複数のヌクレオチド配列を受信することを含み、複数のヌクレオチド配列は、生成モデルによって生成された、実施形態 6 6 に記載の実施形態。

【 0 2 4 2 】

実施形態 6 8 :

ヌクレオチド配列がコアプロモーターであるという決定に基づいて、一つまたは複数の

10

20

30

40

50

基準に従ってヌクレオチド配列をフィルタリングすることをさらに含む、実施形態 66 ~ 67 のいずれかに記載の実施形態。

【0243】

実施形態 69 :

一つまたは複数の基準が、GC 含量またはモチーフのうちの一つまたは複数を含む、実施形態 68 に記載の実施形態。

【0244】

実施形態 70 :

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点 (TSS) を含む、遺伝的データを受信することと、閾値を満たす関連発現スコアに基づいて、第一の複数のヌクレオチド配列から複数の TSS を決定することと、複数の TSS に基づいて、複数のサミットヌクレオチド塩基を決定することと、複数のサミットヌクレオチド塩基の各サミットヌクレオチド塩基に対して、関連する複数の周辺塩基を決定することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列として、各サミットヌクレオチド塩基および関連する複数の周辺塩基を保存することと、第二の複数のヌクレオチド配列の各ヌクレオチド配列に対して、関連する複数のシフト塩基を決定することと、コアプロモーターではないとして標識された第三の複数のヌクレオチド配列として、各関連する複数のシフト塩基を保存することと、コアプロモーターとして標識された第二の複数のヌクレオチド配列、およびコアプロモーターではないとして標識された第三の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、訓練データセットに基づいて、予測モデルに対する複数の特徴を決定することと、訓練データセットの第一の部分に基づいて、複数の特徴に従って予測モデルを訓練することと、訓練データセットの第二の部分に基づいて、予測モデルを試験することと、および試験に基づいて、予測モデルを出力することと、をさらに含む、実施形態 66 ~ 69 のいずれかに記載の実施形態。

【0245】

実施形態 71 :

(a) ヌクレオチド配列と配列長を受信すること、(b) 訓練された生成モデルに、ヌクレオチド配列を提供すること、(c) 生成モデルに基づいて、ヌクレオチド配列に関連付けられた次のヌクレオチドを決定すること、(d) ヌクレオチド配列に次のヌクレオチドを付与すること、(e) ヌクレオチド配列の長さが配列長に等しくなるまで b ~ d を繰り返すこと、および (f) ヌクレオチド配列をコアプロモーター配列として出力すること、を含む方法。

【0246】

実施形態 72 :

コアプロモーター配列に基づいてプロモーターを操作することをさらに含む、実施形態 71 に記載の実施形態。

【0247】

実施形態 73 :

プロモーターを核酸構築物に挿入することをさらに含む、実施形態 71 ~ 72 のいずれかに記載の実施形態。

【0248】

実施形態 74 :

プロモーターを核酸構築物に挿入することは、導入遺伝子上流の核酸構築物にプロモーターを挿入して、導入遺伝子の発現を駆動することを含む、実施形態 73 に記載の実施形態。

【0249】

実施形態 75 :

核酸構築物を含む、アデノ随伴ウイルスまたはレンチウイルスを作製することをさらに含む、実施形態 73 ~ 74 のいずれかに記載の実施形態。

10

20

30

40

50



## 【 0 2 5 0 】

実施形態 7 6 :

配列長は、約 5 0 ヌクレオチド ~ 約 1 0 0 ヌクレオチドである、実施形態 7 1 ~ 7 5 のいずれかに記載の実施形態。

## 【 0 2 5 1 】

実施形態 7 7 :

遺伝的データが、第一の複数のヌクレオチド配列を含み、複数のヌクレオチド配列の各ヌクレオチド配列が、関連発現スコアを有する少なくとも一つの転写開始点 ( T S S ) を含む、遺伝的データを受信することと、第一の複数のヌクレオチド配列に基づいて、コアプロモーターとして標識された第二の複数のヌクレオチド配列を決定することと、閾値を満たす関連発現スコアに基づいて、第二の複数のヌクレオチド配列から第三の複数のヌクレオチド配列を決定することと、第三の複数のヌクレオチド配列に基づいて、コアプロモーターではないとして標識された第四の複数のヌクレオチド配列を決定することと、コアプロモーターとして標識された第三の複数のヌクレオチド配列およびコアプロモーターではないとして標識された第四の複数のヌクレオチド配列に基づいて、訓練データセットを生成することと、および訓練データセットに基づいて、生成モデルを訓練することと、をさらに含む、実施形態 7 1 ~ 7 6 のいずれかに記載の実施形態。

10

## 【 0 2 5 2 】

実施形態 7 8 :

一つまたは複数の基準に従ってヌクレオチド配列をフィルタリングすることをさらに含む、実施形態 7 1 ~ 7 7 のいずれかに記載の実施形態。

20

## 【 0 2 5 3 】

実施形態 7 9 :

一つまたは複数の基準が、 G C 含量またはモチーフのうちの一つまたは複数を含む、実施形態 7 8 に記載の実施形態。

## 【 0 2 5 4 】

実施形態 8 0 :

実施形態 1 ~ 7 9 に記載のいずれかを行うため構成された、装置。

実施形態 8 1 :

装置が実施形態 1 ~ 7 9 に記載のいずれかを行うよう構成された、プロセッサが実行可能な指示実施形態を有する、コンピュータ可読媒体。

30

## 【 0 2 5 5 】

当業者は、通常の実験だけを用いることで、本明細書に記載の方法および組成物の特定の実施形態の多数の同等物を認識し、または確認できる。かかる同等物は、以下の特許請求の範囲に包含されることが意図される。

40

50

【図面】  
【図 1】

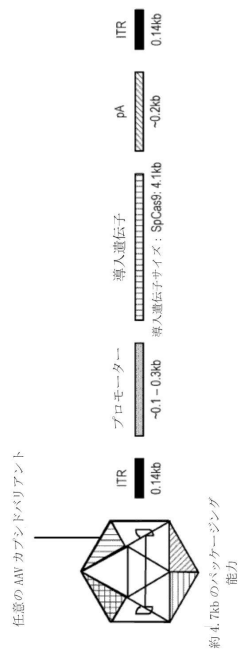
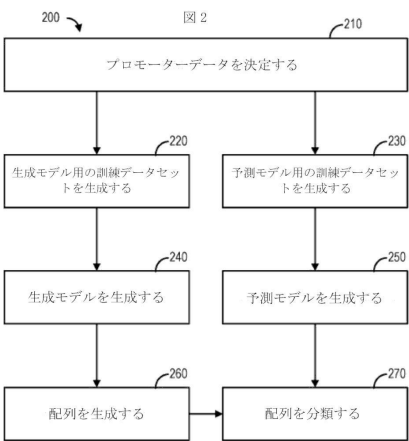


図 1

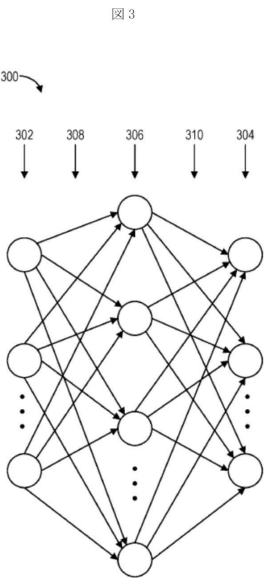
【図 2】



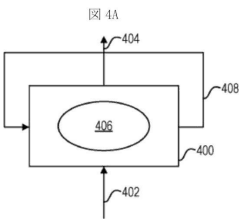
10

20

【図 3】



【図 4 A】

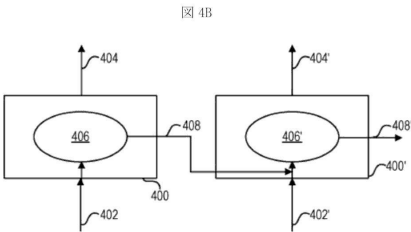


30

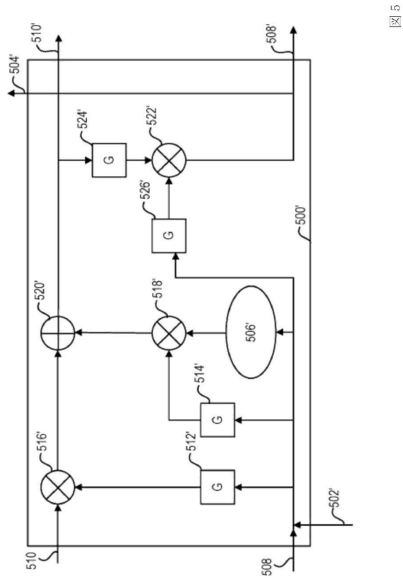
40

50

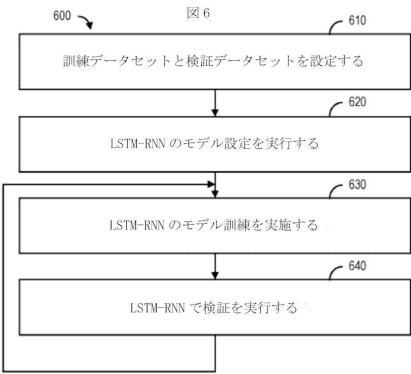
【図 4 B】



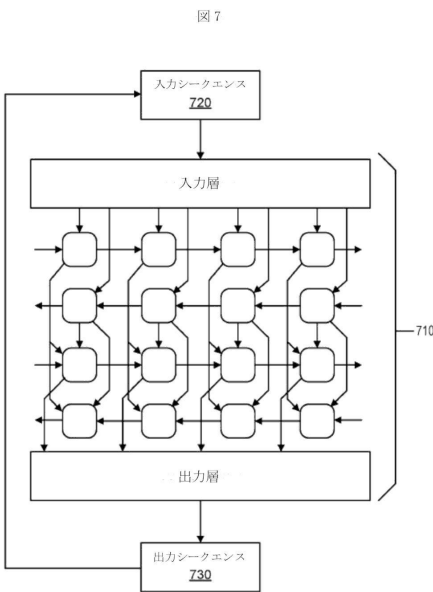
【図 5】



【図 6】



【図 7】



10

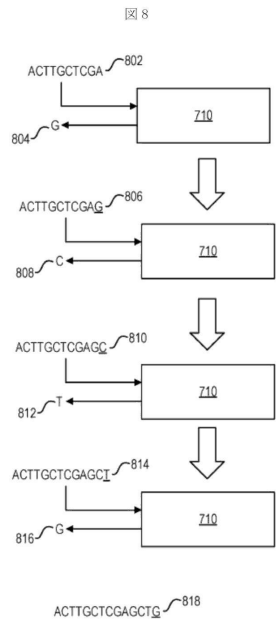
20

30

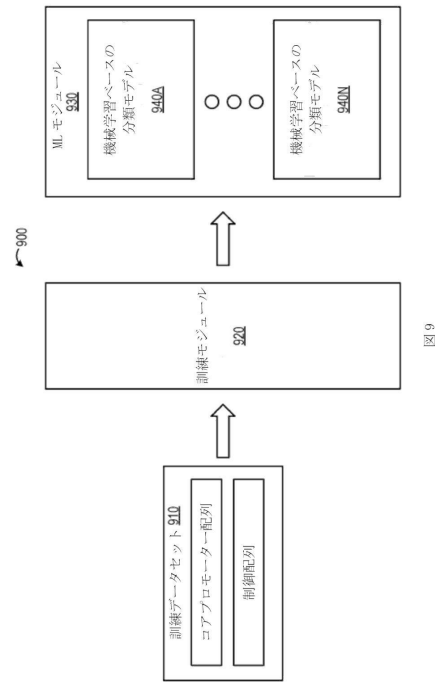
40

50

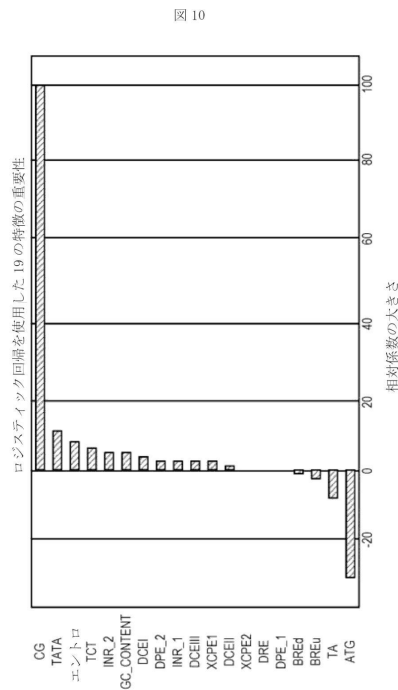
【図 8】



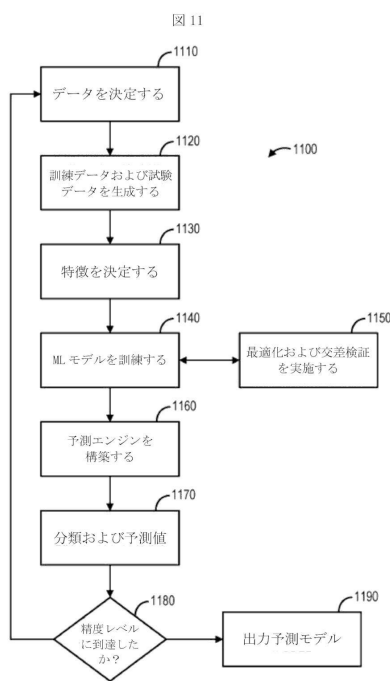
【図 9】



【図 10】



【図 11】



10

20

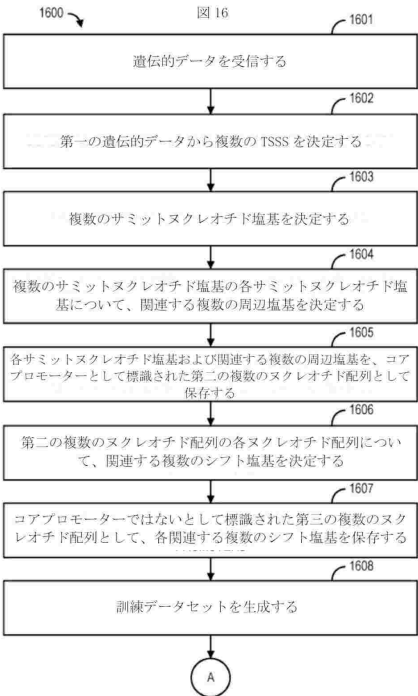
30

40

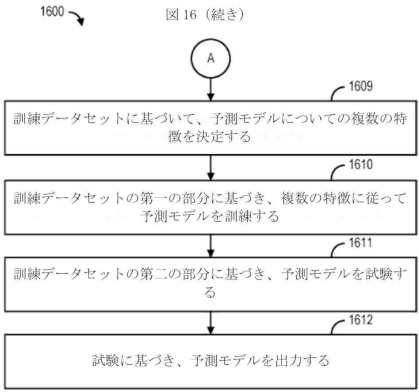
50



【図 16 - 1】



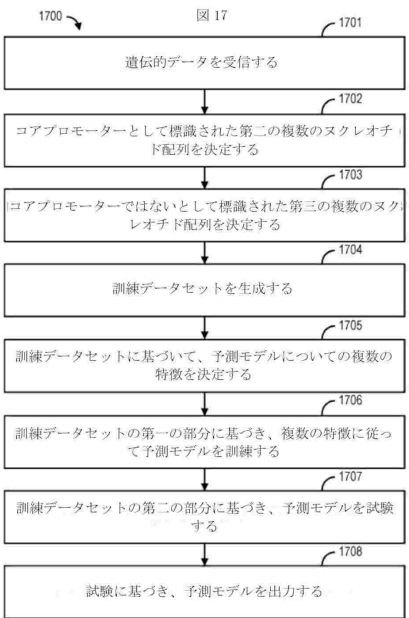
【図 16 - 2】



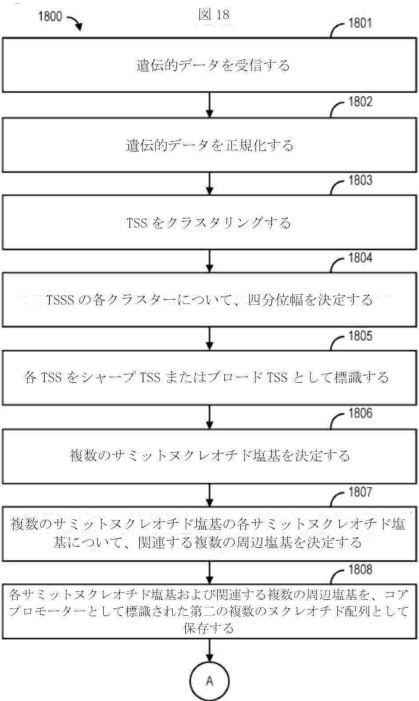
10

20

【図 17】



【図 18 - 1】

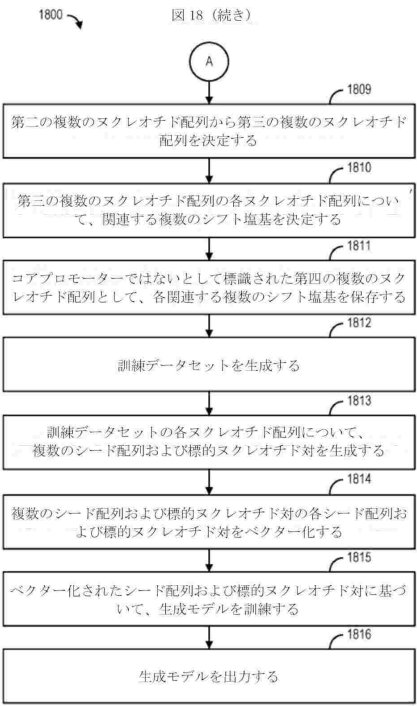


30

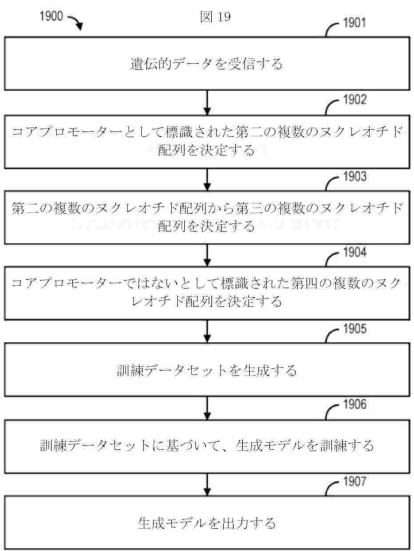
40

50

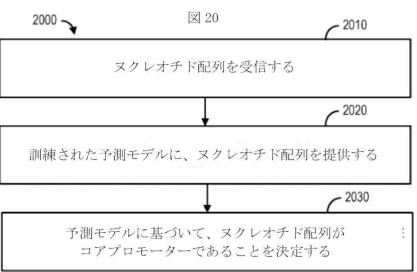
【図 18 - 2】



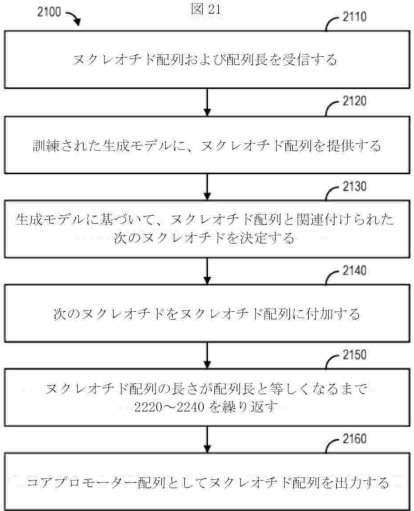
【図 19】



【図 20】



【図 21】



10

20

30

40

50

【配列表】

0007583153000001.app

10

20

30

40

50



フロントページの続き

弁理士 中村 美樹

(72)発明者 ミュルター、フェリクス  
アメリカ合衆国 1 0 5 9 1 ニューヨーク州 タリータウン オールド ソー ミル リバー ロード  
7 7 7

(72)発明者 シェーンヘル、クリストファー  
アメリカ合衆国 1 0 5 9 1 ニューヨーク州 タリータウン オールド ソー ミル リバー ロード  
7 7 7

審査官 山田 倍司

(56)参考文献 米国特許出願公開第 2 0 1 9 / 0 0 7 3 4 4 3 ( U S , A 1 )  
国際公開第 1 9 9 9 / 0 6 6 3 0 2 ( W O , A 2 )  
Ye Wang et al. , Synthetic Promoter Design in Escherichia coli based on Generative Adversarial Network , bioRxiv , 2019年04月25日 , <https://doi.org/10.1101/563775>  
Georgios K. Georgakilas et al. , Solving the transcription start site identification problem with ADAPT-CAGE: a Machine Learning algorithm for the analysis of CAGE data , Scientific Reports , 2020年01月21日 , Vol.10 , No.877 , <https://doi.org/10.1038/s41598-020-57811-3>  
Mhaned Oubounyt et al. , DeePromoter: Robust Promoter Predictor Using Deep Learning , Frontiers in Genetics , 2019年04月05日 , Vol.10 , No.286 , <https://doi.org/10.3389/fgenet.2019.00286>  
David R. Kelley , Cross-species regulatory sequence activity prediction , PLoS Computational Biology , 2020年07月20日 , Vol.16 , No.7 , <https://doi.org/10.1371/journal.pcbi.1008050>

(58)調査した分野 (Int.Cl. , D B 名)

G 1 6 B 5 / 0 0 - 9 9 / 0 0  
C 1 2 N 1 5 / 0 0 - 1 5 / 9 0  
C 1 2 Q 1 / 0 0 - 3 / 0 0